

---

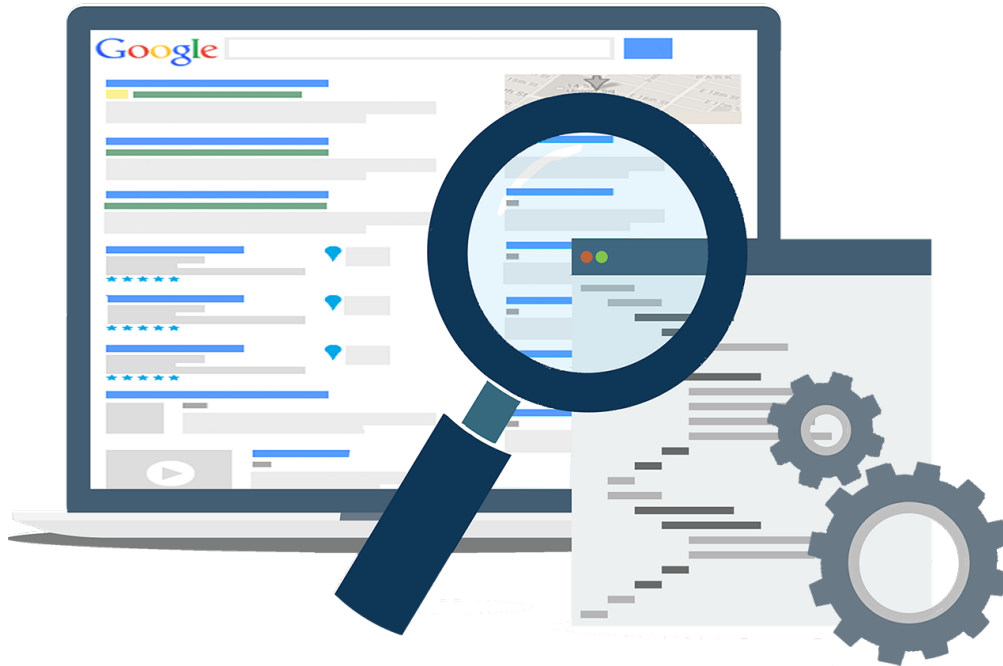
# Búsqueda y Minería de Información

## P2: Motores de Búsqueda e Indexación

---

Javier Alejandro Lougedo Lorente [javier.lougedo@estudiante.uam.es](mailto:javier.lougedo@estudiante.uam.es)

Alba Ramos Pedroviejo [alba.ramosp@estudiante.uam.es](mailto:alba.ramosp@estudiante.uam.es)



Búsqueda y Minería de Información Group 2461  
Profesor: Alejandro Bellogín Kouki  
Grado de Ingeniería Informática - Universidad Autónoma de Madrid  
5 de abril de 2021, Madrid, España

## Resumen

En la presente práctica implementaremos distintos motores de búsqueda y la indexación de documentos por distintas vías, como son las listas de *postings*, diccionarios y otros modelos almacenables en *RAM* y en memoria, así como la implementación de *PageRank*. Esto se hace con la ayuda de las librerías *Whoosh*, *BeautifulSoup4*, *ZipFile* y de alguna librerías *Python*, cuya instalación y uso están comentados en el fichero *readme\_lib.txt*.

## Índice

<b>1. Ejercicio 1: Un modelo vectorial eficiente</b>	<b>2</b>
1.1. Método Orientado a Términos - <b>TermBasedVSMSearcher</b>	2
1.2. Método Orientado a Documentos - <b>DocBasedVSMSearcher</b>	3
1.3. Heap de Ránking - <b>SearchRanking</b>	4
<b>2. Ejercicio 2: Índice en RAM</b>	<b>5</b>
2.1. Estructura de índice - <b>RAMIndex</b>	5
2.2. Construcción del índice - <b>RAMIndexBuilder</b>	6
<b>3. Ejercicio 3: Índice en disco</b>	<b>7</b>
3.1. Estructura de índice - <b>DiskIndex</b>	7
3.2. Construcción del índice - <b>DiskIndexBuilder</b>	9
<b>4. Ejercicio 4: Motor de búsqueda proximal</b>	<b>10</b>
4.1. Implementación - <b>ProximitySearcher</b>	10
<b>5. Ejercicio 5: Índice posicional</b>	<b>12</b>
5.1. Estructura de índice - <b>PositionalIndex</b>	12
5.2. Construcción del índice - <b>PositionalIndexBuilder</b>	12
<b>6. Ejercicio 6: PageRank</b>	<b>13</b>
6.1. Implementación - <b>PagerankDocScorer</b>	13
<b>7. Diagrama de Clases</b>	<b>15</b>
<b>8. Estudio y análisis del rendimiento</b>	<b>15</b>
<b>9. Conclusiones</b>	<b>17</b>

## 1. Ejercicio 1: Un modelo vectorial eficiente

En este ejercicio vamos a mejorar drásticamente la implementación que ya realizamos previamente, dejándola lenta al lado de las nuevas, aplicando los algoritmos que hemos aprendido en teoría, utilizando listas de postings en lugar de un índice forward.

### 1.1. Método Orientado a Términos - TermBasedVSMSearcher

En primer lugar, vamos a implementar el modelo vectorial coseno por el método orientado a términos.

Para ello, implementamos la clase de la siguiente forma:

Código Python 1: TermBasedVSMSearcher

```
1  class TermBasedVSMSearcher(Searcher):
2      def __init__(self, index, parser=BasicParser()):
3          super().__init__(index, parser)
4
5      def search(self, query, cutoff):
6          query_terms = self.parser.parse(query)
7          ranking = SearchRanking(cutoff)
8          posting_values = {}
9          scores = []
10
11         for term in query_terms:
12             for doc_id, _ in self.index.postings(term):
13                 if doc_id not in posting_values.keys():
14                     posting_values[doc_id] = []
15                     posting_values[doc_id].append(term)
16
17         for doc_id in posting_values.keys():
18             score = self.score(doc_id, posting_values[doc_id])
19             if score:
20                 ranking.push(self.index.doc_path(doc_id), score)
21
22         return ranking
23
24     def score(self, doc_id, query_terms):
25         product = 0
26         for term in query_terms:
27             product += tf(self.index.term_freq(term, doc_id)) * \
28                 idf(self.index.doc_freq(term), self.index.ndocs())
29
30         module = self.index.doc_module(doc_id)
31         return (product / module) if module else 0
```

Implementa, ya que en este caso lo vimos más cómodo de esta manera, dos funciones, como de costumbre, una para calcular el *score* y otra para calcular la puntuación y el *ranking* en sí de la colección, tal y como se nos enseña en las clases de teoría. No tuvimos problemas en la implementación y nos resultó sencilla y rápida, aunque intuimos por los tiempos obtenidos que hay alguna ligera mejora posible todavía.

## 1.2. Método Orientado a Documentos - DocBasedVSMSearcher

En este primer ejercicio opcional se pide llevar a cabo la implementación del método orientado a documentos, con el heap de postings que se implementa en el siguiente apartado.

### Código Python 2: DocBasedVSMSearcher

```

1  class DocBasedVSMSearcher(Searcher):
2      def __init__(self, index, parser=BasicParser()):
3          super().__init__(index, parser)
4
5      def search(self, query, cutoff):
6          ndocs = self.index.ndocs()
7          qterms = self.parser.parse(query)
8          indices = [0] * len(qterms) # Cosecuentials
9          postings = []
10         acc = [0] * ndocs # Accumulators
11         ranking = SearchRanking(cutoff)
12
13         for q in qterms:
14             postings += [[p for p in self.index.postings(q)]]
15
16         heap = []
17         heapify(heap)
18         for q in range(len(qterms)):
19             heappush(heap, (postings[q][indices[q]], q))
20             indices[q] += 1 # iterate cosecuentially
21
22         currDoc = nsmallest(1, heap)[0][0][0]
23         try:
24             while(1):
25                 element = heappop(heap)
26                 docid = element[0][0]
27                 freq = element[0][1]
28                 q = element[1]
29                 score = tf(freq) * idf(self.index.doc_freq(qterms[q]), ndocs)
30
31                 if docid != currDoc:
32                     acc[currDoc] /= self.index.doc_module(currDoc)
33                     # insert only if not less than top if full
34                     if len(ranking.ranking) < ranking.cutoff:
35                         ranking.push(self.index.doc_path(
36                             currDoc), acc[currDoc])
37                     elif score > nsmallest(1, ranking.ranking)[0][0]:
38                         ranking.push(self.index.doc_path(
39                             currDoc), acc[currDoc])
40
41                 currDoc = docid
42                 acc[currDoc] += score
43
44                 if indices[q] < len(postings[q]):
45                     heappush(heap, (postings[q][indices[q]], q))
46                     indices[q] += 1
47
48         except Exception: # El heap se ha vaciado
49             acc[currDoc] /= self.index.doc_module(currDoc)
50             if len(ranking.ranking) < ranking.cutoff:
51                 ranking.push(self.index.doc_path(currDoc), acc[currDoc])
52             elif score > nsmallest(1, ranking.ranking)[0][0]:
53                 ranking.push(self.index.doc_path(currDoc), acc[currDoc])
54
55         return ranking

```

En este caso, decidimos no implementar una función de *score* y realizarlo de esta manera porque nos resultaba más sencillo visualizarlo y no vimos la necesidad. En cuando a decisiones de diseño e implementación, no tuvimos muchos problemas tras asistir a las lecciones de teoría y haber realizado la práctica anterior.

Hemos seguido así la teoría lo máximo posible, creando una lista de acumuladores de *score* coseno para cada documento, y una lista de índices para cada término de la consulta, lo que nos permite realizar una iteración cosecuencial sobre sus listas de *postings*:

1. Obtenemos las listas y creamos un *heap* de tantos nodos como términos en la consulta, insertando el primer elemento de cada *qterm*.
2. Tras eso, mientras queden elementos, extraemos la cima (el *heap* se recolocará solo) y calculamos *tf-idf* (ya que la lista de *postings* viene con frecuencias).
3. Si el elemento pertenece al mismo *doc\_id* que el que se extrajo la última vez, se acumula su *tf-idf*, si no, se calcula el *score* coseno del *doc\_id* anterior, dividiendo su acumulado entre su módulo, y se inserta al *ranking*, acumulando además el *tf-idf* del elemento extraído.
4. Si quedan elementos en la lista de *postings* de ese *qterm*, se insertará el siguiente.
5. Sabremos que el *heap* se vacía porque *heappop* lanzará una excepción, en cuyo caso calculamos el *score* coseno del último elemento y se inserta al *ranking*.

Todas las inserciones se realizan comprobando si el *heap* de *ranking* se ha llenado de *cutoff*, en cuyo caso solo se inserta el elemento si tiene un *score* mayor al de la cima.

### 1.3. Heap de Ránking - SearchRanking

Se pide reimplementar la clase *SearchRanking* para utilizar un *heap* de *ranking*, de la que nos podremos aprovechar en ambos métodos implementados en este primer apartado. Este será un *minHeap* de tuplas (*score*, *id*), implementado con ayuda de la librería **heapq**, que parte de una lista y la transforma en un *minHeap* mediante *heapify*. Nos permite insertar nodos con *heappush*, donde insertaremos tuplas como las anteriormente descritas. Para el *score*, emplearemos un método *iter* que devolverá un iterador, formado por una lista de tuplas que deben estar ordenadas en función de su *score*, en orden decreciente (cabe recordar en este punto la teoría, donde creamos un *minHeap* y luego el *ranking* era el reverso de este). Para ordenarlas, usamos la función *nlargest*, que devolverá los N elementos más grandes del *heap* en forma de una lista de tuplas. Este N es, como cabe suponer, el *cutoff*, que resulta útil para emplear este *heap* en otros modelos, y no solo en el orientado a documentos. Finalmente, el *main* espera las tuplas como (*doc\_id*, *score*), así que se invierten en el propio iterador.

La implementación realizada es la que se encuentra en la página a continuación.

## Código Python 3: SearchRanking

```
1 class SearchRanking:
2     def __init__(self, cutoff):
3         self.ranking = []
4         heapify(self.ranking)
5         self.cutoff = cutoff
6
7     def push(self, docid, score):
8         heappush(self.ranking, (score, docid))
9
10    def __iter__(self):
11        return iter([tuple(reversed(e)) for e in nlargest(
12            min(len(self.ranking), self.cutoff),
13            self.ranking)])
```

## 2. Ejercicio 2: Índice en RAM

En este ejercicio llevaremos a cabo la implementación de un índice propio que tenga las mismas funciones que la que ya definimos en la práctica anterior, basada en **Whoosh**. Para implementarlo en RAM, se nos recomienda el uso de la librería **pickle**, a fin de guardarlo en disco y leerlo serializado. Se divide este ejercicio en dos apartados: la implementación del índice en sí y de su *builder*.

### 2.1. Estructura de índice - RAMIndex

Para la implementación del índice en RAM hemos utilizado un diccionario igual que los que hemos visto en la teoría: las claves son los términos, y el valor de estas será la lista de *postings* (en formato *(doc\_id, frecuencia)* de dicho término. De ese modo, para obtener la lista de *postings* de un determinado término, basta con indexar a través de dicho término en el diccionario y ya tendremos, de manera rápida y sencilla, su lista. Obtener el vocabulario es además tan fácil como obtener las claves del diccionario, con lo cual muchos procesos se simplifican bastante, y la implementación resulta sencilla, como veremos a continuación.

De esta forma, para implementar el índice en RAM, creamos un Índice base con dos variables nuevas, que serán un diccionario con los términos y sus *postings* de la forma y con la estructura ya mencionada. Redefinimos las dos funciones de *postings* y *all\_terms* del índice, y añadimos los apartados adicionales de los ficheros a las funciones de *open* y *save*, que llamarán adicionalmente a las del objeto base. Así, obtenemos un índice funcional de forma rápida y sencilla, pues las implementaciones de *postings* y *all\_terms* son rápidas y directas, a través de este diccionario, y las de *open* y *save* bastante lógicas e inmediatas, tomando como referencia las base y la ayuda de **pickle**. No es necesario implementar nada más en este punto.

La implementación se encuentra así en la siguiente página.

## Código Python 4: RAMIndex

```

1  class RAMIndex(Index):
2      def __init__(self, dir_path=None):
3          self.terms_dict = {} # Indice en si: {term: postings}
4          self.postings_path = os.path.join(dir_path, Config.POSTINGS_FILE)
5          super().__init__(dir_path)
6
7      def postings(self, term):
8          return self.terms_dict[term]
9
10     def all_terms(self):
11         return self.terms_dict.keys()
12
13     def save(self, dir_path):
14         with open(os.path.join(dir_path, Config.PATHS_FILE), 'wb') as f:
15             pickle.dump(self.docmap, f)
16         with open(os.path.join(dir_path, Config.DICTIONARY_FILE), 'wb') as f:
17             self.terms_dict = collections.OrderedDict(
18                 sorted(self.terms_dict.items())) # Orden alfabético
19             pickle.dump(self.terms_dict, f)
20         super().save(dir_path)
21
22     def open(self, dir_path):
23         super().open(dir_path)
24
25         try:
26             with open(os.path.join(dir_path, Config.PATHS_FILE), 'rb') as f:
27                 self.docmap = pickle.load(f)
28             with open(os.path.join(dir_path, Config.DICTIONARY_FILE), 'rb') as f:
29                 self.terms_dict = pickle.load(f)
30         except OSError:
31             # the file may not exist the first time
32             pass

```

## 2.2. Construcción del índice - RAMIndexBuilder

Para construirlo, modificamos la función de *commit*, donde se guardan las estructuras del índice tal y como se indica más adelante, así como la de *index\_document*. Aquí, cada vez que llega un nuevo documento, almacenamos su *path* y contamos los documentos que llevamos, a fin de ir asignando un *doc\_id* actualizado a cada documento. Se crea un diccionario auxiliar de término-frecuencia donde vamos almacenando la información que tendrá el índice, y posteriormente *parseamos* el documento, vía *BasicParser* (o cualquier otro indicado), y recorremos cada término realizando el conteo en el diccionario y actualizando su información. Tras procesar el documento, recorremos los términos de este diccionario y se va añadiendo al índice la lista de *postings* de dicho documento, pues ya conocemos la frecuencia.

En cuanto a la función del *commit*, su implementación resulta sencilla, pues viendo lo que hacía el código fuente de *Whoosh* dedujimos que bastaba con hacer un *dump* del índice en sí para posteriormente cargarlo, así como llamar al *save* del índice.

El resto de funciones se heredan del *Builder* base, pues son iguales y no necesitan ningún tipo de cambio o actualización para funcionar correctamente.

## Código Python 5: RAMIndexBuilder

```
1 class RAMIndexBuilder(Builder):
2     def __init__(self, dir_path, parser=BasicParser()):
3         super().__init__(dir_path, parser)
4         self.path = dir_path
5         self.index = RAMIndex(dir_path)
6
7     def commit(self):
8         self.index.save(self.path)
9         with open(self.path + Config.INDEX_FILE, "wb") as f:
10             pickle.dump(self.index, f)
11
12     def index_document(self, path, text):
13         new_doc_id = self.index.ndocs()
14         self.index.add_doc(path)
15         text = BasicParser.parse(text)
16         doc_term_freq = {} # Diccionario term:freq auxiliar en este documento
17         for pos, term in enumerate(text):
18             if self.index.terms_dict.get(term) is None:
19                 doc_term_freq[term] = 0
20                 self.index.terms_dict[term] = []
21             elif doc_term_freq.get(term) is None:
22                 doc_term_freq[term] = 0
23             doc_term_freq[term] += 1
24         for term in doc_term_freq.keys():
25             self.index.terms_dict[term].append(
26                 [new_doc_id, doc_term_freq[term]])
```

### 3. Ejercicio 3: Índice en disco

En este segundo ejercicio opcional, se nos pide llevar a cabo la implementación del índice en **disco**. Para ello, se nos ocurrió incluso que heredase de **RAMIndex**, pero decidimos evitar esto para tener una implementación algo más clara y flexible, llegado el caso, heredando directamente de **Index** en su lugar.

#### 3.1. Estructura de índice - DiskIndex

En primer lugar vamos a hablar de la estructura del índice. En esta ocasión, lo creamos con un diccionario de términos y *postings*, otro de términos y sus *offsets* para cuando sea necesario ir a memoria, otro con sus *paths*, como el de **RAMIndex**, y un conteo de los documentos, que facilita las funciones de *ndocs*. En cuanto a la función de *postings*, en esta ocasión es mucho más compleja. Lo que hacemos es abrir el fichero en el que tenemos todo almacenado de la manera más eficiente a la par que sencilla y clara que hemos encontrado, se va a la posición que indique el offset para dicho término y lee dicha fila del documento, cargando los *postings*. En cuanto a la función de *doc\_freq*, a fin de obtener resultados más eficientes, la reimplementamos también, pues en esta ocasión en lugar de cargar los *postings* y medir su longitud, lo que hace es leer un número que hemos introducido al guardar los *postings* en la función *save*, que es la longitud de todos ellos. De esta forma, esta función pasa a ser mucho más eficiente. En cuanto al *save* y al *open*, se aprovechan en gran medida de *pickle*,



guardando el fichero en un formato cómodo para la función de *postings*, que hemos intentado que adicionalmente, sea legible.

### Código Python 6: DiskIndex

```

1  class DiskIndex(Index):
2      def __init__(self, dir_path=None):
3          self.postings_dict = {} # term:[postings]
4          self.terms_dict = {} # {term: offsetPostings}
5          self.postings_path = os.path.join(dir_path, Config.POSTINGS_FILE)
6          self.count = 0
7          super().__init__(dir_path)
8
9      def postings(self, term):
10         res_postings = []
11         with open(self.postings_path, 'r') as f:
12             f.seek(self.terms_dict[term])
13             doc_amount, raw_postings = f.readline().strip().split("=")
14             raw_postings = raw_postings.split(" ")
15             for i in range(0, int(doc_amount) * 2, 2):
16                 res_postings.append(
17                     [int(raw_postings[i]), int(raw_postings[i+1])])
18         return res_postings
19
20     def all_terms(self):
21         return self.terms_dict.keys()
22
23     def doc_freq(self, term):
24         with open(self.postings_path, 'r') as f:
25             f.seek(self.terms_dict[term])
26             result = f.readline().split("=")[0]
27         return int(result)
28
29     def save(self, dir_path):
30         with open(os.path.join(dir_path, Config.PATHS_FILE), 'wb') as f:
31             pickle.dump(self.docmap, f)
32         with open(os.path.join(dir_path, Config.POSTINGS_FILE), 'w') as f:
33             for term, postings in self.postings_dict.items():
34                 self.terms_dict[term] = f.tell()
35                 n_postings = len(postings)
36                 f.write(str(n_postings) + "=")
37                 for p in postings:
38                     f.write(str(p[0]) + " ")
39                     f.write(str(p[1]) + " ")
40                 f.write("\n")
41         with open(os.path.join(dir_path, Config.DICTIONARY_FILE), 'wb') as f:
42             self.terms_dict = collections.OrderedDict(
43                 sorted(self.terms_dict.items())) # Orden alfabético
44             pickle.dump(self.terms_dict, f)
45         super().save(dir_path) # Guarda módulos
46
47     def open(self, dir_path):
48         super().open(dir_path) # Carga módulos
49         self.postings_dict = {}
50         try:
51             with open(os.path.join(dir_path, Config.PATHS_FILE), 'rb') as f:
52                 self.docmap = pickle.load(f)
53             with open(os.path.join(dir_path, Config.DICTIONARY_FILE), 'rb') as f:
54                 self.terms_dict = pickle.load(f)
55         except OSError:
56             # the file may not exist the first time
57             pass

```

Con esta estructura, obtener los *postings* de un término implica leer el fichero de *postings* y desplazarse al offset indicado. Una vez ahí, leer la línea y recorrer los *docid-freq* que hemos almacenado. Por otro lado, obtener la cantidad de documentos donde aparece un término es tan sencillo como mirar el fichero de *postings*, desplazarse al offset de dicho término y leer el primer número. Obtener el vocabulario se haría igual que **RAMIndex**. El resto de funcionalidad es igual que en **Index**.

Adicionalmente, cuando nos encontramos con un problema en cuanto a la eficiencia que no fuimos capaces de ubicar inicialmente, implementamos una variante, *CachedDiskIndex*, donde guardamos una pequeña caché en RAM para evitar ir muchísimas veces a disco a por el mismo término cuando estemos realizando búsquedas. Esto optimiza bastante los tiempos, como pudimos observar en su momento, pese a tener una implementación ineficiente del guardado y lectura del fichero de *postings*. Posteriormente, mejoramos estas funciones y manera de guardarlo obteniendo una mayor eficiencia y haciendo que no fuese necesaria esta pequeña caché, que sencillamente lo que hacía era almacenar en un diccionario aquellos *postings* que había leído recientemente. Esto se podría implementar, con algún tipo de *heap* o similar, que en cuanto pasase de un determinado número de elementos en la caché, o de una determinada cantidad de memoria, eliminase el elemento por el que se consultó hace más tiempo, para cargar uno más consultado de memoria.

### 3.2. Construcción del índice - DiskIndexBuilder

En este caso, es exactamente igual que la de **RAMIndex**, pues el tratamiento es el mismo, excepto por la línea 5, donde se crea un **DiskIndex** en lugar de un **RAMIndex**.

Código Python 7: DiskIndexBuilder

```
1  class DiskIndexBuilder(Builder):
2      def __init__(self, dir_path, parser=BasicParser()):
3          super().__init__(dir_path, parser)
4          self.path = dir_path
5          self.index = DiskIndex(dir_path)
6
7      def commit(self):
8          self.index.save(self.path) # Guardamos las estructuras del indice
9          with open(self.path + Config.INDEX_FILE, "wb") as f:
10             pickle.dump(self.index, f) # Guardamos el indice en si
11
12      def index_document(self, path, text):
13          new_doc_id = self.index.ndocs()
14          self.index.add_doc(path)
15          text = BasicParser.parse(text)
16          doc_term_freq = {} # Diccionario term:freq en este documento
17          for pos, term in enumerate(text):
18             if self.index.postings_dict.get(term) is None:
19                 doc_term_freq[term] = 0
20                 self.index.postings_dict[term] = []
21             elif doc_term_freq.get(term) is None:
22                 doc_term_freq[term] = 0
23             doc_term_freq[term] += 1
24          for term in doc_term_freq.keys():
25             self.index.postings_dict[term].append(
26                 [new_doc_id, doc_term_freq[term]])
```

## 4. Ejercicio 4: Motor de búsqueda proximal

Este ejercicio nos costó especialmente, a la hora de implementar el algoritmo del *score* de manera adecuada y eficiente. Finalmente, gracias a un pequeño estudio de dos de los *papers* que encontramos en Internet, logramos implementarlo con éxito. En él se pide, básicamente, implementar un motor de búsqueda proximal.

### 4.1. Implementación - ProximitySearcher

Como ya hemos comentado, para implementarlo nos ayudamos de un par de *papers* de Internet, así como las diapositivas en teoría. Es un algoritmo que en papel y en los ejercicios teóricos resulta sencillo y rápido de aplicar, pero que a la hora de programar y traducir a código nos dio problemas. Sin embargo, finalmente conseguimos una implementación exitosa, como se podrá observar más abajo, en la cual no hay comentarios explicativos de código (como si los hay en el código fuente en sí mismo, a fin de reducir el tamaño de la figura).

De esta forma, lo que hacemos es, en primer lugar, es *parsear* la *query*, crear un *ranking* de *cutoff* elementos y observar cuales son los documentos válidos con todas las consultas, a fin de calcular su *score*. Esto lo hacemos cogiendo los documentos que contienen el primer elemento de la consulta y aplicando la intersección con los que contienen el segundo elemento, la intersección del resultado con los que contienen el tercero, etc. Así, acabamos con un *set* de *doc\_ids* con los documentos que contienen todos y cada uno de los términos, independientemente del orden de aparición.

Tras esto, lo que haremos será obtener y limpiar toda la información vía *positional postings*. Lo que hacemos es, básicamente, llamar a *positional postings*, eliminando todos y cada uno de los documentos que no figuren en el *set* de documentos válidos.

Cuando ya lo tenemos listo, para cada uno de estos documentos **higienizados**, llamamos a la función de *score*, que obtendrá B y A de cada grupo de elementos y calculará el *score*. Esta es la parte que más dificultades tuvimos en implementar, ya que no terminábamos de entender muy bien este algoritmo. Su pseudocódigo es el siguiente:

```

score ← 0
for j=1 to |q| do p[j] ← 0
b ← maxj posList[j][p[j]]
while b < ∞ do
  i ← 0
  for j=1 to |q| do
    // Avanzar todas las listas hasta posición b
    while posList[j][p[j]+1] ≤ b do p[j]++
    // Nos queremos quedar en i con la posición más baja de las |q| listas
    if posList[j][p[j]] < posList[i][p[i]] then i ← j
  a ← posList[i][p[i]]
  score ← score + 1/(b - a - |q| + 2)
  b ← posList[i][++p[i]]
return score

```

Figura 1: Pseudocódigo de la función evaluadora

Cabe mencionar que nuestra implementación, gracias a estas higienizaciones y limpiezas, resulta extremadamente eficiente, como se puede observar comparando tiempos de ejecución.

## Código Python 8: ProximitySearcher

```

1  class ProximitySearcher(Searcher):
2      def __init__(self, index, parser=BasicParser()):
3          super().__init__(index, parser)
4
5      def search(self, query, cutoff):
6          query_terms = self.parser.parse(query)
7          ranking = SearchRanking(cutoff)
8
9          suitable_doc_ids = set()
10         for doc_id, _ in self.index.postings(query_terms[0]):
11             suitable_doc_ids.add(doc_id)
12
13         if len(query_terms) > 1:
14             for term in query_terms[1:]:
15                 current_doc_ids = set()
16                 for doc_id, _ in self.index.postings(term):
17                     current_doc_ids.add(doc_id)
18                 suitable_doc_ids &= current_doc_ids
19
20         cleaned_positions = {}
21         for term in query_terms:
22             term_clean_positions = {}
23             term_positions = self.index.positional_postings(term)
24             for doc_positions in term_positions:
25                 if doc_positions[0] in suitable_doc_ids:
26                     term_clean_positions[doc_positions[0]] = doc_positions[1]
27             cleaned_positions[term] = term_clean_positions
28
29         for doc_id in suitable_doc_ids:
30             ranking.push(self.index.doc_path(doc_id),
31                         self.score(cleaned_positions, doc_id, query_terms))
32
33         return ranking
34
35     def score(self, cleaned_positions, doc_id, query_terms):
36         posList = []
37         for term in query_terms:
38             posList.append(cleaned_positions[term][doc_id] + [math.inf])
39         q_len = len(query_terms)
40         if q_len == 1:
41             return float(len(posList[0]) - 1)
42         a = - 1
43         score = 0
44         p = [0] * q_len
45         b = max([t_list[0] for t_list in posList])
46
47         while b != math.inf:
48             i = 0
49             for j in range(q_len):
50                 while posList[j][p[j]+1] <= b:
51                     p[j] += 1
52                 if posList[j][p[j]] < posList[i][p[i]]:
53                     i = j
54             a = posList[i][p[i]]
55             score += 1 / (b - a - q_len + 2)
56             b = posList[i][p[i]+1]
57
58         return score

```

## 5. Ejercicio 5: Índice posicional

En este cuarto y último ejercicio voluntario, se nos pide implementar una variable adicional que extienda los índices con la inclusión de soporte para posiciones en las listas de *postings*. En nuestro caso, este heredaría de la implementación en *RAM*, ya que para los datos que manejamos es más rápido y tiene unas estructuras similares.

### 5.1. Estructura de índice - PositionalIndex

Conocer la lista de *postings* como hasta ahora implica, para un término, iterar sobre las claves de su diccionario (*docids*) y ver la longitud de cada lista de posiciones (*freq*). La lista de *postings* posicionales simplemente implica, para un término, iterar sobre su diccionario y obtener la lista de posiciones de cada *docid*. Lo demás sería igual que en **RAMIndex**, que es la superclase.

También probamos, y podemos, hacer que herede de *DiskIndex*. Sin embargo, nos parecía ligeramente más lento y poco ágil en este sentido. El resultado final de relaciones de herencias entre los distintos índices se puede observar más adelante en el diagrama de clases.

Código Python 9: PositionalIndex

```
1  class PositionalIndex(RAMIndex):
2      def __init__(self, dir_path=None):
3          self.terms_dict = {}
4          super().__init__(dir_path)
5
6      def positional_postings(self, term):
7          dict = self.terms_dict[term]
8          res = []
9          for doc, pos in dict.items():
10             res += (doc, pos) # [ [ doc1, [posiciones] ], [ doc2, [posiciones] ], ...]
11
12         return res
13
14     def postings(self, term):
15         dict = self.terms_dict[term]
16         res = []
17         for doc, pos in dict.items():
18             res.append((doc, len(pos)))
19         return res
```

### 5.2. Construcción del índice - PositionalIndexBuilder

Para construir el índice, no hemos de almacenar frecuencias esta vez, sino posiciones, con lo que vamos a tener un índice que sea un diccionario término-diccionario, donde este segundo diccionario será uno con la forma *docid*-posiciones. Así, el diccionario auxiliar que teníamos antes con la forma término-frecuencia, se convierte en uno de la forma término-posiciones, con lo que podremos conocer las posiciones de manera sencilla, vía la función *enumerate* al iterar.

## Código Python 10: PositionalIndexBuilder

```
1 class PositionalIndexBuilder(RAMIndexBuilder):
2     def __init__(self, dir_path, parser=BasicParser()):
3         super().__init__(dir_path, parser)
4         self.index = PositionalIndex(dir_path)
5
6     # Vamos creando el índice en RAM
7     def index_document(self, path, text):
8         new_doc_id = self.index.ndocs()
9         self.index.add_doc(path)
10        text = BasicParser.parse(text)
11        doc_term_pos = {} # Diccionario term: [pos1, pos2...]
12        for pos, term in enumerate(text):
13            if self.index.terms_dict.get(term) is None:
14                doc_term_pos[term] = []
15                self.index.terms_dict[term] = {}
16            elif doc_term_pos.get(term) is None:
17                doc_term_pos[term] = []
18            doc_term_pos[term].append(pos)
19
20        for term in doc_term_pos.keys():
21            self.index.terms_dict[term][new_doc_id] = doc_term_pos[term]
```

## 6. Ejercicio 6: PageRank

En el último ejercicio de esta práctica, se nos pide implementar el algoritmo *PageRank* en la clase **PageRankDocScorer**, que devolverá un *ranking* de manera similar a como lo hace *Searcher*.

Cuando creamos este *scorer*, leemos el documento de grafos y vamos almacenando todos los nodos del grafo (sin repetición), así como un diccionario con los nodos a los que apunta cada nodo (origen-destinos) y otro con los nodos que apuntan a cada nodo (destino-orígenes). Los sumideros se detectan porque tendrán la lista de destinos vacía, y los podemos guardar en una lista aparte para acelerar el proceso de *ranking* posterior.

### 6.1. Implementación - PagerankDocScorer

Para implementar *PageRank* se ha seguido el algoritmo visto en clase. Inicialmente se les da a todos los nodos un *PageRank* de 1, y luego para cada nodo se calcula su *PageRank* con respecto a los valores de los nodos que le apuntan (más los sumideros) en la iteración anterior. El valor de la iteración anterior se almacena en un diccionario *PageRankPrev* una vez hechos los cálculos de esa iteración. Esto se realiza para tantas iteraciones como nos especifiquen. Finalmente, los valores se ordenan de mayor a menor y se devuelven como tuplas (nodo, *PageRank*) hasta alcanzar el *cutoff*.

## Código Python 11: PagerankDocScorer

```

1  class PagerankDocScorer():
2      def __init__(self, graphfile, r, n_iter):
3          self.n_iter = n_iter
4          with open(graphfile, "r") as f:
5              lines = f.read().splitlines()
6          self.nodes = set() # Elementos no ordenados y sin repetición
7          self.outs = {} # Clave: nodo, valor: lista de a quienes apunta
8          self.to_from_nodes = {} # Clave: nodo, valor: lista de quienes le apuntn
9          for line in lines:
10             ele = line.split('\t')
11             self.nodes.add(ele[0])
12             self.nodes.add(ele[1])
13             try:
14                 self.outs[ele[0]].append(ele[1])
15             except:
16                 self.outs[ele[0]] = [ele[1]]
17             try:
18                 self.to_from_nodes[ele[1]].append(ele[0])
19             except:
20                 self.to_from_nodes[ele[1]] = [ele[0]]
21             try:
22                 if self.outs[ele[1]]:
23                     pass
24             except:
25                 self.outs[ele[1]] = []
26             try:
27                 if self.to_from_nodes[ele[0]]:
28                     pass
29             except:
30                 self.to_from_nodes[ele[0]] = []
31         self.sumideros = []
32         for node in self.outs.keys():
33             if len(self.outs[node]) == 0:
34                 self.sumideros.append(node)
35         self.n = len(self.nodes)
36         self.r = r
37
38     def rank(self, cutoff):
39         p1 = self.r / self.n
40         p2 = 1 - self.r
41         res = []
42         pageRankPrev = {}
43         pageRankAct = {}
44         for n in self.nodes:
45             pageRankPrev[n] = 1
46         for epoc in range(self.n_iter):
47             for n in self.nodes:
48                 sum = 0
49                 for origin in self.to_from_nodes[n]:
50                     sum += pageRankPrev[origin] / len(self.outs[origin])
51                 for i in self.sumideros:
52                     sum += pageRankPrev[i] / self.n
53                 pageRankAct[n] = p1 + p2 * sum
54                 for n in pageRankAct.keys():
55                     pageRankPrev[n] = pageRankAct[n]
56         ordered = dict(sorted(pageRankAct.items(), key=lambda item: item[1], reverse=True))
57         cont = 0
58         for k, v in ordered.items():
59             res.append((k, v))
60             cont += 1
61             if cont == cutoff:
62                 break
63         return res

```

## 7. Diagrama de Clases

Hemos realizado, tal y como se nos pedía en la entrega, un diagrama de clases **sencillo** en el que observar y entender las relaciones jerárquicas y herencias de las distintas clases. Hemos tratado de reducir al máximo la información mostrada para quedarnos con los detalles más importantes, así como simplificar las relaciones a gran escala para no sobrecomplicar el diagrama. El diagrama es este:

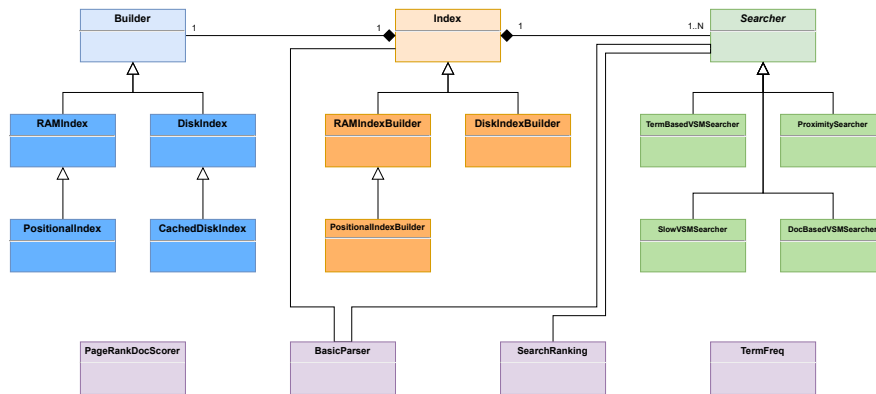


Figura 2: Diagrama de clases de los buscadores implementados.

En él, podemos observar los siguientes detalles relevantes: los *Builders* han de tener un Índice al que dan lugar (cada uno con el suyo propio, p.e. *RAMIndex* con *RAMIndexBuilder*), y lo mismo ocurre con los *Searchers* (que además es abstracta y deberá ser implementada). Los *indices* y *searchers* tienen *parsers*, en este caso solo *BasicParser*, y los *searchers* hacen uso del *ranking*, en este caso de nuevo solo uno. El resto de relaciones son claras.

## 8. Estudio y análisis del rendimiento

Para este estudio, se ha realizado la ejecución en un ordenador *Mac Mini* con *chipset* M1, con las siguientes características (más información [aquí](#)):

### Hardware Overview:

Model Name:	Mac mini
Model Identifier:	Macmini9,1
Chip:	Apple M1
Total Number of Cores:	8 (4 performance and 4 efficiency)
Memory:	16 GB
System Firmware Version:	6723.61.3
Serial Number (system):	C07F30VEQ6NY
Hardware UUID:	69582B71-A8D4-5869-B740-6ED666CC3DC6
Provisioning UDID:	00008103-000C684E21EA001E
Activation Lock Status:	Disabled

Figura 3: Información del entorno de ejecución



Los resultados obtenidos han sido los siguientes:

RAMIndex	Construcción del Índice			Carga del Índice	
Colección	Tiempo de Indexado	Consumo máx. RAM	Espacio en Disco	Tiempo de carga	Consumo máx. RAM
<b>Toy1</b>	0.001466989517211914 seconds	89-92 %	1312 space	0.00014901161193847656 seconds	89-92 %
<b>Toy2</b>	0.0008068084716796875 seconds	89-92 %	1199 space	0.0001399517059326172 seconds	89-92 %
<b>Urls</b>	1.9586701393127441 seconds	89-92 %	139734 space	0.028786897659301758 seconds	89-92 %
<b>1K</b>	32.373591899871826 seconds	89-92 %	7264736 space	1.4944581985473633 seconds	89-92 %
<b>10K</b>	216.32021570205688 seconds	89-92 %	51154094 space	13.291764974594116 seconds	89-92 %

Figura 4: Datos obtenidos de RAMIndex

DiskIndex	Construcción del Índice			Carga del Índice	
Colección	Tiempo de Indexado	Consumo máx. RAM	Espacio en Disco	Tiempo de carga	Consumo máx. RAM
<b>Toy1</b>	0.0039370059967041016 seconds	89-92 %	1203 space	0.000102996826171875 seconds	89-92 %
<b>Toy2</b>	0.004070281982421875 seconds	89-92 %	1102 space	0.0001049041748046875 seconds	89-92 %
<b>Urls</b>	2.15397310256958 seconds	89-92 %	125906 space	0.0013191699981689453 seconds	89-92 %
<b>1K</b>	37.925820112228394 seconds	89-92 %	5978775 space	0.03686094284057617 seconds	89-92 %
<b>10K</b>	241.544331073761 seconds	89-92 %	42537198 space	0.15534687042236328 seconds	89-92 %

Figura 5: Datos obtenidos de DiskIndex

PositionalIndex	Construcción del Índice			Carga del Índice	
Colección	Tiempo de Indexado	Consumo máx. RAM	Espacio en Disco	Tiempo de carga	Consumo máx. RAM
<b>Toy1</b>	0.0008111000061035156 seconds	89-92 %	1331 space	0.00010991096496582031 seconds	89-92 %
<b>Toy2</b>	0.0012600421905517578 seconds	89-92 %	1210 space	0.00011301040649414062 seconds	89-92 %
<b>Urls</b>	1.817927360534668 seconds	89-92 %	209804 space	0.004964113235473633 seconds	89-92 %
<b>1K</b>	33.027628898620605 seconds	89-92 %	11889007 space	1.0943729877471924 seconds	89-92 %
<b>10K</b>	220.53490924835205 seconds	89-92 %	79017238 space	16.21553111076355 seconds	89-92 %

Figura 6: Datos obtenidos de PositionalIndex

Comentando los mismos, el primer detalle a tener en cuenta, y es uno bastante relevante y que por suerte tenemos localizado, es el uso de la RAM. No hay muchas herramientas para tomar dicha medida de forma fiable, más allá de crear un proceso *sniffer* que estudie y mida cuanta RAM va gastando a cada momento. Por ello, optando por una solución más sencilla, decidimos intentar medir la RAM nosotros en todo momento, lo cual dio lugar a problemas, como los que observamos en la toma de datos, debido a que en los ordenador con *Mac OS*, concretamente este *Mac M1*, se hace un uso generalmente muy dinámico de la RAM en todos los programas, e independientemente de que modelo ejecutásemos, obteníamos siempre los mismos resultados. Es por ello que, a modo de por lo menos confirmar nuestras sospechas, estudiamos desde el Monitor de Sistema los valores que iba tomando el proceso *Python* en función de el modelo empleado. Nuestras sospechas fueron confirmadas, obteniendo valores muy altos para *PositionalIndex*, seguido de cerca por *RAMIndex* y muy de lejos por *DiskIndex*, que gastaba menores cantidades.

En cuanto a los tiempos, donde hemos podido llevar a cabo una mucho mejor toma de datos, observamos los siguientes detalles: en primer lugar, como es de esperar, *DiskIndex* es mucho más lento en construcción, aunque no excesivamente, tan solo un 15 % más lento que *RAMIndex*. *PositionalIndex*, para nuestra sorpresa, resulta ser muy próximo a *RAMIndex*, con diferencias poco notables en construcción. En cuanto a tiempos de carga, sin embargo, es donde llega el momento de arrasar de *DiskIndex*. Al tener que cargar una cantidad minúscula de datos, en comparación con los otros dos, podemos observar como es mucho más rápido (del orden de casi 100 veces más rápido) que los otros dos. Cabe mencionar, sin embargo, que estos tiempos pueden resultar engañosos, pues luego en búsqueda *DiskIndex* tardará algo más que estos dos, aunque si implementamos una caché intermedia como la que mencionamos anteriormente, obtendremos unos muy buenos resultados.

Por último, en lo que respecta al espacio en disco, nos sorprendimos muchísimo con los resultados observados, puesto que eran menos en *DiskIndex* y *RAMIndex*. Utilizamos para esta medida la funcionalidad que se facilitaba en el programa *main*, con lo que intuimos que funcionaría adecuadamente. De esta forma, como nos resultaba fácil de suponer, *PositionalIndex* ocupa más que *RAMIndex* (es lógico, al almacenar más información). Sin embargo, *DiskIndex* no era, para nuestra sorpresa, el que más información almacenaba, sino todo lo contrario. Creemos que esto se puede deber a que, al almacenar la información y volver a cargarla de nuevo, usando *pickle* en lugar de almacenar los *postings* de manera algo más eficiente, es posible que demos lugar a datos de un mayor volumen. Hemos comprobado estas mediciones investigando manualmente en la carpeta *index* generada por el programa y confirmando que, efectivamente, ocupa menos la de *DiskIndex*.

De esta forma, *DiskIndex*, que inicialmente no nos convencía porque creíamos que sería mucho más lento, ha demostrado ofrecer muy grandes y llamativas ventajas en comparación con los otros dos.

## 9. Conclusiones

Esta práctica nos ha permitido ahondar mucho más en los distintos motores de búsqueda que hemos ido implementando, a fin de crear distintos buscadores, modelos e índices, que nos han permitido entender plenamente su funcionamiento y analizarlo de cerca. Concretamente, nos ha abierto los ojos en dos detalles muy concretos, y es que *DiskIndex* de primeras no parecía especialmente bueno, de hecho, nos generaba nuestras dudas, pero con la Caché intermedia planteada nos puede permitir un muy buen funcionamiento, y mejor organizado todavía más aún; y por otro lado, que la búsqueda proximal es excesivamente compleja en algún aspecto, pero nos permite destacar, como hemos visto en los *score*, documentos que pueden ser muy aptos para nuestra búsqueda, permitiéndonos trabajar esa proximidad de los datos.