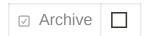
Notas 2do parcial - Ema



TEÓRICO

https://prod-files-secure.s3.us-west-2.amazonaws.com/206ee924-3842-491f-9926-b87d27a727f3/1c1e61f5-eeab-40ca-a69d-c500003fbeb0/Captulo_4_-Almacenamiento_de_tablas_en_archivos_e_ndices.pdf

https://prod-files-secure.s3.us-west-2.amazonaws.com/206ee924-3842-491f-9926-b87d27a727f3/46156097-ad3e-4ed7-afcc-fa23bcd801a7/Captulo_4_-Procesamiento_de_consultas.pdf

https://prod-files-secure.s3.us-west-2.amazonaws.com/206ee924-3842-491f-9926-b87d27a727f3/330490cb-5898-483c-8f90-f7a6ee963957/Captulo_4_-Factor_de_selectividad_y_operadores_fsicos.pdf

https://prod-files-secure.s3.us-west-2.amazonaws.com/206ee924-3842-491f-9926-b87d27a727f3/cde7c58b-4432-4288-b906-040fc0c6ec1d/Captulo_5_-Normalizacin.pdf

Práctico - RETORNO DE LA INFORMACIÓN

Pasado por una compañera:

Hola! Esto es lo que logré hacer del práctico de RI. Las únicas aclaraciones que te hago es que me faltan los ejercicio 8 y 10 (el 9 no hay que hacerlo) pero el 10 se respondió en clases y es una pavada. Y que el ejerció 7 me faltó la partecita esa de dar algoritmos y el ejercicio dos no le des mucha bola porque me ayude con chatGPT.

1	Ejercicio o.
	La recuperación de info. es el resultado de una consulta cugos datos del
	resultado son relevantes para la misma. En el contexto de la web el resultado de la consulta incluye enlaces a
	paginas web que sean relevantes para la consulta t un pequens cexto
	que describe el contenido de la página. Ejercicio 1.
	Punc Pind_documents (index, guery_terms, x) E
	result = CJ
	Por term in gery-terms If term in index docs-containing-terms-index [term]
	if not result:
	results = docs_containing_terms
	results = list (set (results). intersection (docs-containing-terms))
	Pinal-results = [doc Por doc in results if results.count(doc) >= K]
	return Anal-result.
	Ejercicio 2.
	Los componentes principales de un sist de EI (1) Obtener colecciones
	(2) Construir vocabulario (3) Resumir docs (construir vectores). (4) Construir indice invertido
	(5) Procesador de consultar, se puede subdividir: Ordenar resultados por relevancia
	Obtener relevancia ej calcular similia
	Ejercicio 3. Calcular resultados.
	El modelo del espacio vectorial de RI es una tecnica que nos ayuda a encontrar documentos relevantes en respuesta a una consulta. Un vector tiene la forma: (termino, peso) donde el termino es alguno de los terminos de la consulta y peso es alguna métrica de relevancia.
	como p.ej: la frewencia, el pess en terminos de TF-IDF, etc.
	Ejeraao 4.
sares	Un tesauro es un diccionario de sinonimos y antonimos. Es beneficioso para el el porque foicilità la busqueda de documento estan relacionados con los terminos de la busqueda, mo solo con los terminos en si sino con polobres que moto en los puedes, mo solo con los terminos en si sino con polobres que moto en los susqueda.
	terminos en si sim con palabras que están relaciona das con esos te

```
No, en este contexto no tendria sentido usar la frewencia como metrica
ua que al ser textos pequeños la probabilidad de que aporrescon multiples veces
n mismo termino es muy baja (a là sumo podra aparecer 1 o dos
veces).
  En cambio la Precuencia inversa de docs. si tendria sentido por
que en el coso de los tablos podramos ver en cuantos registros aparece una palabra si aparece en pocos registros la relevancia del termino sera alta de la contraria sera poco relevante.

Una manera definir la similitud de dos tuplos usando TF-IDF es con el cosens del angulo entre los vectores de ambos tuplos.
      Formula de la
      Similitud del coxeno | NAII. II BII
 Ejercicio 6.
    las computadoras pueden aprender de los datos el aprendiçaje automotico es un tipo de IA
   los algoritmos de aprendizaje automático predicen resultados
   las computadoros pueden aprender el lenguaje
El procesamiento del lenguaje natural es un campo de IA
    los algoritmos de procesamiento de lenguaje natural analizan texto
  la ordenadores pueden ver como los humanos
la visión del ordenador es un campo de la IA
los algoritmos de visión artificial reconocen objetos en imágenes
     la inteligencia artificial es un campo de la informatica los algoritmos de inteligencia artificial resuelven problemas
  Dor 5
El deep learning aprende patrones complejos
El deep learning potencia el reconocimiento de imploenes
El deep learning potencia el procesa miento del lenguaje natural
  Indice invertida
    term Downent, frecuencia computadoros 1:1, 2:1, 4:1
    IA 1:1, 2:1
lenguaje 2:3, 5:1
algoritmos 1:1, 2:1, 3:1, 4:1
```

```
Ejercicio 7.
   tara construir el resultado de una consulta
   - crear una lista con los terminos de la consulta
  - Para coda termino, se búsca en el indice invertido para encontrar todos los docs en los que el termino aparece - luego estos docs se ordenan sean algún criteno
   - Por ultimo la lista de los docs ordenados se devuelven como
   el resultado de la consulto.
   Di las consultas son basa das en proximidad entances en el indice
invertido no solo se quarda el doc en donde sale el termino sino
tambien en que posición el doc se encuentra dicho termino.
   near: dos terminos son "cercanos" si estan como mucho a cua-
tro terminos de distancia.
    adjacent: dos terminos están adjacentes si ambos terminos
aparecen uno seguido del otro sin ningun termino de por medio
after: un termino esta después que otro si ambos terminos aparecen uno seguido del atro con al menos un termino de por
medio.
    sintaxis: near (term1, term2), adjacent (term1, term2)
      after (terms, terms).
```

PageRank - EMA

Es algoritmo de ordenamiento de "popularidad" de los sitios web. Básicamente se trata de ordenar los sitios web en función de cuál es más probable de que visites si hacés una caminata aleatoria por la web (estás en una página, entonces decidís irte a otra de forma random o seguir alguno de los links que están puestos en la pág.)

Dado esto, sabemos que cada página tiene un "rango" (lo que abarca, la popularidad, etc). También, sabemos que las págs tienen enlaces hacia afuera y hacia dentro (los 1eros son los links que tengo y los 2dos son los links de otras páginas hacia mí)

Dado esto, entonces, la popularidad de un sitio web es mayor mientras más enlaces hacia dentro tenga. Cada uno de estos tiene un peso, el cual es la "popularidad" del sitio web que tiene ese link. Por ello, un sitio web con más rango que me tiene como link me hace aumentar mucho más mi popularidad que si no abarca mucho.

Además, decimos que mi página tiene un rango alto si la suma de los rangos de las páginas que me apuntan tiene un valor alto

Básicamente:

- Las páginas apuntadas por muchas páginas web es más probable que sean visitadas y van a tener un PageRank más alto
- Las páginas apuntadas por páginas web con alto PageRank van a tener una mayor probabilidad de ser visitadas y entonces van a tener un mayor PageRank

DATAZOS

Retorno de la información

- El retorno de la información es el proceso de retornar documentos a partir de una colección de documentos en respuesta a una consulta.
 - Los documentos suelen estar en lenguaje natural no estructurado.

Bases de Datos Relacionales	Sistemas de retorno de la información	
Datos estructurados	Datos no estructurados	
Dirigidos por esquemas relacionales	No hay esquemas fijos	
Modelo de consultas estructurado	Modelos de consulta libres de forma	
Operaciones sobre metadatos	Operaciones sobre datos	
Las consultas retornan datos	Las búsquedas retornan lista de punteros a documentos.	
Los resultados se basan en correspondencia exacta y son siempre correctos.	Los resultados se basan en correspondencia aproximada y medidas de efectividad.	
Trabajan con transacciones	No trabajan con transacciones	

• La relevancia se basa en factores como:

- Frecuencia de términos: frecuencia de ocurrencia de término de una consulta en un documento.
- Frecuencia inversa de documentos: ¿en cuántos documentos ocurre la palabra? Si ocurre en menos documentos se le da más importancia a la palabra.
- Enlaces a documentos: Si hay más enlaces a un documento, el documento es más importante.

· La mayoría de SRI agregan:

- Las palabras que ocurren en el título, lista de autores, títulos se les da más importancia.
- Las palabras cuya primera ocurrencia es tarde en el documento se les da poca importancia.
- Proximidad: si las palabras de una consulta aparecen cerca entre sí en el documento, el documento tiene más importancia que si las palabras ocurren bien lejos unas de otra.
- Otra posibilidad es usar TF-IDF (frecuencia de término-frecuencia inversa de documento).
- TF-IDF es usado para evaluar la importancia de una palabra en una colección de documentos.
- Idea de TF-IDF: términos que capturan la esencia de un documento aparecen frecuentemente en el mismo; pero para que un término sea bueno para discriminar un documento de los demás, entonces debe ocurrir en unos pocos documentos de la colección.

Índices invertidos

Document 1

This example shows an example of an inverted index.

Document 2 Inverted index is a data structure for associating terms to documents.

Document 3 Stock market index is used for capturing the sestiments of the financial market.

ID	Term	Document: position
1.	example	1:2, 1:5
2.	inverted	1:8, 2:1
3.	index	1:9, 2:2, 3:3
4.	market	3:2, 3:13