

Búsqueda por Similitud

Profesor Heider Sanchez

P1. Búsqueda por Rango

Implementar en cualquier lenguaje de programación el algoritmo lineal de búsqueda por rango, el cual recibe como parámetro el objeto de consulta y un **radio de cobertura**. Luego usando la distancia Euclidiana (ED) se retorna todos los elementos que son cubiertos por el radio.

Algorithm RangeSearch(Q, r)	
1.	result = []
2.	for all objects C_i in the collection
3.	dist = ED (Q, C_i)
4.	if dist < r
5.	append(result, C_i)
6.	endif
7.	endfor
8.	return result

- Aplique la búsqueda para 3 elementos de la colección (Q_{15} , Q_{82} , Q_{121}) y para tres valores de radio ($r_1 < r_2 < r_3$).
- El objeto de consulta debe ser retirado de la colección antes de aplicar la búsqueda.
- Para saber qué valores de radio seleccionar, debe primero realizar un análisis de la distribución de las distancias computando N veces la distancia entre dos elementos aleatorios de la colección.
- Para evaluar la efectividad del resultado se debe usar la medida de Precisión ¿Cuántos de los objetos recuperados pertenecen a la misma categoría de la consulta?:

$$PR = \frac{\#ObjetosRelevantesRecuperados}{\#ObjetosRecuperados}$$

A continuación, se proporciona el cuadro que debe ser llenado por el alumno.

PR	Q_{300}	Q_{900}	Q_{1337}
$r_1 = 2$	0.9090909090909091	1.0	1.0
$r_2 = 6$	0.5429447852760736	0.5192307692307693	0.5858895705521472
$r_3 = 9$	0.46841294298921415	0.5299600532623169	0.5022288261515602

P2. Búsqueda KNN

Usando los mismos objetos de consulta del ejercicio anterior, implementar y aplicar el algoritmo lineal de búsqueda de los k vecinos más cercano (KNN) variando el k entre $\{2, 4, 8, 16, 32\}$.

Algorithm KnnSearch(Q, k)	
1.	result = []
2.	for all objects C_i in the collection
3.	dist = ED (Q, C_i)
4.	append(result, { C_i , dist})
5.	endfor
6.	orderByDist(result)
7.	return result[1:k]

**** Otra forma de implementación es gestionando la lista de resultado como una cola de prioridad máxima ¿Habría alguna mejora?**

Si, ya que no se tendría que ordenar la lista para poder devolverla sino que se van ordenando mientras se van insertando, lo que reduce el costo computacional.

PR	Q_{300}	Q_{900}	Q_{1337}
$k = 2$	1.0	1.0	1.0
$k = 4$	1.0	1.0	1.0
$k = 8$	0.875	1.0	1.0
$k = 16$	0.8125	1.0	1.0
$k = 32$	0.5625	0.84375	0.84375

Preguntas:

- 1- ¿Cuál de los dos métodos de búsqueda es más fácil de implementar si tengo todos los parámetros definidos? ¿Y cuál es más eficiente?

En mi opinión, la búsqueda por rango fue más fácil de implementar porque lo único que se usa es la ED y se compara con un r fijo, en cambio, en el KNN tuve que buscar la forma de ordenar la lista donde cada elemento es un par. Gracias a los resultados, pude notar que el KNN es más eficiente porque nos devuelve resultados con alta precisión en la mayoría de los casos

- 2- ¿Cuál de los dos métodos de búsqueda usted usaría en un ambiente real de recuperación de la información? Sustente su respuesta.

Aunque el KNN sea más eficaz al momento de recuperar data relevante, me parece que podría caer en el error de devolver elementos que no quiero en el caso de que cada elemento esté bien separado del otro. Siento que con la búsqueda de rango, limito a mi gusto el espacio de búsqueda y así recuperar elementos con la similaridad que yo quiera y no elementos que, si bien son los más cercanos, pueden estar a una distancia considerable.

