

Session 3: User Relevance Feedback

Q1 2023/2024

Àlex Domínguez Rodríguez

Mario Fernández Simón

Introducción

En esta práctica tendremos 2 objetivos principales. En primer lugar, programaremos en Elasticsearch un ciclo simple de “User Relevance Feedback”. En segundo lugar, experimentaremos y evaluaremos esta estrategia sobre un conjunto de documentos analizando distintos parámetros.

We will, we will Rocchio you

Después de seguir los pasos que nos marca el enunciado para implementar el archivo de python empezaremos con la experimentación, pero antes explicaremos porque hemos utilizado diccionarios en vez de vectores ordenados.

El motivo detrás de esta decisión reside en una mejor eficiencia al sumar diccionarios de gran tamaño que vectores. En el caso de los vectores, suponiendo que un vector tiene tamaño x_1 y otro tiene tamaño x_2 , sumarlos daría un coste de $O(x_1+x_2)$ y si repetimos este proceso k veces nos daría un coste de $O(k * (x_1 + x_2))$.

Por otro lado, el juntar 2 diccionarios no depende del número de elementos que haya sino del número de elementos únicos existentes. Esto origina que se reduzca drásticamente el tamaño de los elementos a juntar y más aún si son 2 textos de temas parecidos, ya que compartirán más palabras entre los 2. Por este motivo, hemos decidido utilizar diccionarios.

Experimentación

Para la experimentación iremos probando distintos valores para todos los parámetros y observaremos cómo cambia el resultado. Para el estudio hemos utilizado el corpus de *20_newsgroups/sci.med*.

Empezaremos viendo cómo cambia al tener valores de alfa y beta iguales, cuando alfa es más grande y en el caso contrario.

❖ Alfa = 1, Beta = 1, nrounds = 5, k = 4, r = 4, query = blood doctor

```
19 Documents
['blood^1.872782634297978', 'astemizole^1.7025951077952668', 'doctor^1.5832466020239246', 'sheryl^1.2779205480539872']
2 Documents
['astemizole^3.4051902155905336', 'sheryl^2.5558410961079745', 'blood^2.387867467654162', 'doctor^1.9870327111174109']
2 Documents
['astemizole^5.107785323385801', 'sheryl^3.833761644161962', 'blood^2.902952301010346', 'doctor^2.3908188202108973']
2 Documents
['astemizole^6.810380431181067', 'sheryl^5.111682192215949', 'blood^3.41803713436653', 'doctor^2.7946049293043838']
2 Documents
['astemizole^8.512975538976335', 'sheryl^6.389602740269936', 'blood^3.933121967722714', 'doctor^3.19839103839787']
```

Figura 1: Resultado de consulta.

❖ Alfa = 2, Beta = 0.5, nrounds = 5, k = 4, r = 4, query = blood doctor

```
19 Documents
['blood^2.4363913171489893', 'doctor^2.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936']
2 Documents
['blood^5.1303250509760705', 'doctor^4.785139656570668', 'astemizole^2.5538926616929003', 'sheryl^1.916880822080981']
2 Documents
['blood^10.518192518630233', 'doctor^9.772172367688079', 'astemizole^5.959082877283434', 'sheryl^4.472721918188956']
2 Documents
['blood^21.293927453938558', 'doctor^19.7462377899229', 'astemizole^12.769463308464502', 'sheryl^9.584404110404906']
2 Documents
['blood^42.84539732455521', 'doctor^39.694368634392546', 'astemizole^26.390224170826638', 'sheryl^19.807768494836804']
```

Figura 2: Resultado de consulta.

Vemos cómo al ser alfa más grande que beta, se mantiene tanto *blood* como *doctor* en las primeras posiciones, ya que estos resultados dependen más del resultado inicial y menos de los pesos calculados.

❖ Alfa = 0.5, Beta = 2, nrounds = 5, k = 4, r = 4, query = blood doctor

```
19 Documents
['astemizole^3.4051902155905336', 'sheryl^2.5558410961079745', 'potassium^2.270126810393689', 'blood^2.245565268595956']
0 Documents
['astemizole^1.7025951077952668', 'sheryl^1.2779205480539872', 'potassium^1.1350634051968445', 'blood^1.122782634297978']
0 Documents
['astemizole^0.8512975538976334', 'sheryl^0.6389602740269936', 'potassium^0.5675317025984222', 'blood^0.561391317148989']
0 Documents
['astemizole^0.4256487769488167', 'sheryl^0.3194801370134968', 'potassium^0.2837658512992111', 'blood^0.28069565857449']
0 Documents
['astemizole^0.21282438847440835', 'sheryl^0.1597400685067484', 'potassium^0.14188292564960556', 'blood^0.140347829287']
```

Figura 3: Resultado de consulta.

Aquí, por el contrario, vemos como *doctor* desaparece y aunque *blood* se mantiene se ve superado por otros términos como *potassium*.

❖ Alfa = 1, Beta = 0.5, nrounds = 3, k = 4, r = 4, query = blood doctor

```
19 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936']
2 Documents
['astemizole^1.7025951077952668', 'blood^1.693933733827081', 'doctor^1.493516355587057', 'sheryl^1.2779205480539872']
2 Documents
['astemizole^2.5538926616929003', 'blood^1.951476150505173', 'sheryl^1.916880822080981', 'doctor^1.6954094101054489']
```

Figura 4: Resultado de consulta.

La figura anterior muestra cómo con pocas iteraciones los términos blood y doctor van perdiendo fuerza poco a poco pero se mantienen. Esto también se debe a que hemos suavizado los valores de alfa y beta con unos que hemos visto que responden mejor después de probar diferentes opciones.

❖ Alfa = 1, Beta = 0.5, nrounds = 7, k = 4, r = 4, query = blood doctor

```
19 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936']
2 Documents
['astemizole^1.7025951077952668', 'blood^1.693933733827081', 'doctor^1.4935163555587057', 'sheryl^1.2779205480539872']
2 Documents
['astemizole^2.5538926616929003', 'blood^1.951476150505173', 'sheryl^1.916880822080981', 'doctor^1.6954094101054489']
2 Documents
['astemizole^3.4051902155905336', 'sheryl^2.5558410961079745', 'blood^2.2090185671832647', 'doctor^1.897302464652192']
2 Documents
['astemizole^4.256487769488167', 'sheryl^3.194801370134968', 'blood^2.4665609838613567', 'doctor^2.0991955191989353']
2 Documents
['astemizole^5.107785323385801', 'sheryl^3.8337616441619615', 'blood^2.7241034005394487', 'doctor^2.3010885737456785']
2 Documents
['astemizole^5.959082877283434', 'sheryl^4.472721918188955', 'blood^2.9816458172175406', 'doctor^2.5029816282924218']
```

Figura 5: Resultado de consulta.

Al aumentar el número de iteraciones, los términos iniciales sí que pierden importancia pero manteniendo un cierto peso.

❖ Alfa = 1, Beta = 0.5, nrounds = 4, k = 2, r = 4, query = blood doctor

```
19 Documents
['astemizole^1.7025951077952668', 'blood^1.515084833356184', 'doctor^1.4037861090934862', 'sheryl^1.2779205480539872']
2 Documents
['astemizole^3.4051902155905336', 'sheryl^2.5558410961079745', 'blood^2.030169666712368', 'doctor^1.8075722181869724']
2 Documents
['astemizole^5.107785323385801', 'sheryl^3.833761644161962', 'blood^2.545254500068552', 'doctor^2.211358327280459']
2 Documents
['astemizole^6.810380431181067', 'sheryl^5.111682192215949', 'blood^3.0603393334247357', 'doctor^2.6151444363739453']
```

Figura 6: Resultado de consulta.

❖ Alfa = 1, Beta = 0.5, nrounds = 4, k = 6, r = 4, query = blood doctor

```
19 Documents
['blood^1.3875059510202772', 'doctor^1.3037742718874608', 'astemizole^0.5675317025984222', 'potassium^0.54388454832348']
0 Documents
['blood^1.3875059510202772', 'doctor^1.3037742718874608', 'astemizole^0.5675317025984222', 'potassium^0.54388454832348']
0 Documents
['blood^1.3875059510202772', 'doctor^1.3037742718874608', 'astemizole^0.5675317025984222', 'potassium^0.54388454832348']
0 Documents
['blood^1.3875059510202772', 'doctor^1.3037742718874608', 'astemizole^0.5675317025984222', 'potassium^0.54388454832348']
```

Figura 7: Resultado de consulta.

Como se deduce de la fórmula, cuando la k es pequeña los pesos de los términos en el conjunto de k vectores tienen más importancia y si con un k grande pasa al revés.

❖ Alfa = 1, Beta = 0.5, nrounds = 4, k = 2, r = 2, query = blood doctor

```
19 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624']
19 Documents
['blood^1.872782634297978', 'doctor^1.5832466020239249']
19 Documents
['blood^2.309173951446967', 'doctor^1.8748699030358873']
19 Documents
['blood^2.745565268595956', 'doctor^2.1664932040478497']
```

Figura 8: Resultado de consulta.

❖ Alfa = 1, Beta = 0.5, nrounds = 4, k = 2, r = 6, query = blood doctor

```
19 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936',
sium^0.5675317025984222', 'seas.gwu.edu^0.47922020552024525']
0 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936',
sium^0.5675317025984222', 'seas.gwu.edu^0.47922020552024525']
0 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936',
sium^0.5675317025984222', 'seas.gwu.edu^0.47922020552024525']
0 Documents
['blood^1.436391317148989', 'doctor^1.2916233010119624', 'astemizole^0.8512975538976334', 'sheryl^0.6389602740269936',
sium^0.5675317025984222', 'seas.gwu.edu^0.47922020552024525']
```

Figura 9: Resultado de consulta.

Por último, cuando tenemos una r pequeña no se tiene la posibilidad de obtener nuevos términos o ver como los que hay se van ajustando con cada ronda. En cambio, al poner una r más grande no ocurre esto. Lo que sí observamos es que a partir de la segunda iteración no encontramos ningún documento. Esto puede deberse a que quizás las palabras buscadas en la query no tienen tanta relación entre ellas.

Conclusiones

Valores Alfa y Beta:

Para valores grandes de Alfa, vemos como el resultado se relaciona más con la consulta inicial, ya que este valor se multiplica en la fórmula por el valor del resultado inicial. Si es pequeño pasa todo lo contrario.

Para valores grandes de Beta, vemos como el resultado se relaciona con el peso de los vectores tf-idf de los documentos, esto provoca que aparezcan nuevos términos con más facilidad. Si es pequeño pasa todo lo contrario.

Por lo tanto, vemos como ambos parámetros están inversamente relacionados.

Valor Nrounds:

Con valores pequeños, se observa como el resultado de la primera consulta es muy parecido al resultado final, lo cual es lógico porque se ha aplicado menos veces la regla de Rocchio. Para *nrounds* más grandes vemos como sí que cambia hasta cierta ronda, donde a partir de ahí, ya los valores apenas se modifican.

Valores R y K:

Para valores de R grandes vemos como la consulta no devuelve ningún documento, ya que la Query se trata con una AND y al tener más términos que tener en cuenta es difícil que encuentre textos apropiados y menos en temas como medicina donde hay una gran variedad de textos que no tienen nada que ver entre ellos.

Por otro lado, se aprecia como la K está muy relacionada con el recall. Debido a que cuanto más grande sea su valor más documentos se recuperarán. No obstante, vemos como en nuestros experimentos ocurre todo lo contrario, se obtienen menos documentos que con K más pequeñas. Esto puede deberse a valores de otros parámetros y los temas de los que tratan los documentos. En nuestro caso, tenemos un $K = 6$ y $R = 4$, quizás con valores más pequeños de R podríamos haber encontrado más documentos al buscar menos términos en común.