

Scoring Regime	DCG	AP	RBP $p = 0.8$
TF-IDF	7.5691	0.8770	0.9057
BM25 $k_1=1.2$ $b=0.75$	7.6555	0.9059	0.9279
BM25 $k_1=2$ $b=0.75$	7.6454	0.9010	0.9248
BM25 $k_1=10$, $b=0.5$	7.3047	0.7944	0.8477

Hasil

Skema BM25 dengan parameter $k_1=1.2$ $b=0.75$ menghasilkan *result* yang paling bagus dibandingkan dengan skema TF-IDF. Akan tetapi *result* terburuk juga didapatkan oleh BM25 dengan parameter $k_1=10$, $b=0.5$.

Karena BM25 memiliki dua parameter yang bisa diubah, hasil dari scoring sangat bergantung pada kedua parameter ini. Pada skema BM25, k_1 mempengaruhi term frequency dan b mempengaruhi document length normalization. Bisa dilihat bahwa $k_1 = 1.2$ dan $b = 0.75$ memberikan hasil yang terbaik untuk BM25.

Untuk mendapatkan k_1 dan juga b yang optimal, dokumen juga perlu diperhatikan. Beberapa dokumen memiliki kata yang banyak berulang sehingga k_1 tidak bisa terlalu tinggi. Di sisi lain, beberapa kata yang berulang juga dapat membantu hasil search sehingga k_1 juga tidak boleh terlalu kecil. Ada juga panjang dokumen yang bervariasi sehingga variabel b juga perlu disesuaikan dengan panjang dokumen ini.

Ada juga skema TF-IDF yang dinormalisasi. Skema ini memberikan hasil uji sebagai berikut

Scoring Regime	DCG	AP	RBP $p = 0.8$
TF-IDF Normalized	0.7779	6.9490	0.6784

Bila kita memasukan normalized TF-IDF ke dalam skema yang lain, maka TF-IDF unnormalized memberikan hasil yang terburuk. Hal ini karena skema *normalized* membagi score dengan panjang document. Ini menjadi masalah karena suatu dokumen panjang yang memiliki string sama persis dengan query bisa saja kalah dengan dokumen pendek yang hanya memiliki beberapa kata yang sama, menghasilkan list hasil search yang kurang tepat.