Appendix

Table I shows the QQ metrics selected as features from the paper Predicting Query Quality for Applications of Text Retrieval to Software Engineering Tasks text[1]. Although we only selected nine of the QQ metrics as features, the two approaches we proposed still achieved good results.

TABLE I Selected Query Quality metrics

Measure	Description	Formula
AvgTFIDF	Average of the term frequency- inverse document frequency (tf- idf) ¹ values over all terms in the document	$\frac{1}{ Q_d } \sum_{q \in Q_d} t fidf(q, d)$
MAXTFIDF	Maximum of the term frequency- inverse document frequency (tf- idf) values over all terms in the document	$\max_{q \in Q_d} tfid\!f(q,d)$
DevTFIDF	The standard deviation of the term frequency-inverse document frequency (tf-idf) values over all terms in the document	$\sqrt{\frac{1}{ Q_d } \sum_{q \in Q_d} (tfidf(q,d) - AvgTFIDF)}$
${\bf AvgLogEntropy}$	Average LogEntropy ² values over all terms in the document	$\frac{1}{ Q_d } \sum_{q \in Q_d} LogEntropy(q, d)$
MedLogEntropy	Median LogEntropy values over all terms in the document	$egin{aligned} median \ LogEntropy(q,d) \ q \in Q_d \end{aligned}$
${\bf DevLogEntropy}$	The standard deviation of the LogEntropy values over all terms int the document	$\sqrt{\frac{1}{ Q_d } \sum_{q \in Q_d} (LogEntropy(q, d) - AvgLogEntropy)}$
SumSCQ	The sum of the collection-query similarity $(SCQ)^3$ over all terms in the document	$\sum_{q \in Q_d} (SCQ(q))$
AvgSCQ	The average of the collection- query similarity (SCQ) over all terms in the document	$\frac{1}{ Q_d } \sum_{q \in Q_d} SCQ(q)$
MaxSCQ	The maximum of the collection- query similarity (SCQ) over all terms in the document	$\frac{1}{ Q_d } \max_{q \in Q_d} SCQ(q)$

$$\begin{split} ^{1}tfidf(t,d) &= tf(t,d) \cdot \log \left(\frac{|D|}{|D_{t}|}\right) \\ ^{2}LogEntropy(t,d) &= \log(tf(t,d)+1) \cdot \left(1 + \frac{\sum_{d \in D_{t}} \frac{tf(t,d)}{tf(t,D)} \cdot \log \frac{tf(t,d)}{tf(t,D)}}{\log(|D|+1)}\right) \\ ^{3}SCQ(t) &= (1 + \log((tf,D))) \cdot \log \left(\frac{|D|}{|D_{t}|}\right) \end{split}$$

q - a term in the document;

 Q_d - the set of terms in the document d;

tf(t,d) - the frequency of term t in d;

D - the set of documents in the collection;

 D_t - the set of documents containing term t;

d - a document in the document collection D;

tf(t,D) - the frequency of term t in all docs;

References

[1] Chris Mills, Gabriele Bavota, Sonia Haiduc, Rocco Oliveto, Andrian Marcus, and Andrea De Lucia. 2017. Predicting Query Quality for Applications of Text Retrieval to Software Engineering Tasks. ACM Trans. Softw. Eng. Methodol. 26, 1, Article 3 (May 2017), 45 pages.