

## The 21 Pre-Retrieval Query Measures

Property	Measure	Description	Formula
Specificity	<i>AvgIDF</i>	Average of the Inverse Document Frequency ( <i>idf</i> ) values over all query terms	$\frac{1}{ Q } \sum_{q \in Q} idf(q)$
	<i>MaxIDF</i>	Maximum of the Inverse Document Frequency ( <i>idf</i> ) values over all query terms	$\max_{q \in Q} (idf(q))$
	<i>DevIDF</i>	The standard deviation of the Inverse Document Frequency ( <i>idf</i> ) values over all query terms	$\sqrt{\frac{1}{ Q } \sum_{q \in Q} (idf(q) - avgIDF)^2}$
	<i>AvgICTF</i>	Average Inverse Collection Term Frequency ( <i>ictf</i> ) values over all query terms	$\frac{1}{ Q } \sum_{q \in Q} ictf(q)$
	<i>MaxICTF</i>	Maximum Inverse Collection Term Frequency ( <i>ictf</i> ) values over all query terms	$\max_{q \in Q} (ictf(q))$
	<i>DevICTF</i>	The standard deviation of the Inverse Collection Term Frequency ( <i>ictf</i> ) values over all query terms	$\sqrt{\frac{1}{ Q } \sum_{q \in Q} (ictf(q) - avgICTF)^2}$
	<i>AvgEntropy</i>	Average <i>entropy</i> values over all query terms	$\frac{1}{ Q } \sum_{q \in Q} entropy(q)$
	<i>MedEntropy</i>	Median <i>entropy</i> values over all query terms	$median_{q \in Q} (entropy(q))$
	<i>MaxEntropy</i>	Maximum <i>entropy</i> values over all query terms	$\max_{q \in Q} (entropy(q))$
	<i>DevEntropy</i>	The standard deviation of the <i>entropy</i> values over all query terms	$\sqrt{\frac{1}{ Q } \sum_{q \in Q} (entropy(q) - avgEntropy)^2}$
	<i>QS</i>	Query Scope – the percentage of documents in the collection containing at least one of the query terms	$\frac{ \bigcup_{q \in Q} D_q }{ D }$
	<i>SCS</i>	Simplified Clarity Score – the Kullback-Leiber divergence of the query language model from the collection language model	$\sum_{q \in Q} p_q(Q) \cdot \log\left(\frac{p_q(Q)}{p_q(D)}\right)$
Coherency	<i>AvgVAR</i>	Average of the variances of the query term weights over the documents containing the query term ( <i>VAR</i> ), over all query terms	$\frac{1}{ Q } \sum_{q \in Q} VAR(q)$
	<i>MaxVAR</i>	Maximum of the variances of the query term weights over the documents containing the query term ( <i>VAR</i> ), over all query terms	$\max_{q \in Q} (VAR(q))$
	<i>SumVAR</i>	Sum of the variances of the query term weights over the documents containing the query term ( <i>VAR</i> ), over all query terms	$\sum_{q \in Q} VAR(q)$
	<i>CS</i>	Coherence Score – the average of the pairwise similarity between all pairs of documents containing one of the query terms ( <i>cs</i> ) among all	$\frac{1}{ Q } \sum_{q \in Q} cs(q)$
Similarity	<i>AvgSCQ</i>	The average of the collection-query similarity ( <i>SCQ</i> ) over all query terms	$\frac{1}{ Q } \sum_{q \in Q} SCQ(q)$
	<i>MaxSCQ</i>	The maximum of the collection-query similarity ( <i>SCQ</i> ) over all query terms	$\max_{q \in Q} (SCQ(q))$
	<i>SumSCQ</i>	The sum of the collection-query similarity ( <i>SCQ</i> ) over all query terms	$\sum_{q \in Q} SCQ(q)$
Term relatedness	<i>AvgPMI</i>	Average Pointwise Mutual Information ( <i>PMI</i> ) over all pairs of terms in the query	$\frac{2( Q -2)!}{( Q )!} \sum_{q_1, q_2 \in Q} PMI(q_1, q_2)$
	<i>MaxPMI</i>	Maximum Pointwise Mutual Information ( <i>PMI</i> ) over all pairs of terms in the query	$\max_{q_1, q_2 \in Q} (PMI(q_1, q_2))$
$idf(t) = \log\left(\frac{ D }{ D_t }\right)$		$p_t(d) = \frac{tf(t,d)}{ d }$	$w(t, d) = \frac{1}{ d } \log(1 + tf(t, d)) \cdot idf(t)$
$ictf(t) = \log\left(\frac{ D }{tf(t,D)}\right)$		$p_t(D) = \frac{tf(t,D)}{ D }$	$SCQ(t) = (1 + \log(ctf(t, D))) \cdot idf(t)$
$\bar{w}_t = \frac{1}{ D_t } \sum_{d \in D_t} w(t, d)$		$p_t(Q) = \frac{tf(t,Q)}{ Q }$	$cs(t) = \frac{\sum_{(d_i, d_j) \in D_t} sim(d_i, d_j)}{ D_t  \cdot ( D_t  - 1)}$
			$entropy(t) = \sum_{d \in D_t} p_t(d) \cdot \log_{ D } p_t(d)$
			$PMI(t_1, t_2) = \log \frac{p_{t_1, t_2}(D)}{p_{t_1}(D) \cdot p_{t_2}(D)}$
			$VAR(t) = \sqrt{\frac{\sum_{d \in D_t} (w(t, d) - \bar{w}_t)^2}{df(t)}}$

$Q$  – the set of query terms;  $q$  – a term in the query;  $D$  – 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;  $tf(t, d)$  – the frequency of term  $t$  in  $d$   
 $tf(t, Q)$  – the frequency of term  $t$  in the query;  $sim(d_i, d_j)$  – the cosine similarity between the vector-space representations of  $d_i$  and  $d_j$