

## The 7 Post-Retrieval Query Measures

Property	Measure	Description	Formula or Algorithm
<b>Robustness</b>	<i>Subquery Overlap</i>	Captures the extent of the overlap between the result set retrieved by the entire query and the result sets retrieved by individual query terms. The lower the standard deviation, the better the query.	i) Run the original query $q$ , and obtain the result list $R$ ii) Run each individual query term $q_i$ in the original query as a separate query and obtain the result list $R_i$ . iii) For each individual query term $q_i$ , compute the overlap between the first $k$ ( $k=10$ ) documents in $R$ and the first $k$ documents in $R_i$ (i.e., number of documents found in both result lists) iv) The overall score of the query is the standard deviation of the values of the overlap considered for each term in the query.
	<i>Robustness Score</i>	The terms' weights in the top relevant documents are slightly perturbed and the resulting documents are re-ranked. The correlation between the initial rank and that after modification is considered. The higher the robustness score, the better the query.	i) Run the original query $q$ , and obtain the result list $R$ ii) Take the top 50 documents in $R$ and consider them as ranked list $L$ iii) For each document $d$ in $L$ , get a perturbed document $d'$ from $d$ in the following way: <ul style="list-style-type: none"> <li>• All terms <math>t</math> from the corpus that do not appear in document <math>d</math>, will not be included in <math>d'</math> neither.</li> <li>• All terms <math>t</math> from the corpus that appear in document <math>d</math> with frequency <math>f</math>, but do not appear in the query will appear in document <math>d'</math> with the same frequency <math>f</math>.</li> <li>• Each term <math>t</math> that appears in <math>d</math> with frequency <math>f</math> and appears also in the query <math>q</math> will appear also in <math>d'</math>, but with a frequency <math>f'</math>, which is a random number obtained from a Poisson distribution <math>P(\lambda)</math> with <math>\lambda = f</math></li> </ul> iv) The new 50 documents obtained are ranked according to the query $q$ , resulting in a second ranked list $L'$ , where each document corresponds to a document in $L$ . v) Compute the Spearman rank correlation between the positions of the 50 documents in $L$ and the positions of their corresponding perturbed documents in $L'$ and record the correlation obtained. vi) Repeat steps iii) – v) 100 times, and the final robustness score is the average Spearman correlation between the 100 runs.
	<i>First Rank Change</i>	Captures the probability of a document found on the first position in the list of results to still remain on the first position after a perturbation is applied to it. The higher the score, the better the query.	i) Run the original query $q$ , and obtain the result list $R$ ii) Take the top 50 documents in $R$ and consider them as ranked list $L$ iii) For each document $d$ in $L$ , get a perturbed document $d'$ from $d$ in the following way: <ul style="list-style-type: none"> <li>• All terms <math>t</math> from the corpus that do not appear in document <math>d</math>, will not be included in <math>d'</math> neither</li> <li>• All terms <math>t</math> from the corpus that appear in document <math>d</math> with frequency <math>f</math>, but do not appear in the query will appear in document <math>d'</math> with the same frequency <math>f</math>.</li> <li>• Each term <math>t</math> that appears in <math>d</math> with frequency <math>f</math> and appears also in the query <math>q</math> will appear also in <math>d'</math>, but with a frequency <math>f'</math>, which is a random number obtained from a Poisson distribution <math>P(\lambda)</math> with <math>\lambda = f</math></li> </ul> iv) The new 50 documents obtained are ranked according to the query $q$ , resulting in a second ranked list $L'$ , where each document corresponds to a document in $L$ . v) record a 1 if the top ranked document in $L$ is also the top ranked document (after perturbation) in $L'$ , and record 0 otherwise. vi) Repeat steps iii) – v) 100 times, and the final score is the sum of the values (0 or 1) obtained in step v) for all the 100 runs.
	<i>Clustering Tendency</i>	Measures the cohesion of the top retrieved documents as the textual similarity between them. The higher the clustering tendency the better the query.	$CT = Mean \left( \frac{Sim_{query}(d_{mp}, d_{nn} q)}{Sim_{query}(p_{sp}, d_{mp} q)} \right) * \frac{1}{T} \sum_{i=1}^T (x_i - y_i)$ <p> <math>q</math> - the query  <math>p_{sp}</math> – the sampled point, i.e., a randomly chosen document from the corpus, which does not appear in the top 100 documents  <math>d_{mp}</math> - the marked point, i.e., the document, inside the top 100 documents in the ranked list, with largest similarity with the sampled point  <math>d_{nn}</math> - the nearest neighbor of the marked point within the top 100 documents </p>

			<p>in the ranked list</p> <p><math>x_i</math> - the maximum weight for a term <math>i</math> across the top 100 retrieved documents</p> <p><math>y_i</math> - the minimum weight for a term <math>i</math> across the top 100 retrieved documents</p> <p>The mean in the CT formula is computed for 100 randomly sampled points (i.e., the similarity formulas are computed 100 times, each time with a different randomly sampled point).</p> $Sim_{query}(d_i, d_j   q) = \frac{\sum_{k=1}^T d_{ik} d_{jk}}{\sqrt{\sum_{k=1}^T d_{ik}^2 \sum_{k=1}^T d_{jk}^2}} * \frac{\sum_{k=1}^T c_k q_k}{\sqrt{\sum_{k=1}^T c_k^2 \sum_{k=1}^T q_k^2}}$ <p><math>d_i</math> and <math>d_j</math> - the two documents</p> <p><math>T</math> - the number of unique terms in the collection</p> <p><math>q</math> - the query (with weight <math>q_k</math> for term <math>k</math>) - the weight of a term in the query or a document is its tfidf</p> <p><math>c</math> - the vector of terms common to both <math>d_i</math> and <math>d_j</math> with weights <math>c_k</math> being the average of <math>d_{ik}</math> and <math>d_{jk}</math>.</p>
	<i>Spatial Autocorrelation</i>	<p>Changes the retrieval-scores of each top relevant document as the average of the scores of its most similar documents. Then, the linear correlation of the new scores with original ones is used. The higher the spatial autocorrelation the better the query.</p>	<p>i) Run the original query <math>q</math>, and obtain the result list <math>R</math></p> <p>ii) Take the top 50 documents in <math>R</math> and consider them as ranked list <math>L</math></p> <p>iii) For each document <math>d</math> in <math>L</math>, compute the cosine similarity between <math>d</math> and the rest of the documents in <math>L</math>, using tf-idf as the weight of the terms in the document vectors.</p> <p>iv) Among the documents in <math>L</math>, select the 5 documents that are most similar to <math>d</math> according to the cosine similarity.</p> <p>v) Let <math>s</math> be the score of document <math>d</math> in <math>L</math>. Assign a new score to <math>d</math>, which is the average score of the 5 most similar documents to it according to the cosine similarity.</p> <p>vi) Perform the above steps for each document <math>d</math> in <math>L</math></p> <p>vii) The Pearson correlation between the original scores of the documents in <math>L</math> and the derived scores of those documents represents the index of spatial autocorrelation.</p>
<b>Score distribution</b>	<i>Weighted Information Gain (WIG)</i>	<p>Measures the divergence between the mean retrieval score of top-ranked documents and that of the entire corpus. The hypothesis is that the more similar these documents are to the query, with respect to the query similarity exhibited by a general non-relevant document (i.e., the corpus), the more effective the retrieval. The higher the weighted information gain, the better the query.</p>	$WIG(q) = \frac{1}{k} \sum_{d \in D_q^k} \sum_{t \in q} \lambda(t) \log \frac{\Pr(t d)}{\Pr(t D)}$ <p><math>q</math> - query</p> <p><math>t</math> - a term in the query <math>q</math></p> <p><math>D</math> - set of all documents in corpus</p> <p><math>D_q</math> - the set of documents in the result set to query <math>q</math></p> <p><math>D_q^k</math> - the top <math>k</math> documents in the result list</p> <p><math>k</math> - the number of top documents to consider</p> <p><math> q </math> - number of terms in the query <math>q</math></p> $\lambda(t) = \frac{1}{\sqrt{ q }}$
	<i>Normalized Query Commitment (NCQ)</i>	<p>Measures the standard deviation of retrieval scores in <math>D_q^k</math>, normalized by the score of the whole collection. The higher NCQ, the better the query.</p>	$NCQ = \frac{\sqrt{\frac{1}{k} \sum_{d \in D_q^k} (Score(d) - \mu)^2}}{Score(D_q)}$ <p><math>k</math> - the number of top documents to consider. Best performance was obtained with <math>k=100</math>.</p> <p><math>D_q^k</math> - The top <math>k</math> documents from the result list returned in response to query <math>q</math>.</p> <p><math>Score(d)</math> - the score obtained by document <math>d</math> in <math>D_q^k</math></p> <p><math>D_q</math> - the set of all results returned in response to query <math>q</math></p> <p><math>Score(D_q)</math> - the sum of the scores of all the documents in the result list returned by the IR technique</p> $\mu = \frac{1}{k} \sum_{d \in D_q^k} Score(d)$ <p>NCQ - represents the normalized standard deviation of the retrieval scores in <math>D_q^k</math></p>