

# Normalized Query Commitment Revisited

Haggai Roitman  
IBM Research - Haifa  
haggai@il.ibm.com

## ABSTRACT

We revisit the Normalized Query Commitment (NQC) query performance prediction (QPP) method. To this end, we suggest a scaled extension to a discriminative QPP framework and use it to analyze NQC. Using this analysis allows us to redesign NQC and suggest several options for improvement.

## ACM Reference Format:

Haggai Roitman. 2019. Normalized Query Commitment Revisited. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '19)*, July 21–25, 2019, Paris, France. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3331184.3331334>

## 1 INTRODUCTION

Query performance prediction (QPP) is a core IR task whose primary goal is to assess retrieval quality in the absence of relevance judgements [1]. In this work, we revisit Shtok et al.'s *Normalized Query Commitment* (NQC) QPP method [10]. NQC is a state-of-the-art post-retrieval QPP method [1], based on document retrieval scores variance analysis. Nowadays, the NQC method serves as a common competitive baseline in many QPP works [4, 6, 8, 11].

We first present the NQC method and shortly discuss the motivation behind it. We then shortly present Roitman et al.'s discriminative QPP framework [8]. Using a scaled extension to this framework, we analyze the NQC method, “deconstructing” it into its most basic parts. Based on this analysis, we revise NQC's design and suggest several options for improvement. Using an extensive evaluation over common TREC corpora, we demonstrate that, by redesigning NQC, where we extend it with more proper calibration and scaling, we are able to significantly improve its prediction accuracy.

## 2 BACKGROUND

### 2.1 Normalized Query Commitment

Let  $q$  denote a query and let  $D$  denote the respective ranked-list of top- $k$  documents in corpus  $C$  with the highest retrieval scores. Let  $s(d)$  further denote document  $d$ 's ( $\in D$ ) retrieval score with respect to  $q$ . The NQC method estimates the query's performance according to the standard-deviation of  $D$ 's document retrieval scores,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

SIGIR '19, July 21–25, 2019, Paris, France

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6172-9/19/07...\$15.00

<https://doi.org/10.1145/3331184.3331334>

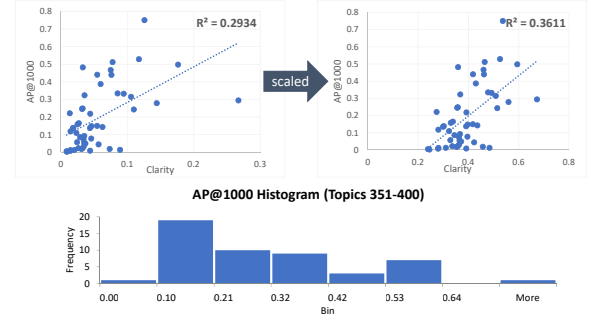


Figure 1: Example of scaling Clarity's values (TREC Robust queries 351-400,  $\gamma = 0.3$ )

further normalized by the corpus score<sup>1</sup>  $s(C)$ , formally:

$$NQC(D|q) \stackrel{def}{=} \frac{\sqrt{\frac{1}{k} \sum_{d \in D} (s(d) - \hat{\mu}_D)^2}}{s(C)}, \quad (1)$$

where  $\hat{\mu}_D \stackrel{def}{=} \frac{1}{k} \sum_{d \in D} s(d)$  is  $D$ 's mean document score. The key idea behind NQC is the assumption that the mean score  $\hat{\mu}_D$  may serve as a *pseudo-ineffective reference (score) point* of the retrieval. The more the retrieval scores deviate from this point, the less chance is assumed for a query-drift within  $D$ ; and hence, a more qualitative retrieval will be predicted [10]. The corpus score  $s(C)$  further serves as a query-sensitive normalization term, allowing to compare NQC values across different queries [10].

### 2.2 Discriminative QPP

In this work, we build on top of Roitman et al.'s discriminative QPP framework [8]. In [8], the authors have shown that many of the previously suggested post-retrieval QPP methods (e.g., Clarity [2], WIG [12] and SMV [11]) share the following basic form:

$$WPM(D|q) \stackrel{def}{=} \frac{1}{k} \sum_{d \in D} s(d) \cdot \phi_F(d), \quad (2)$$

where  $\phi_F(d) \stackrel{def}{=} \prod_j (f_j(d))^{\alpha_j}$  is a *Weighted Product Model* (WPM)

discriminative calibrator; with  $f_j(d)$  represents some retrieval quality feature and  $\alpha_j \geq 0$  denotes its relative importance [8]. Within this framework,  $\phi_F(d)$  calibrates each document  $d$ 's ( $\in D$ ) retrieval score  $s(d)$  according to the likelihood of  $d$  being a relevant response to query  $q$  [8]. To this end,  $\phi_F(d)$  may encode various retrieval quality properties, such as properties of  $q$ ,  $C$ ,  $D$  and the document  $d$  itself. As Roitman et al. have pointed out, some of such properties may be complementing each other (e.g., query vs. corpus quality

<sup>1</sup>Such a score is obtained by treating the corpus as a single (large) document.

effects), and therefore, tradeoffs in the design of general QPP methods should be properly modeled [8].  $\phi_F(d)$ , therefore, models such tradeoffs (i.e., using the weights  $\alpha_j$ ).

### 2.3 Scaled Weighted Product Model

While most of existing QPP methods are designed to predict a given quality measure (with AP@1000 being the most commonly used measure [1]), the relationship between a given predictor's estimates and actual quality numbers may not be necessarily linear. As an example, the bottom part of Figure 1 depicts an histogram of the AP@1000 values obtained by a query-likelihood (QL) based-retrieval method, which was applied over TREC Robust topics 351-400. As we can observe, there is a high variability in quality values, having query difficulty non-uniformly distributed.

To address such variability during prediction, we now suggest a simple, yet effective extension to [8], which scales the calibrated-mean estimator defined in Eq. 2, as follows:

$$SWPM(D|q) \stackrel{\text{def}}{=} \left( \frac{1}{k} \sum_{d \in D} s(d) \cdot \phi_F(d) \right)^\gamma, \quad (3)$$

where  $\gamma \geq 0$  is a scaling parameter. Note that, whenever  $\gamma < 1$ , higher variability in prediction values is encouraged. Going back to Figure 1, its upper part further illustrates the relationship between Clarity's [2] predicted values and actual quality numbers before and after applying scaling ( $\gamma = 0.3$ ). As we can observe, after scaling, Clarity's prediction accuracy has significantly improved.

### 3 NQC REVISITED

We now show that, the NQC method can be derived as a scaled calibrated-mean estimator. To this end, we first rewrite NQC (defined in Eq. 1), as follows:

$$NQC(D|q) \stackrel{\text{def}}{=} \left[ \frac{1}{k} \sum_{d \in D} s(d) \left( \frac{1}{s(C)} \right)^2 \left( \frac{s(d) - \hat{\mu}_D}{\sqrt{s(d)}} \right)^2 \right]^{0.5} \quad (4)$$

Using Eq. 4, it now becomes apparent that, NQC may be treated as an instance of the generic scaled estimator defined in Eq. 3. In NQC's case, we have:  $f_1(d) = \frac{1}{s(C)}$  and  $\alpha_1 = 2$ ,  $f_2(d) = \frac{s(d) - \hat{\mu}_D}{\sqrt{s(d)}}$  and  $\alpha_2 = 2$ , and a scaling parameter  $\gamma = 0.5$ .

Reformulating the NQC method as a scaled calibrated-mean estimator has two main advantages. First, being represented as a discriminative QPP method, we can potentially improve its prediction accuracy by applying supervised-learning to better tune its calibration (i.e.,  $\alpha_1$  and  $\alpha_2$ ) and scaling (i.e.,  $\gamma$ ) parameters. As we shall later demonstrate, such a more fine-tuned calibration may indeed result in a better prediction quality.

Second, and more important, similar to [8], such a representation allows to analyze the existing design of NQC's calibration features (i.e.,  $f_1(d)$  and  $f_2(d)$ ), potentially even redesigning some of them. Among these two features,  $f_1(d)$  was previously examined in [8], where it was shown to be utilized as a calibration feature within the SMV method [11]. We, therefore, next focus our attention on the second calibration feature of  $f_2(d)$ .

### 3.1 Analysis of $f_2(d)$

Using the new NQC formulation, we now try to explain the motivation behind its second calibration feature  $f_2(d) = \frac{s(d) - \hat{\mu}_D}{\sqrt{s(d)}}$ . To this end, we "break"  $f_2(d)$  into two main parts. The first, given by  $s(d) - \hat{\mu}_D$ , measures the difference between a given document  $d$ 's ( $\in D$ ) score and that of the pseudo-ineffective reference (score) point, considered by NQC as the mean document score  $\hat{\mu}_D$ . The larger the difference is, the more will document  $d$ 's own score be estimated as an informative evidence for a true view of its relevance (i.e.,  $s(d) > \hat{\mu}_D$ ) or irrelevance (i.e.,  $s(d) < \hat{\mu}_D$ ). The second part, given by  $\sqrt{s(d)}$ , basically serves as a scaler, allowing to measure score differences with respect to some comparable scale. In NQC's case, scaling is simply performed using the document score itself.

### 3.2 Redesigning $f_2(d)$

We now generalize the basic form of  $f_2(d)$ , and reformulate it as follows:

$$f_2(d) \stackrel{\text{def}}{=} \frac{s(d) - s_{neg}}{\sqrt{Z_d}}, \quad (5)$$

where  $s_{neg}$  represents some pseudo-ineffective reference (score) point for comparison, and  $Z_d$  denotes some scaler value. Therefore, for the original NQC method we have the configuration of  $s_{neg} = \hat{\mu}_D$  and  $Z_d = s(d)$ .

We next suggest several other configurations, based on alternative  $s_{neg}$  and  $Z_d$  instantiations. Starting with  $Z_d$ , as a first alternative, we suggest  $s_{max} \stackrel{\text{def}}{=} \max_{d \in D} s(d)$ . As a second alternative, motivated by Ozdemiray et al.'s **ScoreRatio** QPP method [6], we also suggest  $s_{min} \stackrel{\text{def}}{=} \min_{d \in D} s(d)$ . As a third alternative, we suggest a range-sensitive scaler:  $s_{max} - s_{min} + \epsilon$ , where  $\epsilon \stackrel{\text{def}}{=} 10^{-10}$  is used to avoid dividing by zero. Noting that both  $s_{max}$  and  $s_{min}$  may actually represent outlier values, as a fourth and final alternative, which tries to reduce such outliers, we suggest the following inter-percentile range scaler  $Q_{.95} - Q_{.05} + \epsilon$ . Here,  $Q_{.x}$  represents the value of the  $x\%$  percentile of  $D$ 's observed score distribution.

Next, we suggest several alternative estimates of  $s_{neg}$ . As a first alternative, we suggest  $s_{min}$ , which represents the document with the lowest relevance score, and hence, presumably the least effective sample reference. As a second alternative, following [4] we consider  $s_{neg}$  to be the mean score of the documents at the lower ranks of  $D$ . Formally, for a given  $l \leq k$ , we estimate  $\hat{\mu}_{D_{neg}} \stackrel{\text{def}}{=} \frac{1}{k-l} \sum_{i=l+1}^k s(d_i)$ , where  $d_i$  denotes the document that is ranked at position  $i$  in  $D$ .

Our third and final alternative is motivated by Diaz's **Autocorrelation** QPP method [3]. Using the opposite of the log<sup>2</sup> that was suggested in [3], which is based on the *Cluster-Hypothesis* in IR [5], we expect an informative document score  $s(d)$  to be quite different from those scores of documents that are the most dissimilar to  $d$ . Formally, let  $dist(d, d')$  denote the distance between two documents in  $D$  and let  $KFN(d)$  be the set of  $K$ -farthest neighbors

<sup>2</sup>The Autocorrelation method treats similar documents as pseudo-effective reference points, and hence, the score of a document is actually supposed to be similar to the scores of its closest neighbors.

of  $d$  in  $D$  according to  $\text{dist}(d, d')$ . We then define  $s_{neg}$  as the mean score of documents in  $KFN(d)$ .

## 4 EVALUATION

### 4.1 Experimental setup

**4.1.1 Datasets.** For our evaluation we have used the following TREC corpora and topics: **TREC4** (201-250), **TREC5** (251-300), **AP** (51-150), **WSJ** (151-200), **ROBUST** (301-450, 601-700), **WT10g** (451-550) and **GOV2** (701-850). These datasets were used by many previous QPP works [1–3, 7–12]. Titles of TREC topics were used as queries, except for TREC4, where we used topic descriptions instead [8]. We used the Apache Lucene<sup>3</sup> open source search library for indexing and searching documents. Documents and queries were processed using Lucene’s English text analysis (i.e., tokenization, lowercasing, Porter stemming and stopwords). For retrieval, we used Lucene’s language-model-based cross-entropy scoring with Dirichlet smoothed document language models, where the smoothing parameter was set to 1000 [8, 10].

**4.1.2 Baseline predictors.** We compared the original NQC method (see again Eq. 1) and its proposed extensions to several different baseline QPP methods. As a first line of baselines, we compared with the following state-of-the-art post-retrieval QPP methods:

- **Clarity** [2]: estimates query performance relatively to the divergence between the relevance model induced from a given (retrieved) list and the corpus background model.
- **WIG** [12]: estimates query performance according to the difference between the average document retrieval score in a given list and that of the corpus  $s(C)$ .
- **WEG** [4]: a WIG-like alternative that uses  $\hat{\mu}_{D_{neg}}$  as the pseudo-ineffective reference point instead of  $s(C)$ .
- **Autocorrelation** [3]: based on the *Cluster-Hypothesis in IR*, it estimates query performance according to the Pearson’s- $r$  correlation between  $D$ ’s original document scores and those estimated by interpolating the scores of each document  $d$ ’s ( $\in D$ )  $K$ -nearest neighbors relatively to that document’s similarity with each neighbor.
- **ScoreRatio** [6]: simply estimates query performance according to the ratio  $\frac{s_{max}}{s_{min}}$ .
- **SMV** [11]: is a direct alternative to NQC that also considers score magnitude and variance, estimated as:

$$SMV(D|q) \stackrel{\text{def}}{=} \frac{1}{k \cdot s(C)} \sum_{d \in D} s(d) \left| \ln \frac{s(d)}{\hat{\mu}_D} \right|.$$

Using our derivation of NQC as a scaled calibrated QPP method (see Section 3), we further evaluated various alternatives of NQC, as follows. We first evaluated various  $f_2(d)$  configurations within NQC, i.e., **NQC**( $s_{neg}, Z_d$ ). To this end, we instantiated  $s_{neg}$  and  $Z_d$  according to our proposed alternatives.

Next, using the original NQC configuration (i.e.,  $s_{neg} = \hat{\mu}_D$  and  $Z_d = s(d)$ ), we also evaluated a calibration-only version of NQC (**C-NQC**), where we only tuned its  $\alpha_1$  and  $\alpha_2$  parameters, while still fixing the scaling parameter to  $\gamma = 0.5$ . In a similar manner, we evaluated a scaled-only version of NQC (**S-NQC**), where we fixed  $\alpha_1 = 2$  and  $\alpha_2 = 2$  and only tuned the scaling parameter  $\gamma$ . We

further evaluated a combined predictor (**SC-NQC**), where all the three parameters were tuned. Finally, we evaluated a pre-tuned **SC-NQC** employed with the best  $f_2(d)$  configuration learned for the non-calibrated/non-scaled NQC versions, denoted **SC-NQC**(best).

**4.1.3 Setup.** On each setting, we predicted the performance of each query with respect to its top-1000 retrieved documents [1]. We assessed prediction over queries quality according to the correlation between the predictor’s values and the actual average precision (AP@1000) values calculated using TREC’s relevance judgments [1]. To this end, we report the Pearson’s- $r$  (P- $r$ ) and Kendall’s- $\tau$  (K- $\tau$ ) correlations, which are the most commonly used measures [1].

Most of the methods that we have evaluated required to tune some free parameters. First, following the common practice [1], for each QPP method, we tuned  $k$  – the number of documents used for prediction<sup>4</sup>; with  $k \in \{5, 10, 20, 50, 100, 150, 200, 500, 1000\}$ . For **Clarity**, we induced a relevance model using the top- $m$  ( $\in \{5, 10, 20, 50, 100\}$ ) documents in  $D$  and further applied clipping at the top- $n$  terms cutoff, with  $n \in \{10, 20, 50, 100\}$ . For the **Autocorrelation** and **NQC**( $KFN(K), \cdot$ ) baselines, we tuned  $K \in \{3, 5, 10, 20, 30\}$ , further using the *Bhattacharyya* distance between the unsmoothed language models of the documents in  $D$  as the (dis)similarity measure  $\text{dist}(\cdot)$  of choice. To realize  $\hat{\mu}_{D_{neg}}$  in **SMV** and **NQC**, we further tuned  $l \in \{5, 10, 20, 50, 100, 200, 500\}$ .

To learn the calibration feature weights (i.e.,  $\alpha_1$  and  $\alpha_2$ ) and scaling parameter (i.e.,  $\gamma$ ) of the NQC variants, following [7–9], we used a *Coordinate Ascent* approach. To this end, we selected the feature weights over the grid of  $[0, 5] \times [0, 5] \times [0, 1]$ , in steps of 0.1 per dimension. Following [7–10], we trained and tested all methods using a holdout (2-fold cross validation) approach. On each dataset, we generated 30 random splits of the query set; each split had two folds. We used the first fold as the (query) train set. We kept the second fold for testing. We recorded the average prediction quality over the 30 splits. Finally, we measured statistical significant differences of prediction quality using a two-tailed paired t-test with  $p < 0.05$  computed over all 30 splits.

### 4.2 Results

We report our evaluation results in Table 1. As a “stand-alone” QPP method, even the original NQC method (i.e., **NQC**( $\hat{\mu}_D, s(d)$ )) already provides highly competitive prediction quality results compared to the other state-of-the-art QPP methods. We, therefore, next evaluate the impact of our proposed NQC extensions.

**4.2.1 Evaluation of various  $f_2(d)$  configurations.** We start with a **qualitative** examination of the relative contribution of each of the two parts of  $f_2(d)$  (i.e.,  $s_{neg}$  and  $Z_d$ ). To this end, we count the relative number of cases, per dataset and quality measure (i.e., P- $r$  or K- $\tau$ ), in which using a specific configuration part has resulted in a better prediction accuracy than that of the original NQC’s configuration. Overall, in 136 out of the  $14 \times 19 = 266$  possible cases, utilizing one of the alternative configurations has resulted in a better prediction.

Among the three alternative  $s_{neg}$  options, the relative preference was:  $KFN(d)$  (45/70),  $\hat{\mu}_{D_{neg}}$  (38/70) and  $s_{min}$  (29/70). This

<sup>3</sup><http://lucene.apache.org>

<sup>4</sup>All NQC variants were tuned with the **same** value of  $k$ .

**Table 1: Evaluation results.** “greenish” and “reddish” colored values represent an improvement or a decline in prediction accuracy compared to the original NQC method. A statistical significant difference between a given NQC variant and the original NQC are marked with \* ( $p < .05$ ).

	TREC4		TREC5		AP		WSJ		ROBUST		WT10g		GOV2	
	P-r	K- $\tau$	P-r	K- $\tau$	P-r	K- $\tau$	P-r	K- $\tau$	P-r	K- $\tau$	P-r	K- $\tau$	P-r	K- $\tau$
Autocorrelation	.456	.366	.188	.136	.348	.226	.586	.496	.385	.321	.299	.198	.247	.198
ScoreRatio	.420	.444	.344	.204	.231	.166	.620	.411	.405	.370	.226	.275	.295	.246
Clarity	.401	.357	.314	.208	.413	.265	.500	.355	.404	.329	.289	.203	.206	.164
WIG	<b>.533</b>	<b>.502</b>	.347	.252	.613	.417	.685	.463	.560	.399	.221	.323	.498	.352
WEG	.349	.454	.296	.201	.499	.357	.696	.482	.562	.400	.462	<b>.383</b>	.419	.337
SMV	.458	.386	.414	.332	.536	.379	.713	.484	.534	.391	.466	.279	.419	.310
NQC( $\hat{\mu}_D, s(d)$ ) (Shtok et al. [10])	.491	.383	.454	.275	.619	.405	.722	.510	.580	.422	.522	.330	.378	.253
NQC( $\hat{\mu}_D, \max$ )	.484	.371	.455	.277	.628	.413	.723	.494	.582	<b>.424</b>	.518	.324	<b>.386*</b>	.256
NQC( $\hat{\mu}_D, \min$ )	.486	.403	.468	.265	<b>.592*</b>	<b>.390*</b>	<b>.689*</b>	<b>.486*</b>	.540	.425	<b>.538*</b>	<b>.337*</b>	<b>.350*</b>	<b>.240*</b>
NQC( $\hat{\mu}_D, \max - \min$ )	<b>.436*</b>	<b>.325*</b>	<b>.417*</b>	<b>.301*</b>	<b>.645*</b>	<b>.414*</b>	<b>.706*</b>	<b>.460*</b>	<b>.585*</b>	<b>.405*</b>	<b>.407*</b>	<b>.237*</b>	<b>.460*</b>	<b>.320*</b>
NQC( $\hat{\mu}_D, Q_{.95} - Q_{.05}$ )	<b>.451*</b>	<b>.335*</b>	<b>.406*</b>	<b>.297*</b>	<b>.643*</b>	<b>.386*</b>	.715	.496	<b>.440*</b>	<b>.409*</b>	<b>.384*</b>	<b>.274*</b>	<b>.421*</b>	<b>.294*</b>
NQC( $\hat{\mu}_{Dneg}, s(d)$ )	.482	.393	.503*	.304*	.625*	.423*	.729*	.504	.567*	.412*	.523	.333	.416*	.292*
NQC( $\hat{\mu}_{Dneg}, \max$ )	<b>.474*</b>	.381	<b>.495*</b>	.314	<b>.632*</b>	<b>.422*</b>	<b>.734*</b>	.513	<b>.568*</b>	<b>.409*</b>	.516	.325	<b>.427*</b>	<b>.300*</b>
NQC( $\hat{\mu}_{Dneg}, \min$ )	.483	.412*	.518*	.306*	.615	.407	.705*	.499*	.541*	.420*	<b>.539*</b>	<b>.334*</b>	<b>.384*</b>	<b>.280*</b>
NQC( $\hat{\mu}_{Dneg}, \max - \min$ )	<b>.415*</b>	<b>.320*</b>	<b>.427*</b>	<b>.299*</b>	<b>.635*</b>	.409	<b>.734*</b>	<b>.496*</b>	<b>.551*</b>	<b>.377*</b>	<b>.399*</b>	<b>.233*</b>	<b>.507*</b>	<b>.355*</b>
NQC( $\hat{\mu}_{Dneg}, Q_{.95} - Q_{.05}$ )	<b>.443*</b>	<b>.344*</b>	<b>.438*</b>	<b>.336*</b>	.678*	.406	.736*	.502*	<b>.465*</b>	<b>.394*</b>	<b>.383*</b>	<b>.273*</b>	<b>.504*</b>	<b>.363*</b>
NQC( $\min, s(d)$ )	<b>.466*</b>	.393	<b>.506*</b>	<b>.312*</b>	<b>.602*</b>	<b>.423*</b>	.715	.496	<b>.538*</b>	<b>.393*</b>	<b>.406*</b>	<b>.277*</b>	<b>.445*</b>	<b>.306*</b>
NQC( $\min, \max$ )	<b>.459*</b>	<b>.374</b>	<b>.498*</b>	<b>.312*</b>	.617	<b>.430*</b>	.723	.494	<b>.536*</b>	<b>.388*</b>	<b>.391*</b>	<b>.265*</b>	<b>.456*</b>	<b>.318*</b>
NQC( $\min, \min$ )	.484	<b>.425*</b>	<b>.522*</b>	<b>.317*</b>	<b>.591*</b>	<b>.424*</b>	.714	.500	<b>.529*</b>	<b>.408*</b>	<b>.429*</b>	<b>.291*</b>	<b>.412*</b>	<b>.296*</b>
NQC( $\min, \max - \min$ )	<b>.380*</b>	<b>.311*</b>	<b>.430*</b>	<b>.304*</b>	<b>.630*</b>	<b>.412*</b>	<b>.706*</b>	<b>.460*</b>	<b>.493*</b>	<b>.329*</b>	<b>.245*</b>	<b>.176*</b>	<b>.528*</b>	<b>.358*</b>
NQC( $\min, Q_{.95} - Q_{.05}$ )	<b>.407*</b>	<b>.338*</b>	<b>.442*</b>	<b>.330*</b>	<b>.674*</b>	<b>.412*</b>	.715	.496	<b>.457*</b>	<b>.353*</b>	<b>.255*</b>	<b>.205*</b>	<b>.532*</b>	<b>.367*</b>
NQC( $KFN(d), s(d)$ )	<b>.475*</b>	<b>.422*</b>	<b>.524*</b>	<b>.297*</b>	<b>.638*</b>	<b>.445*</b>	<b>.730*</b>	<b>.561*</b>	.580	<b>.413*</b>	<b>.506*</b>	<b>.333*</b>	<b>.432*</b>	<b>.327*</b>
NQC( $KFN(d), \max$ )	<b>.470*</b>	<b>.418*</b>	<b>.519*</b>	<b>.306*</b>	<b>.643*</b>	<b>.449*</b>	<b>.739*</b>	<b>.561*</b>	.583	<b>.414</b>	<b>.500*</b>	<b>.324*</b>	<b>.443*</b>	<b>.327*</b>
NQC( $KFN(d), \min$ )	.488	<b>.434*</b>	<b>.530*</b>	<b>.286*</b>	<b>.635*</b>	<b>.428*</b>	<b>.691*</b>	<b>.523*</b>	<b>.556*</b>	.422	<b>.528*</b>	<b>.351*</b>	<b>.400*</b>	<b>.310*</b>
NQC( $KFN(d), \max - \min$ )	<b>.415*</b>	<b>.349</b>	<b>.474*</b>	<b>.317*</b>	<b>.646*</b>	<b>.452*</b>	<b>.760*</b>	<b>.558*</b>	<b>.572*</b>	<b>.394*</b>	<b>.385*</b>	<b>.253*</b>	<b>.515*</b>	<b>.374*</b>
NQC( $KFN(d), Q_{.95} - Q_{.05}$ )	<b>.448*</b>	<b>.371</b>	<b>.479*</b>	<b>.330*</b>	<b>.688*</b>	<b>.453*</b>	<b>.744*</b>	<b>.582*</b>	<b>.504*</b>	<b>.403*</b>	<b>.380*</b>	<b>.272*</b>	<b>.509*</b>	<b>.373*</b>
C-NQC( $\hat{\mu}_D, s(d)$ )	.492	<b>.412*</b>	.455	<b>.310*</b>	<b>.622*</b>	<b>.412*</b>	<b>.737*</b>	.510	.580	<b>.412</b>	<b>.559*</b>	<b>.377*</b>	<b>.507*</b>	<b>.360*</b>
S-NQC( $\hat{\mu}_D, s(d)$ )	.496	.383	<b>.473*</b>	.275	.620	.405	<b>.728*</b>	<b>.502*</b>	.580	.422	.522	.330	<b>.398*</b>	.253
SC-NQC( $\hat{\mu}_D, s(d)$ )	<b>.504*</b>	<b>.412*</b>	<b>.473*</b>	<b>.310*</b>	<b>.623*</b>	<b>.412*</b>	<b>.737*</b>	.510	.580	.422	<b>.559*</b>	<b>.377*</b>	<b>.527*</b>	<b>.360*</b>
SC-NQC(best)	<b>.510*</b>	<b>.457*</b>	<b>.560*</b>	<b>.317*</b>	<b>.691*</b>	<b>.454*</b>	<b>.770*</b>	<b>.567*</b>	<b>.585*</b>	<b>.405*</b>	<b>.561*</b>	<b>.380*</b>	<b>.566*</b>	<b>.401*</b>

demonstrates that, a better choice of a pseudo-ineffective reference point within NQC should be one that is more sensitive to the document in mind  $d$  (i.e.,  $KFN(d)$ ). This in comparison to a point that is more generally estimated (i.e.,  $\hat{\mu}_{Dneg}$  and  $s_{\min}$ ). Moreover, among the latter two alternatives, considering more than one sample from the lower ranks of  $D$  is a better choice. In a similar manner, among the four alternative  $Z_d$  options, the relative preference was:  $s_{\max}$  (34/56),  $s_{\min}$  (29/56),  $s_{\max} - s_{\min}$  (25/56) and lastly  $Q_{.95} - Q_{.05}$  (23/56). This further demonstrates that, a better choice for a scaler is one that depends on a single point (having  $s_{\max}$  a better choice than  $s_{\min}$ ), rather than a range of values. Finally, **quantitatively**, by more proper configuration of  $f_2(d)$ , NQC’s prediction accuracy has improved in all datasets, up to 38% and 46% improvement in P-r and K- $\tau$ , respectively. This demonstrates the usefulness of analyzing NQC using our scaled extension to Roitman et al.’s discriminative QPP framework [8].

**4.2.2 Effect of calibration and scaling.** The four bottom rows in Table 1 further report the effect of NQC’s calibration and scaling. First, we observe that, by better tuning of NQC’s calibration features within C-NQC( $\hat{\mu}_D, s(d)$ ), its prediction accuracy has improved in most cases (up to 42% better). Next, scaling (i.e., S-NQC( $\hat{\mu}_D, s(d)$ )) by itself also improves NQC’s accuracy (up to 5%). Combining both calibration and scaling (i.e., SC-NQC( $\hat{\mu}_D, s(d)$ )), has resulted in most cases in a further improved accuracy (up to 42%). Finally, further using the best  $f_2(d)$  configuration together with calibration and scaling (i.e., SC-NQC(best)) provides the best prediction strategy for NQC (up to 60% improvement).

## 5 CONCLUSIONS

We introduced a simple, yet highly effective, extension to Roitman et al.’s discriminative QPP framework [8]. Our main focus was on the NQC method, where using our proposed extension, we were able to redesign it and suggest several options for improvement.

## REFERENCES

- [1] David Carmel and Oren Kurland. Query performance prediction for ir. In *Proceedings of SIGIR '12*.
- [2] Steve Cronen-Townsend, Yun Zhou, and W. Bruce Croft. Predicting query performance. In *Proceedings of SIGIR '02*.
- [3] Fernando Diaz. Performance prediction using spatial autocorrelation. In *Proceedings of SIGIR '07*.
- [4] Ahmad Khwileh, Andy Way, and Gareth J. F. Jones. Improving the reliability of query expansion for user-generated speech retrieval using query performance prediction. In *CLEF*, 2017.
- [5] Oren Kurland. The cluster hypothesis in information retrieval. In *Advances in Information Retrieval*, pages 823–826. Springer International Publishing, 2014.
- [6] A. M. Ozdemiray and Ismail S. Altıngöve. Query performance prediction for aspect weighting in search result diversification. *Proceedings of CIKM '14*.
- [7] Haggai Roitman. An extended query performance prediction framework utilizing passage-level information. In *Proceedings of ICTIR*, pages 35–42, New York, NY, USA, 2018. ACM.
- [8] Haggai Roitman, Shai Erera, Oren Sar-Shalom, and Bar Weiner. Enhanced mean retrieval score estimation for query performance prediction. In *Proceedings of ICTIR '17*.
- [9] Haggai Roitman, Shai Erera, and Bar Weiner. Robust standard deviation estimation for query performance prediction. In *Proceedings of ICTIR '17*.
- [10] Anna Shtok, Oren Kurland, David Carmel, Fiana Raiber, and Gad Markovits. Predicting query performance by query-drift estimation. *ACM Trans. Inf. Syst.*, 30(2):11:1–11:35, May 2012.
- [11] Yongquan Tao and Shengli Wu. Query performance prediction by considering score magnitude and variance together. In *Proceedings of CIKM '14*.
- [12] Yun Zhou and W. Bruce Croft. Query performance prediction in web search environments. In *Proceedings of SIGIR '07*.