

Kullback-Leibler Divergence Revisited

Fiana Raiber
Yahoo Research, Israel
fiana@yahoo-inc.com

Oren Kurland
Technion, Israel
kurland@ie.technion.ac.il

ABSTRACT

The KL divergence is the most commonly used measure for comparing query and document language models in the language modeling framework to ad hoc retrieval. Since KL is rank equivalent to a specific weighted geometric mean, we examine alternative weighted means for language-model comparison, as well as alternative divergence measures. The study includes analysis of the inverse document frequency (IDF) effect of the language-model comparison methods. Empirical evaluation, performed with different types of queries (short and verbose) and query-model induction approaches, shows that there are methods that often outperform the KL divergence in some settings.

KEYWORDS

language models, weighted geometric mean

1 INTRODUCTION

Comparing a language model induced from the query with that induced from the document is a standard ranking approach in the language modeling framework to ad hoc document retrieval [20]. The Kullback-Leibler (KL) divergence has been the most commonly used measure for language-model comparison, as it is a natural choice for comparing probability distributions.

The KL divergence is rank equivalent to the cross entropy measure [20] which is in turn rank equivalent to a specific weighted geometric mean [19, 20]: that of the probabilities assigned to terms in the support of the query model¹ by the document language model; the probabilities assigned to these terms by the query language model serve as weights in the mean. Given the rank equivalence between the KL divergence and this weighted geometric mean, we study alternative weighted means for comparing the query and document language models; namely, the arithmetic and harmonic means, integrations of these with the geometric mean, and generalized versions of the means. In addition, we study alternative divergence measures.

Since using the KL divergence for ranking has an IDF (inverse document frequency) effect, we study this effect for the alternative language-model-comparison methods that we consider. The IDF effect holds if the impact of a term on the retrieval score becomes smaller when its corpus frequency increases.

¹The support is the set of terms assigned a non-zero probability by the language model.

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.

ICTIR '17, October 1–4, 2017, Amsterdam, The Netherlands

© 2017 ACM. ISBN 978-1-4503-4490-6/17/10...\$15.00

DOI: <http://dx.doi.org/10.1145/3121050.3121062>

We perform extensive empirical evaluation of the various measures for comparing query and document language models using five TREC datasets. We vary the types of queries used (short titles vs. verbose descriptions+titles) and the query-model induction approach: unsmoothed maximum likelihood estimate (MLE) vs. pseudo-feedback-based query-model induction; specifically, we use the relevance model [21] and the mixture model [31].

We found that there are measures that can often outperform the KL divergence for certain types of queries and query-model induction methods. For example, the Power mean [6], which generalizes the harmonic and arithmetic means, and which under certain conditions converges to the geometric mean, outperforms the KL divergence when using MLEs induced from short title queries or pseudo-feedback-based query models. When using MLEs induced from verbose queries, KL is the best performing.

Our key contributions can be summarized as follows:

- We study various weighted means and divergence measures for comparing query and document language models; our study includes analysis of the IDF effect.
- We perform an extensive empirical evaluation and demonstrate that in some settings there are measures that often outperform the commonly used KL divergence.

2 RELATED WORK

Most work on using language models for ad hoc retrieval has focused on improving the language models induced from the query and a document rather than on the measure used to compare them [30]. There are studies of using various divergence measures to compare language models in natural-language-processing tasks; e.g. [23]. However, we are not aware of such in-depth studies for the ad hoc retrieval task. In contrast, document-query similarity estimates were compared in the vector space model [33].

The use of KL divergence for ad hoc retrieval was formally supported in [25]. KL was also shown to outperform a document likelihood approach for relevance-model-based ranking [21]. We show that there are measures more effective than KL for that end.

Axiomatic analysis of retrieval methods in the language modeling framework focused on the query likelihood model [14], translation models [18] and pseudo-feedback-based query models and their smoothing [9, 15]. For example, term frequency, document frequency and document length axioms have been devised and their empirical merits were demonstrated. In contrast, we evaluate different measures for comparing a query model with a document model and analyze whether they have an IDF effect.

3 COMPARING QUERY AND DOCUMENT LANGUAGE MODELS

We address the ad hoc retrieval task: ranking documents d in corpus \mathcal{D} in response to query q . Let θ_q and θ_d denote the unigram

language models induced from q and d , respectively. Henceforth, these are referred to as query and document (language) models, respectively. Various language-model induction methods have been used in work on ad hoc retrieval [30].

To rank documents in the corpus, the document and query language models can be compared [20]. Since the language models are probability distributions defined over the vocabulary, a natural similarity measure, which is the most commonly used in work on ad hoc retrieval, is the Kullback-Leibler (KL) divergence. As higher KL values correspond to decreased similarity, the negative KL divergence is used:

$$-KL(\theta_q, \theta_d) \stackrel{def}{=} - \sum_w p(w|\theta_q) \log \frac{p(w|\theta_q)}{p(w|\theta_d)} = H(\theta_q) - CE(\theta_q, \theta_d). \quad (1)$$

$H(\theta_q) \stackrel{def}{=} - \sum_w p(w|\theta_q) \log p(w|\theta_q)$ is the entropy of the query language model and $CE(\theta_q, \theta_d) \stackrel{def}{=} - \sum_w p(w|\theta_q) \log p(w|\theta_d)$ is the cross entropy between the query and document models. (As is the case for KL, higher values of CE correspond to decreased similarity.) Hereinafter, summations as that in Eq. 1 are applied only over terms w for which $p(w|\theta_q) > 0$ [21]; i.e., terms in the *support* of the query language model. Since the entropy of the query model does not affect ranking, KL and CE are rank equivalent: $KL(\theta_q, \theta_d) \stackrel{rank}{=} CE(\theta_q, \theta_d)$.

If an unsmoothed maximum likelihood estimate (MLE) is used for the query language model, then ranking produced using -CE (and therefore using -KL) is equivalent to that produced using the query likelihood method [20, 29]. Formally, $p(w|\theta_x^{MLE}) \stackrel{def}{=} \frac{\text{tf}(w \in x)}{|x|}$ is the MLE of term w with respect to the text (or text collection) x where $\text{tf}(w \in x)$ is the number of occurrences of w in x ; $|x| \stackrel{def}{=} \sum_{w \in x} \text{tf}(w \in x)$ is x 's length. Then, the following holds:

$$\begin{aligned} -CE(\theta_q^{MLE}, \theta_d) &= \sum_w p(w|\theta_q^{MLE}) \log p(w|\theta_d) \\ &= \frac{1}{|q|} \log \prod_w p(w|\theta_d)^{\text{tf}(w \in q)}; \end{aligned} \quad (2)$$

$\prod_w p(w|\theta_d)^{\text{tf}(w \in q)}$ is q 's likelihood with respect to d .

3.1 The IDF Effect

An important implication of the rank equivalence between (negative) cross entropy used with a query MLE and the query likelihood method is the *IDF effect*. That is, Zhai and Lafferty [32] showed that for query likelihood, the effect of query terms on ranking is inversely related to their corpus frequency. Hence, this IDF (inverse document frequency) effect also governs ranking using cross entropy (and KL) with a query MLE.

An interesting question which was somewhat overlooked in past literature is whether using KL or CE with query models which are not maximum likelihood estimates results in an IDF effect.² Indeed, while MLE is the standard choice for inducing a query model using only the query terms, there are query models with a much broader support, e.g., pseudo-feedback-based query models [21, 31].

²In contrast, there has been work on the importance of selecting and increasing the probability in the query language model of terms with high IDF values (e.g., [9, 15]).

We now show that regardless of the query model utilized, using CE (and hence KL) for ranking results in an IDF effect. The formal argument is similar to that used to demonstrate the IDF effect for the query likelihood model [32]. We provide the details so as to later on contrast the IDF effect entailed by using KL and CE with that entailed by using alternative measures studied below.

We assume a standard smoothed unigram document language model θ_d . For the two most commonly used smoothing methods in work on ad hoc retrieval, Dirichlet and Jelinek-Mercer [32], the document language model is $\theta_d \stackrel{def}{=} (1 - \alpha_d)\theta_d^{MLE} + \alpha_d\theta_{\mathcal{D}}^{MLE}$; for Dirichlet, $\alpha_d = \frac{\mu}{|d| + \mu}$ where μ is a parameter; for Jelinek-Mercer, α_d is simply a constant. For $w \in d$ we define $p_s(w|\theta_d) \stackrel{def}{=} p(w|\theta_d) = (1 - \alpha_d)p(w|\theta_d^{MLE}) + \alpha_d p(w|\theta_{\mathcal{D}}^{MLE})$. For $w \notin d$, we define $p_u(w|\theta_d) \stackrel{def}{=} p(w|\theta_d) = \alpha_d p(w|\theta_{\mathcal{D}}^{MLE})$; 's' and 'u' stand for "seen" and "unseen", respectively [32]. The negative CE is then:

$$-CE(\theta_q, \theta_d) = \sum_w p(w|\theta_q) \log p(w|\theta_d) \quad (3a)$$

$$= \sum_{w \in d} p(w|\theta_q) \log p_s(w|\theta_d) + \sum_{w \notin d} p(w|\theta_q) \log p_u(w|\theta_d) \quad (3b)$$

$$= \sum_{w \in d} p(w|\theta_q) \log \frac{p_s(w|\theta_d)}{p_u(w|\theta_d)} + \sum_w p(w|\theta_q) \log p_u(w|\theta_d) \quad (3c)$$

$$= \sum_{w \in d} p(w|\theta_q) \log \frac{p_s(w|\theta_d)}{\alpha_d p(w|\theta_{\mathcal{D}}^{MLE})} + \log \alpha_d - CE(\theta_q, \theta_{\mathcal{D}}^{MLE}) \quad (3d)$$

$$\stackrel{rank}{=} \sum_{w \in d} p(w|\theta_q) \log \left(1 + \frac{1 - \alpha_d}{\alpha_d} \frac{p(w|\theta_d^{MLE})}{p(w|\theta_{\mathcal{D}}^{MLE})} \right) + \log \alpha_d \quad (3e)$$

The transitions are based on separating groups of indices (Eq. 3b), ignoring document independent factors for ranking (Eq. 3e) and using definitions and arithmetic manipulations.

The summation in Eq. 3e is over terms in the query-model support that appear in d . It constitutes the "document-query match" which is based on the probability of query terms in the query ($p(w|\theta_q)$), document ($p(w|\theta_d^{MLE})$) and corpus ($p(w|\theta_{\mathcal{D}}^{MLE})$) models. The latter yields the "IDF effect" as it regularizes the impact of term occurrence in the document by occurrence in the corpus. Document length is used as a normalizer in $p(w|\theta_d^{MLE})$. If Dirichlet smoothing is applied, then α_d also (inversely) depends on document length.

3.2 Weighted Means

The negative cross entropy amounts to (cf. [19, 20]):

$$-CE(\theta_q, \theta_d) = \sum_w p(w|\theta_q) \log p(w|\theta_d) \stackrel{rank}{=} \prod_w p(w|\theta_d)^{p(w|\theta_q)}.$$

That is, negative CE is rank equivalent to a weighted geometric mean (henceforth **Geo**), $Geo(\theta_d; \theta_q)$, of the probabilities assigned to terms in θ_q 's support by the document model θ_d ; the probabilities assigned by θ_q serve as weights whose sum is 1.

The view of (negative) CE as a weighted geometric mean gives rise to the question, which to the best of our knowledge this paper is the first to explore, of whether using other means to compare the query and document language models can improve retrieval effectiveness. There are numerous weighted means (e.g., [4, 6]). In what follows we mainly focus on the three classical and most widely

used means: arithmetic, geometric and harmonic, their integrations and generalizations. Table 1 presents the weighted means we study.

All these means, as well as the divergence measures we discuss in Section 3.3, include a “query-document match” component³ which relies on occurrence of terms from the query model support in the document; document length normalizes this occurrence since MLE-based estimates are used. Hence, in addition to discussing general properties of the means and divergence measures, we study whether their query-document match component has an IDF effect. Analyzing more complicated connections between document length, IDF and term frequency (cf. [14]) is left for future work.

The weighted arithmetic mean (**Ari**) is less “conservative” than the weighted geometric mean applied by CE; i.e., it is less affected by outliers. Specifically, a low probability assigned by a document model to one of the terms in the query model support incurs higher penalty in Geo than in Ari. Using manipulations similar to those in Eq. 3, Ari amounts to:

$$\text{Ari}(\theta_d; \theta_q) \stackrel{\text{rank}}{=} (1 - \alpha_d) \sum_{w \in d} p(w|\theta_q)p(w|\theta_d^{MLE}) + \alpha_d \sum_w p(w|\theta_q)p(w|\theta_D^{MLE}).$$

The document-query match (the first summation) has no IDF effect.

The weighted harmonic mean (**Har**), which is more conservative than Geo, is the third classical mean. Using manipulations similar to those in Eq. 3, we arrive at:

$$\text{Har}(\theta_d; \theta_q) \stackrel{\text{rank}}{=} \left(\sum_{w \in d} \frac{(\alpha_d - 1)p(w|\theta_q)p(w|\theta_d^{MLE})}{\alpha_d p_s(w|\theta_d)p(w|\theta_D^{MLE})} + \frac{1}{\alpha_d} \sum_w \frac{p(w|\theta_q)}{p(w|\theta_D^{MLE})} \right)^{-1}.$$

The document-query match (first summation) has an IDF effect since $\alpha_d < 1$ and $p_s(w|\theta_d)$ increases with increasing values of corpus frequency. Note that the second summation has an “inverse” IDF effect (i.e., DF effect) which can be considered a drawback of Har: the more frequent a query term in the corpus, the higher the retrieval score.⁴ However, if α_d is the same for all documents, then the second summation has no effect on ranking, and the entire retrieval score (used to rank documents) exhibits an IDF effect as was the case for the KL divergence in the transition from Eq. 3d to Eq. 3e. Indeed, using Jelinek-Mercer smoothing or assuming equal document lengths in analyzing the IDF effect yields a constant α_d . Using the equal document lengths assumption is in line with Fang and Zhai’s methodology of analyzing the IDF effect [13]: assuming two documents of the same length, the document with more occurrences of query terms that are less frequent in the corpus should receive a higher retrieval score. Accordingly, hereinafter, and as mentioned above, our analysis of means and divergence measures will focus on whether the query-document match has an IDF effect.

The next two means integrate the geometric mean employed by negative CE with the arithmetic mean (**GeoAri**) and the harmonic mean (**GeoHar**). These means are computed in iterations that are guaranteed to converge. In the first iteration, $i = 1$, we initialize $g_1 = \text{Geo}(\theta_d; \theta_q)$, $a_1 = \text{Ari}(\theta_d; \theta_q)$ and $h_1 = \text{Har}(\theta_d; \theta_q)$. Then, for

³Other components are document-independent or α_d ; sometimes, these two interact.
⁴This is also the case for the second summation in Ari.

Table 1: Weighted means. β and γ are free parameters.

M	$M(\theta_d; \theta_q)$	IDF Effect
Ari	$\sum_w p(w \theta_q)p(w \theta_d)$	\times
Har	$\left(\sum_w \frac{p(w \theta_q)}{p(w \theta_d)} \right)^{-1}$	\checkmark
GeoAri	Integration of $\text{Geo}(\theta_d; \theta_q)$ and $\text{Ari}(\theta_d; \theta_q)$	\checkmark
GeoHar	Integration of $\text{Geo}(\theta_d; \theta_q)$ and $\text{Har}(\theta_d; \theta_q)$	\checkmark
Power	$\left(\sum_w p(w \theta_q)p(w \theta_d)^\beta \right)^{\frac{1}{\beta}}$	$\beta < 1$
Lehmer	$\frac{\sum_w p(w \theta_q)p(w \theta_d)^\gamma}{\sum_w p(w \theta_q)p(w \theta_d)^{\gamma-1}}$	$0 \leq \gamma < 1$

GeoAri, in iteration $i + 1$ we set $g_{i+1} = \sqrt{a_i g_i}$ and $a_{i+1} = \frac{1}{2}(a_i + g_i)$ until g_{i+1} and a_{i+1} converge to the same value.⁵ Similarly, for GeoHar, we compute $g_{i+1} = \sqrt{h_i g_i}$ and $h_{i+1} = 2(h_i^{-1} + g_i^{-1})^{-1}$ until convergence.⁶ Since Geo and Har have the IDF effect, and Ari does not, GeoHar has the effect while GeoAri has it to a somewhat limited extent. The following inequality holds for the means considered thus far:

$$\begin{aligned} \text{Ari}(\theta_d; \theta_q) &\geq \text{GeoAri}(\theta_d; \theta_q) \geq \text{Geo}(\theta_d; \theta_q) \\ &\geq \text{GeoHar}(\theta_d; \theta_q) \geq \text{Har}(\theta_d; \theta_q). \end{aligned}$$

The weighted **Power** mean is essentially a family of means that covers a wide range of aggregates bounded by the minimal and maximal values aggregated: $\lim_{\beta \rightarrow -\infty} \text{Power}(\theta_d; \theta_q) = \min_w p(w|\theta_d)$ and $\lim_{\beta \rightarrow \infty} \text{Power}(\theta_d; \theta_q) = \max_w p(w|\theta_d)$. For $\beta = 1$ and $\beta = -1$ Power amounts to the arithmetic and harmonic means, respectively. Furthermore, $\lim_{\beta \rightarrow 0} \text{Power}(\theta_d; \theta_q) = \text{Geo}(\theta_d; \theta_q)$. (Recall that Geo and negative CE are rank equivalent.) Using a similar transition to that from Eq. 3b to Eq. 3c it can be shown that the query-document match component in Power has an IDF effect for $\beta < 1$. In Section 4.2 we show that β close to 0 yields the best retrieval performance which also often transcends that of using Geo (CE).

As additional reference comparison we consider **Lehmer**, which is a special case of the Gini family of means [4]. For $\gamma = 1$ and $\gamma = 0$, Lehmer amounts to the arithmetic and harmonic means, respectively. It can be shown that the query-document match component in Lehmer has an IDF effect for $0 \leq \gamma < 1$, but not for $1 \leq \gamma \leq 2$. The analysis for $\gamma < 0$ and $\gamma > 2$ is quite involved and might require numerical simulation as the query-document match component and document-independent factor have mutual effects.

3.3 Divergence Measures

KL is a measure of the divergence between two probability distributions — query and document language models in our case. We therefore now turn to examine alternative divergence measures. These measures, presented in Table 2, were used in a variety of tasks, including term co-occurrence estimation [23], topic-based document segmentation [5], story link detection [7], query performance prediction [3] and static index pruning [8]. We rank documents by ascending values of the divergence measures.

⁵ g_{i+1} is the geometric mean of g_i and a_i , while a_{i+1} is their arithmetic mean.

⁶ g_{i+1} is the geometric mean of g_i and h_i , while h_{i+1} is their harmonic mean.

Table 2: Divergence measures. $p(w|\theta_{qd}^{[\eta]}) \stackrel{def}{=} \eta p(w|\theta_q) + (1 - \eta)p(w|\theta_d)$; η is a free parameter.

D	$D(\theta_q, \theta_d)$	IDF Effect
Hellinger	$\sqrt{\sum_w (\sqrt{p(w \theta_q)} - \sqrt{p(w \theta_d)})^2}$	✓
TotalVariation	$\sum_w p(w \theta_q) - p(w \theta_d) $	✗
JensenShannon	$KL(\theta_q, \theta_{qd}^{[\eta]}) + KL(\theta_d, \theta_{qd}^{[\eta]})$; $\eta = \frac{1}{2}$	✓
J	$KL(\theta_q, \theta_d) + KL(\theta_d, \theta_q)$	✓
ResistorAverage	$(KL(\theta_q, \theta_d)^{-1} + KL(\theta_d, \theta_q)^{-1})^{-1}$	✓
χ^2 Neyman	$\sum_w \frac{(p(w \theta_q) - p(w \theta_d))^2}{p(w \theta_d)}$	✓
χ^2 Pearson	$\sum_w \frac{(p(w \theta_q) - p(w \theta_d))^2}{p(w \theta_q)}$	✓
χ^2 Symmetric	$\sum_w \frac{(p(w \theta_q) - p(w \theta_d))^2}{p(w \theta_q) + p(w \theta_d)}$	✓
Skew	$KL(\theta_q, \theta_{qd}^{[\eta]})$	✓

Most of the measures we consider are special cases of the f-divergence [11]⁷. Given a convex function f defined over $(0, \infty)$ such that $f(1) = 0$, the f-divergence between θ_q and θ_d is:

$$\sum_w p(w|\theta_d) f\left(\frac{p(w|\theta_q)}{p(w|\theta_d)}\right).$$

Different choices of f result in different divergence measures. For example, setting $f(x) = x \log x$ results in the KL divergence, setting $f(x) = (1 - \sqrt{x})^2$ yields the **Hellinger** divergence, and setting $f(x) = |x - 1|$ results in the **TotalVariation** distance.

Some of the measures we consider are not metrics as they do not satisfy at least one of the following properties: non-negativity, identity of indiscernibles⁸, symmetry and triangle inequality. KL divergence, for example, is not a metric since symmetry and triangle inequality do not hold. TotalVariation, on the other hand, is a metric satisfying all four properties.

Numerous symmetric versions of the KL divergence were proposed. These include **JensenShannon** [24] (defined using the mean language model of θ_q and θ_d), **J** [16] and **ResistorAverage** [17]. The J divergence is not a metric despite being symmetric since the triangle inequality does not hold. The square root of the Jensen-Shannon divergence is a metric.⁹

Three additional instances of the f-divergence are the asymmetric χ^2 **Neyman** [26] and χ^2 **Pearson** [27], where χ^2 *Pearson*(θ_q, θ_d) = χ^2 *Neyman*(θ_d, θ_q), and their symmetric version χ^2 **Symmetric**. The functions f in f-divergence that yield these three measures are $f(x) = \frac{(x-1)^2}{x}$, $f(x) = (1-x)^2$ and $f(x) = \frac{(x-1)^2}{x+1}$, respectively.

The **Skew** divergence [22, 23] was proposed to address cases where $KL(\theta_x, \theta_y)$ is not defined because the support of θ_x is not a subset of the support of θ_y . In our case, $KL(\theta_q, \theta_d)$ is defined since smoothed document language models are used.

We next make a few observations about the IDF effect in the query-document match components of the divergence measures. It can be shown that Hellinger, χ^2 Neyman, χ^2 Pearson and χ^2 Symmetric

exhibit the IDF effect. (More precisely, since we rank by ascending order of divergence values, the negative divergence measures exhibit, or not, the IDF effect.) In contrast, it is easy to verify that TotalVariation does not have an IDF effect.

For the negative J divergence, the negative $KL(\theta_q, \theta_d)$ employs an IDF effect since it is rank equivalent to $-CE(\theta_q, \theta_d)$ (see Eq. 1) and negative CE employs an IDF effect as shown in Eq. 3e. To estimate $-KL(\theta_d, \theta_q)$, we use an unsmoothed MLE for the document model and a query model smoothed via Jelinek-Mercer with the corpus (with parameter α_q), denoted θ_q^s . (See Section 4.1 for further details.) Using the fact that $KL(\theta_d^{MLE}, \theta_q^s) = -H(\theta_d^{MLE}) + CE(\theta_d^{MLE}, \theta_q^s)$, and applying manipulations as in Eq. 3, we get:

$$-KL(\theta_d^{MLE}, \theta_q^s) \stackrel{rank}{=} \sum_w p(w|\theta_d^{MLE}) \log\left(1 + \frac{1 - \alpha_q}{\alpha_q} \frac{p(w|\theta_q)}{p(w|\theta_d^{MLE})}\right) + \log \alpha_q - KL(\theta_d^{MLE}, \theta_D^{MLE});$$

the summation is over w that are in both d and θ_q 's support. High values of $p(w|\theta_d^{MLE})$ reduce the value of the document-query match (the summation), and hence there is an IDF effect. This effect is somewhat counter balanced by the fact that documents with models similar to that of the corpus, i.e., with low $KL(\theta_d^{MLE}, \theta_D^{MLE})$, are rewarded. Overall, since $-KL(\theta_q, \theta_d)$ and $-KL(\theta_d, \theta_q)$ have an IDF effect in their query-document match components, so does (negative) J divergence.

Using similar arguments to those used for the J divergence, we can show that ResistorAverage has an IDF effect. Along the same lines, the KL divergence factors that constitute JensenShannon and Skew have an IDF effect. However, in both measures, higher values of η make $\theta_{qd}^{[\eta]}$ more similar to θ_q . This causes loss of information about the original differences between θ_q and θ_d .

4 EMPIRICAL EXPLORATION

Our next goal is studying the retrieval effectiveness of using the different measures discussed above for comparing document and query language models.

4.1 Experimental Setup

The five datasets used for experiments are specified in Table 3. AP and ROBUST are small, mostly newswire, collections; GOV2 is a crawl of the .gov domain and CW09B is the Category B of the ClueWeb09 collection. An additional dataset, CW09BF, was created for ClueWeb09 Category B by filtering out from the initial document ranking documents assigned with a score below 50 by Waterloo's spam classifier [10]. Further details about the initial ranking are provided below. Unless stated otherwise, topic titles, which are short, served for queries. In Section 4.2.3 we also present evaluation when using verbose queries composed of both the title and description. We applied Krovetz stemming to documents and queries and removed stopwords on the INQUERY list only from queries. The Indri toolkit¹⁰ was used for experiments.

We used a two-phase retrieval approach to compare the effectiveness of the different measures discussed in Section 3. In the first phase, an initial list of documents is retrieved using $KL(\theta_q, \theta_d)$,

⁷Also known as the Ali-Silvey distance [2].

⁸For metric D and a pair of models θ_x and θ_y , $D(\theta_x, \theta_y) = 0$ iff $\theta_x = \theta_y$.

⁹The f functions in f-divergence that yield the J and JensenShannon divergence are $f(x) = (x-1) \log x$ and $f(x) = \frac{1}{2} \log \frac{2-x}{1+x} + \frac{1}{2} x \log \frac{2x}{1+x}$, respectively.

¹⁰www.lemurproject.org/indri

Table 3: TREC data used for experiments.

corpus	# documents	data	queries
AP	242,918	Disks 1 – 3	51 – 150
ROBUST	528,155	Disks 4 – 5 (-CR)	301 – 450, 600 – 700
GOV2	25,205,179	GOV2	701 – 850
CW09B	50,220,423	ClueWeb09 Category B	1 – 200
CW09BF			

where an unsmoothed unigram language model, specifically, a maximum likelihood estimate (MLE), is used to induce θ_q as is standard [20]. In the second phase, $M(\theta_d; \theta_q)$ or $D(\theta_q, \theta_d)$ is used to re-rank the 10000 most highly ranked documents in the initial ranking, where M and D are the weighted mean and divergence measures presented in Section 3, respectively. We take a re-ranking approach since applying numerous configurations of query expansion upon large-scale corpora is computationally expensive. For a uniform experimental setting, this approach is used in all cases even if query expansion is not employed. It is important to note that re-ranking an initial list of documents was shown to yield similar performance to ranking the entire corpus for pseudo-feedback query expansion which we use here, especially for precision-oriented evaluation metrics [12]. We further note that performance patterns similar to those reported here were also observed in experiments with re-ranking a shorter list of 1000 documents. (Actual numbers are omitted as they convey no additional insight.)

We use three query-model induction methods in the second phase. The first utilizes unsmoothed MLE induced from the original query as was the case in the first phase.¹¹ The additional two models are relevance model #3 (RM3 [1, 21]) and the mixture model (MM [31]). For these two models, θ_q is induced from the top 50 ranked documents in the initial list. The resultant query language models are (much) richer than—i.e., their support is much larger than that of—the one induced using MLE from a short query; these query models represent expanded queries of the short title query.

The document model θ_d is in both phases a Dirichlet-smoothed unigram language model with the smoothing parameter $\mu=1000$ [32]. Since the document model is smoothed, $KL(\theta_q, \theta_d)$ is always defined (i.e., the support of θ_q is a subset of the support of θ_d). In contrast, $KL(\theta_d, \theta_q)$, which is used in J and ResistorAverage, might be undefined if the support of θ_d is not a subset of the support of θ_q . We took the following two approaches to address this issue: (i) applied Jelinek-Mercer smoothing [32] with $\lambda = 0.1$ to θ_q , or (ii) considered only terms that appear in the support of both θ_q and θ_d . The document model is unsmoothed in both cases, i.e., θ_d^{MLE} is used.¹² For the J divergence, the first approach was found in our experiments to be more effective for short title queries and is therefore used in Section 4.2.1, while the second was found to be more effective for longer queries with a larger support and is used in Sections 4.2.2 and 4.2.3. For ResistorAverage, the first approach resulted in better performance in all cases. We note that the first approach entails an IDF effect for $-KL(\theta_d^{MLE}, \theta_q)$ (see Section 3.3) while the second does not.

¹¹Note that the document ranking produced using $KL(\theta_q, \theta_d)$ in the second phase is in this case the same as the initial ranking.

¹²This is also the case for $KL(\theta_d, \theta_q^{[\eta]})$ in JensenShannon. We note that JensenShannon is always well defined since θ_q is smoothed with θ_d in $\theta_q^{[\eta]}$.

To estimate retrieval effectiveness, we use mean average precision (MAP), precision of the top 5 ranked documents (p@5), normalized discounted cumulative gain of the top 20 ranked documents (NDCG20) and reliability of improvement (RI) [28]. RI is the difference between the number of queries whose AP performance is improved and the number of queries whose AP performance is hurt compared to the initial ranking divided by the total number of queries. Ten-fold cross validation is used to set free-parameter values where MAP serves for optimization in the training phase. Folds are determined based on query IDs. The reported performance is the average over all queries in a dataset when these serve for testing. (Each query belongs to a single test fold.) Statistically significant differences of retrieval performance between two methods are computed over all the queries in a dataset when these serve for testing. The two-tailed paired t-test with $p \leq 0.05$ is used.

To construct RM3 and MM, the number of terms and the weight of the original query are set to values in {25, 50} and {0.1, 0.2, . . . , 0.9}, respectively. Unsmoothed document language models (MLE) are used to induce RM3. The mixture weight of the corpus in MM is in {0.1, 0.3, 0.5, 0.7, 0.9}. The values of η (in Skew), β (in Power) and γ (in Lehmer) are selected from {0.1, 0.2, . . . , 0.9}, $\{\pm 0.05, \pm 0.15, \pm 0.25, \pm 0.5, \pm 1, \pm 2, \pm 3, \pm 10\}$, and $\{0, \pm 0.05, \pm 0.15, \pm 0.25, \pm 0.5, \pm 1, \pm 2, \pm 3, \pm 10\}$, respectively.

4.2 Experimental Results

4.2.1 Short Queries. Table 4 presents the results of using the unsmoothed MLE query language model θ_q^{MLE} —i.e., only the terms in the original short query are assigned with a non-zero probability. We see that in most cases Power is the best performing measure. Power outperforms KL in 12 out of the 15 relevant comparisons for the MAP, p@5 and NDCG20 evaluation metrics (5 datasets \times 3 evaluation metrics) and in 5 of these the improvement is statistically significant; Power is outperformed (but not statistically significantly) by KL in a single case (NDCG20 for ROBUST). Power is the only measure among those considered that is never statistically significantly outperformed by KL. It is also the only measure with positive RI values across all five datasets, indicating that the number of queries for which average precision performance is improved compared to the initial ranking (KL) is higher than the number of queries for which performance is hurt.

Two additional measures outperforming KL in the majority of the relevant comparisons are GeoHar (in 8 relevant comparisons; 4 are statistically significant) and Lehmer (in 9 relevant comparisons; 5 are statistically significant). All other measures are outperformed by KL in most relevant comparisons and many of the differences are statistically significant.

The descending order of the seven weighted means according to a pairwise comparison of retrieval effectiveness¹³ is Power, Lehmer, GeoHar, KL (Geo), Har, GeoAri and Ari. The poor performance of Ari could be explained by the fact that it has no IDF effect as shown in Section 3.2. Har is worse than Geo potentially due to being too conservative. GeoAri improves over Ari but not over Geo while GeoHar improves over both Geo and Har, which attests to the merits of integrating these two means. Power and Lehmer rely

¹³The order is determined by counting the number of relevant comparisons (out of the 15) in which a method outperforms another method.

Table 4: Short title queries. Using unsmoothed MLE to induce query models. Because the KL divergence is used to produce the initial ranking, its RI is undefined. ‘Init’: initial ranking. ‘k’: statistically significant difference of retrieval effectiveness (MAP, p@5 and NDCG20) with the KL divergence. The best result in a column is highlighted.

	AP				ROBUST				GOV2				CW09B				CW09BF			
	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI
KL (=Init)	21.1	43.6	42.6	–	25.4	48.7	43.9	–	29.2	55.5	44.5	–	17.9	22.7	20.1	–	18.7	34.7	25.8	–
Ari	10.5 _k	24.6 _k	24.4 _k	–69.7	6.4 _k	14.5 _k	12.7 _k	–83.5	7.0 _k	9.6 _k	9.1 _k	–93.2	6.6 _k	5.8 _k	6.0 _k	–66.7	7.3 _k	13.2 _k	9.7 _k	–67.2
Har	15.5 _k	39.4 _k	35.8 _k	–55.6	22.7 _k	47.1	41.3 _k	–36.5	26.5 _k	57.0	44.4	–24.3	19.5 _k	28.6 _k	23.8 _k	27.3	19.1	35.5	26.7	8.6
GeoAri	18.0 _k	41.4	39.1 _k	–35.4	16.0 _k	37.8 _k	32.4 _k	–70.7	17.1 _k	34.5 _k	27.4 _k	–87.8	10.0 _k	8.6 _k	9.2 _k	–66.7	12.4 _k	23.9 _k	16.9 _k	–63.1
GeoHar	17.2 _k	41.0	37.7 _k	–47.5	24.0 _k	47.7	42.1 _k	–24.1	28.0 _k	57.2	45.3	–8.1	19.6 _k	27.6 _k	23.3 _k	35.4	19.4	35.9	26.8 _k	19.7
Power	21.2	45.1	42.6	19.2	25.5	48.9	43.8	5.2	29.2	55.8	44.9 _k	23.0	19.8 _k	27.9 _k	23.6 _k	37.4	19.2	35.6	26.7 _k	19.7
Lehmer	20.8	44.2	41.6 _k	–9.1	25.0 _k	48.8	43.7	–7.6	29.0	55.4	44.7	10.8	19.5 _k	28.0 _k	23.4 _k	31.8	19.4	36.0	26.8 _k	20.7
Hellinger	17.7 _k	41.8	39.3 _k	–30.3	14.2 _k	34.1 _k	28.9 _k	–73.5	14.3 _k	29.7 _k	23.1 _k	–89.2	9.9 _k	10.8 _k	10.0 _k	–66.7	11.9 _k	23.6 _k	16.3 _k	–65.2
TotalVariation	10.5 _k	24.6 _k	24.4 _k	–70.7	6.4 _k	14.5 _k	12.7 _k	–83.1	7.0 _k	9.6 _k	9.1 _k	–92.6	6.6 _k	5.8 _k	6.0 _k	–66.2	7.3 _k	13.2 _k	9.7 _k	–67.2
JensenShannon	10.6 _k	32.3 _k	28.1 _k	–74.7	8.3 _k	22.9 _k	18.9 _k	–80.3	4.3 _k	10.0 _k	8.3 _k	–94.6	4.4 _k	10.1 _k	7.3 _k	–88.4	3.9 _k	13.4 _k	8.3 _k	–85.9
J	19.0 _k	41.0	40.7 _k	–40.4	23.0 _k	47.9	42.0 _k	–30.9	28.5	60.4 _k	46.3	–12.2	17.9	23.9	21.6	5.6	19.4	39.1 _k	28.4 _k	2.5
ResistorAverage	14.7 _k	40.0	36.6 _k	–66.7	15.5 _k	40.2 _k	33.7 _k	–74.3	23.1 _k	56.4	42.6	–51.4	16.8	21.5	19.6	4.0	18.2	37.2	26.7	–4.5
χ^2 Neyman	15.5 _k	39.4 _k	35.8 _k	–55.6	22.7 _k	47.1	41.3 _k	–36.5	26.5 _k	56.9	44.4	–24.3	19.5 _k	28.6 _k	23.8 _k	27.3	19.1	35.5	26.7	8.6
χ^2 Pearson	10.6 _k	25.3 _k	24.6 _k	–69.7	6.5 _k	14.9 _k	13.0 _k	–83.9	7.1 _k	10.1 _k	9.4 _k	–92.6	6.6 _k	6.1 _k	6.2 _k	–66.2	7.4 _k	13.5 _k	9.8 _k	–66.7
χ^2 Symmetric	10.7 _k	26.7 _k	25.4 _k	–67.7	6.6 _k	15.3 _k	13.3 _k	–82.7	7.3 _k	11.1 _k	10.2 _k	–92.6	6.8 _k	6.5 _k	6.4 _k	–66.7	7.5 _k	13.9 _k	10.1 _k	–67.2
Skew	9.1 _k	22.4 _k	21.3 _k	–84.8	6.4 _k	15.7 _k	13.6 _k	–87.6	3.1 _k	6.2 _k	4.3 _k	–95.9	4.1 _k	8.8 _k	6.5 _k	–86.4	3.8 _k	10.7 _k	7.6 _k	–89.4
Cosine	7.5 _k	21.8 _k	19.4 _k	–83.8	10.1 _k	19.2 _k	18.1 _k	–89.6	4.1 _k	13.4 _k	8.7 _k	–95.9	1.5 _k	5.6 _k	3.5 _k	–88.9	1.4 _k	5.8 _k	3.1 _k	–88.9

on a free parameter in contrast to the other means – some of which they generalize as noted in Section 3.2. We further study the effect of β on the retrieval performance of Power in Section 4.2.5.

The descending order of the ten divergence measures according to a pairwise performance comparison is KL, χ^2 Neyman, J, ResistorAverage, Hellinger, χ^2 Symmetric, JensenShannon, χ^2 Pearson, TotalVariation and Skew. TotalVariation is low ranked potentially because it does not have an IDF effect. Skew is low ranked (and to some extent also JensenShannon) potentially due to smoothing the document model with the query model. KL, χ^2 Neyman, J and ResistorAverage are presumably among the top 4 as they all have an IDF effect.

4.2.2 Pseudo-Feedback-Based Query Models. Table 5 presents the results of using RM3 and MM as the query language models. (The initial ranking was attained, as described in Section 4.1, using an unsmoothed MLE induced from the title queries.) These query language models are much richer than the unsmoothed MLE query models explored in Section 4.2.1; i.e., their support is much larger. Two measures that stand out are the J divergence and Power. The former outperforms KL in 8 relevant comparisons for RM3 and in 13 for MM, while the latter outperforms KL in 12 relevant comparisons for RM3 and 10 for MM; some of these improvements are statistically significant. Both measures are statistically significantly outperformed by KL in at most two cases for RM3 and never for MM. The RI of both measures is always positive, whereas for the KL divergence it is positive in all but a single case: MM for CW09B.

Apart from Geo (KL), which is promoted to the second position after Power, the performance order of the weighted means remains the same as in Section 4.2.1. The top 4 divergence measures in descending order of pairwise comparisons for RM3 are J, KL, ResistorAverage and χ^2 Neyman. For MM, χ^2 Neyman and ResistorAverage switch places. These measures were also the top 4 in Section 4.2.1. The bottom 4 measures for both RM3 and MM are χ^2 Pearson, JensenShannon, Skew and TotalVariation. These measures also appeared at the bottom of the list in Section 4.2.1.

4.2.3 Verbose Queries. To further study the effect of the query-model support size on retrieval performance, we repeated the experiment from Section 4.2.1 but now using verbose queries composed of the topic’s title and description. As a result, the support of the query model is (much) larger than that used in Section 4.2.1 where MLE of short title queries is used. The results are presented in Table 6. As was the case thus far, the two measures outperforming the KL divergence are the J divergence and Power. However, unlike the case in Sections 4.2.1 and 4.2.2, both measures are outperformed by KL in at least the same number of relevant comparisons in which they outperform it. Most performance differences between these measures and KL are not statistically significant. The performance-based ordering of the weighted means is as in Section 4.2.1, except that now Geo (KL) is promoted to the first rank. The top 4 divergence measures are KL, J, ResistorAverage and Hellinger, while the bottom 4 are JensenShannon, Skew, χ^2 Pearson and TotalVariation as was the case in Sections 4.2.1 and 4.2.2.

4.2.4 Power vs. J divergence. We saw that Power is the most effective weighted mean while J is the most effective divergence measure among the alternatives we considered. Table 7 contrasts their performance for: MLE query models induced from short (title) or verbose (title+description) queries, and pseudo-feedback-based query models (RM3 or MM). Evidently, Power is more effective than J for short title queries and for RM3, while the reverse holds for verbose queries and MM.

4.2.5 Further Analysis. We showed that Power often outperforms KL for query models induced from short titles using unsmoothed MLE and those induced using pseudo feedback. In Figure 1 we study the effect of β on the MAP performance of Power. Short title queries and the MLE query model are used. We note that similar trends were also observed for RM3, MM and verbose queries. These results are omitted due to space considerations and as they convey no further insight. We see that the performance of Power is the highest for values of β close to 0. (Recall from Section 3.2 that Power converges to Geo as β approaches zero and that it has an IDF effect for $\beta < 1$.) A case in point, we found that

Table 5: Using RM3 and MM to induce query models. ‘Init’: initial ranking, ‘i’ and ‘k’ mark statistically significant differences of retrieval effectiveness (MAP, p@5 and NDCG20) with Init and KL divergence, respectively. The best result in a column per query model is highlighted.

	AP				ROBUST				GOV2				CW09B				CW09BF				
	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	
Init	21.1	43.6	42.6	–	25.4	48.7	43.9	–	29.2	55.5	44.5	–	17.9	22.7	20.1	–	18.7	34.7	25.8	–	
RM3	KL	27.1 ⁱ	48.7 ⁱ	46.5 ⁱ	43.4	28.8 ⁱ	48.8	44.6	23.7	32.8 ⁱ	58.0	47.3 ⁱ	43.9	19.0 ⁱ	26.1 ⁱ	22.5 ⁱ	11.6	19.9 ⁱ	37.2	27.1 ⁱ	13.6
	Ari	11.0 ⁱ	27.5 ⁱ	26.1 ⁱ	–64.6	6.4 ⁱ	14.6 ⁱ	13.1 ⁱ	–84.3	8.2 ⁱ	15.7 ⁱ	13.7 ⁱ	–91.9	6.6 ⁱ	5.7 ⁱ	6.1 ⁱ	–63.1	7.9 ⁱ	15.3 ⁱ	11.3 ⁱ	–63.1
	Har	17.8 ⁱ	40.8 ⁱ	40.0 ⁱ	–41.4	22.3 ⁱ	47.2	40.6 ⁱ	–37.3	23.8 ⁱ	53.5 ⁱ	42.3 ⁱ	–47.3	14.0 ⁱ	28.0 ⁱ	19.8 ⁱ	–28.8	14.9 ⁱ	32.2 ⁱ	22.3 ⁱ	–43.4
	GeoAri	24.1 ⁱ	45.1 ⁱ	45.0 ⁱ	20.2	24.0 ⁱ	45.1 ⁱ	41.1 ⁱ	–24.9	25.7 ⁱ	50.3 ⁱ	40.6 ⁱ	–35.1	12.3 ⁱ	14.4 ⁱ	13.6 ⁱ	–48.5	14.8 ⁱ	30.8 ⁱ	20.5 ⁱ	–42.4
	GeoHar	18.7 ⁱ	41.0 ⁱ	41.0 ⁱ	–26.3	23.9 ⁱ	47.6	42.3 ⁱ	–28.1	26.0 ⁱ	51.1 ⁱ	43.6 ⁱ	–32.4	15.7 ⁱ	29.0 ⁱ	20.9	–21.2	16.7 ⁱ	34.8	24.2 ⁱ	–25.3
	Power	26.8	48.9 ⁱ	46.5 ⁱ	42.4	28.9 ⁱ	49.3	44.7	24.1	33.2 ⁱ	58.8	47.5 ⁱ	39.9	19.4 ⁱ	28.3 ⁱ	23.3 ⁱ	21.2	19.9 ⁱ	37.4	27.2 ⁱ	6.6
	Lehmer	24.4 ⁱ	46.5	46.0 ⁱ	32.3	26.4 ⁱ	49.0	43.9	2.4	30.2 ⁱ	57.2	46.0	10.1	17.9 ⁱ	22.9 ⁱ	20.8 ⁱ	–5.1	19.0 ⁱ	35.8	26.1 ⁱ	–3.0
	Hellinger	25.1 ⁱ	47.1	46.0 ⁱ	22.2	22.8 ⁱ	46.1 ⁱ	40.8 ⁱ	–24.9	23.8 ⁱ	54.1 ⁱ	42.0 ⁱ	–39.2	15.5 ⁱ	27.0 ⁱ	20.1 ⁱ	–35.9	16.7 ⁱ	36.2	24.5 ⁱ	–28.8
	TotalVariation	13.8 ⁱ	43.2 ⁱ	36.6 ⁱ	–60.6	11.3 ⁱ	33.8 ⁱ	26.3 ⁱ	–82.7	11.7 ⁱ	41.1 ⁱ	30.3 ⁱ	–86.5	7.7 ⁱ	19.2 ⁱ	13.2 ⁱ	–78.8	8.8 ⁱ	26.3 ⁱ	16.7 ⁱ	–76.3
	JensenShannon	18.3 ⁱ	43.6	40.2 ⁱ	–34.3	13.9 ⁱ	36.5 ⁱ	30.6 ⁱ	–75.9	12.7 ⁱ	41.8 ⁱ	31.0 ⁱ	–85.1	9.8 ⁱ	23.5	16.1 ⁱ	–71.7	10.1 ⁱ	26.8 ⁱ	17.6 ⁱ	–69.7
	J	26.3 ⁱ	47.7	46.5 ⁱ	42.4	28.1 ⁱ	50.4 ⁱ	45.2 ⁱ	29.7	33.2 ⁱ	60.7 ⁱ	48.6 ⁱ	40.5	18.7	27.0 ⁱ	21.8 ⁱ	13.6	19.9 ⁱ	39.9 ⁱ	27.4 ⁱ	8.1
	ResistorAverage	23.5 ⁱ	47.5	45.1	10.1	22.0 ⁱ	44.8 ⁱ	40.0 ⁱ	–41.4	26.6 ⁱ	58.0	44.5 ⁱ	–29.7	18.5	28.3 ⁱ	21.9 ⁱ	9.1	19.1	40.0 ⁱ	26.7	6.1
	χ^2 Neyman	19.3 ⁱ	41.8 ⁱ	40.9 ⁱ	–16.2	23.5 ⁱ	47.6	41.8 ⁱ	–30.5	26.0 ⁱ	56.1	44.1 ⁱ	–35.1	18.7	29.0 ⁱ	23.2 ⁱ	13.6	18.6 ⁱ	35.5	26.4	–7.6
	χ^2 Pearson	19.4 ⁱ	44.4 ⁱ	42.1 ⁱ	–14.1	15.9 ⁱ	39.4 ⁱ	33.2 ⁱ	–64.3	15.5 ⁱ	46.2 ⁱ	34.3 ⁱ	–74.3	9.7 ⁱ	21.9 ⁱ	15.7 ⁱ	–72.7	11.1 ⁱ	30.6 ⁱ	19.6 ⁱ	–69.7
χ^2 Symmetric	21.2 ⁱ	45.5	43.5 ⁱ	0.0	16.8 ⁱ	40.6 ⁱ	33.9 ⁱ	–60.6	17.4 ⁱ	47.4 ⁱ	35.4 ⁱ	–67.6	11.7 ⁱ	23.6	17.5 ⁱ	–58.1	12.8 ⁱ	32.9 ⁱ	21.2 ⁱ	–60.6	
Skew	17.4 ⁱ	42.4 ⁱ	39.3 ⁱ	–42.4	14.1 ⁱ	35.3 ⁱ	30.0 ⁱ	–75.5	12.3 ⁱ	41.1 ⁱ	29.8 ⁱ	–87.8	9.1 ⁱ	21.4 ⁱ	14.3 ⁱ	–76.3	9.7 ⁱ	25.2 ⁱ	16.8 ⁱ	–72.2	
MM	KL	27.8 ⁱ	50.1 ⁱ	46.9 ⁱ	–22.2	27.1 ⁱ	47.5	43.2	24.1	32.1 ⁱ	58.0	45.4	31.1	18.9	24.8	21.8	–5.1	20.4 ⁱ	38.7 ⁱ	27.4	10.6
	Ari	18.7 ⁱ	38.0 ⁱ	36.6 ⁱ	–24.2	10.8 ⁱ	22.0 ⁱ	19.9 ⁱ	–68.7	10.7 ⁱ	21.1 ⁱ	16.5 ⁱ	–89.2	6.8 ⁱ	5.4 ⁱ	6.2 ⁱ	–65.7	8.0 ⁱ	14.7 ⁱ	11.9 ⁱ	–76.3
	Har	19.4 ⁱ	44.8 ⁱ	41.0 ⁱ	–23.2	22.5 ⁱ	46.4 ⁱ	40.5 ⁱ	–32.5	24.9 ⁱ	54.7	42.1 ⁱ	–44.6	13.5 ⁱ	26.0	18.1 ⁱ	–42.9	16.0 ⁱ	34.4 ⁱ	24.0 ⁱ	–33.3
	GeoAri	25.2 ⁱ	49.1 ⁱ	45.5	16.2	22.8 ⁱ	42.7 ⁱ	39.3 ⁱ	–28.5	24.8 ⁱ	47.4 ⁱ	39.0 ⁱ	–33.8	12.4 ⁱ	13.0 ⁱ	13.7 ⁱ	–53.5	15.4 ⁱ	31.7 ⁱ	21.9 ⁱ	–36.9
	GeoHar	22.1 ⁱ	43.8 ⁱ	41.4 ⁱ	–11.1	23.6 ⁱ	47.2	41.8 ⁱ	–30.1	26.9 ⁱ	56.6	43.7 ⁱ	–20.9	15.6 ⁱ	27.3 ⁱ	19.9	–27.3	17.6 ⁱ	36.6	25.4 ⁱ	–18.7
	Power	27.9 ⁱ	49.3 ⁱ	46.7 ⁱ	20.2	27.3 ⁱ	47.6	43.5	24.9	32.1 ⁱ	58.1	45.6	27.0	20.1 ⁱ	29.9 ⁱ	24.1 ⁱ	27.8	20.3 ⁱ	38.6 ⁱ	27.7 ⁱ	17.2
	Lehmer	27.6 ⁱ	49.9 ⁱ	46.9 ⁱ	22.2	26.0 ⁱ	45.8 ⁱ	42.5	8.8	30.1 ⁱ	57.3	46.4	21.6	17.8 ⁱ	23.3	20.3 ⁱ	–21.2	19.5 ⁱ	38.8 ⁱ	26.9	–6.6
	Hellinger	26.4 ⁱ	47.3	46.4	22.2	20.3 ⁱ	40.5 ⁱ	36.2 ⁱ	–24.1	23.2 ⁱ	50.7 ⁱ	37.8 ⁱ	–40.5	15.6 ⁱ	24.0	19.4 ⁱ	–40.9	17.0 ⁱ	37.4	24.2 ⁱ	–34.3
	TotalVariation	22.1 ⁱ	45.1 ⁱ	43.2 ⁱ	0.0	14.2 ⁱ	33.3 ⁱ	28.1 ⁱ	–54.2	15.2 ⁱ	38.4 ⁱ	28.0 ⁱ	–79.7	7.5 ⁱ	12.0 ⁱ	11.0 ⁱ	–78.8	10.0 ⁱ	24.5 ⁱ	17.3 ⁱ	–69.2
	JensenShannon	24.5 ⁱ	46.7	44.5	14.1	16.9 ⁱ	34.7 ⁱ	31.3 ⁱ	–47.8	13.2 ⁱ	32.4 ⁱ	25.0 ⁱ	–86.5	10.2 ⁱ	21.4	14.9 ⁱ	–74.7	10.7 ⁱ	25.6 ⁱ	16.8 ⁱ	–72.2
	J	28.2 ⁱ	49.7 ⁱ	47.5 ⁱ	21.2	27.3 ⁱ	48.5 ⁱ	43.8 ⁱ	24.9	32.2 ⁱ	58.2	46.5 ⁱ	39.2	19.2 ⁱ	26.1 ⁱ	21.8 ⁱ	12.6	20.8 ⁱ	41.8 ⁱ	29.0 ⁱ	27.8
	ResistorAverage	25.1 ⁱ	47.3	45.3	8.1	20.7 ⁱ	42.7 ⁱ	38.7 ⁱ	–32.5	26.9 ⁱ	54.7	42.5 ⁱ	–20.3	18.4	23.3	21.0	8.1	19.6	38.8	26.6	–1.0
	χ^2 Neyman	21.6 ⁱ	46.5	43.4 ⁱ	–6.1	22.8 ⁱ	44.7 ⁱ	40.3 ⁱ	–31.3	26.3 ⁱ	56.2	44.3	–27.0	18.5	29.5 ⁱ	23.0 ⁱ	4.0	18.7 ⁱ	36.9	26.1	–8.1
	χ^2 Pearson	22.9 ⁱ	46.3 ⁱ	43.5 ⁱ	4.0	14.6 ⁱ	33.7 ⁱ	28.5 ⁱ	–51.0	15.8 ⁱ	40.7 ⁱ	29.2 ⁱ	–71.6	10.5 ⁱ	22.9	16.7 ⁱ	–64.6	11.6 ⁱ	33.4 ⁱ	20.6 ⁱ	–58.1
χ^2 Symmetric	24.0 ⁱ	47.5	45.0	8.1	15.2 ⁱ	33.6 ⁱ	29.1 ⁱ	–45.4	16.6 ⁱ	43.0 ⁱ	31.1 ⁱ	–70.3	12.3 ⁱ	23.5	17.8 ⁱ	–55.6	13.1 ⁱ	34.4 ⁱ	21.7 ⁱ	–48.0	
Skew	23.8 ⁱ	46.7	43.7 ⁱ	14.1	16.8 ⁱ	32.6 ⁱ	29.1 ⁱ	–47.8	14.2 ⁱ	33.5 ⁱ	25.6 ⁱ	–86.5	9.5 ⁱ	18.2 ⁱ	13.3 ⁱ	–78.8	11.0 ⁱ	26.0 ⁱ	17.3 ⁱ	–73.7	

Table 6: Verbose (title+description) queries. Using unsmoothed MLE to induce query models. Because the KL divergence is used to produce the initial ranking, its RI is undefined. ‘Init’: initial ranking, ‘k’: statistically significant difference of retrieval effectiveness (MAP, p@5 and NDCG20) with the KL divergence. The best result in a column is highlighted.

	AP				ROBUST				GOV2				CW09B				CW09BF			
	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI	MAP	p@5	NDCG20	RI
KL (=Init)	23.8	49.1	46.0	–	28.4	54.8	47.3	–	30.0	61.4	47.1	–	18.1	29.2	24.2	–	18.1	36.1	26.6	–
Ari	9.8 _k	29.3 _k	26.6 _k	–77.8	5.6 _k	13.3 _k	12.5 _k	–93.6	7.2 _k	10.9 _k	9.8 _k	–91.9	4.8 _k	3.1 _k	4.2 _k	–63.1	6.2 _k	12.0 _k	8.7 _k	–50.0
Har	10.2 _k	25.5 _k	24.2 _k	–87.9	17.6 _k	36.9 _k	32.0 _k	–82.7	15.5 _k	45.5 _k	31.9 _k	–78.4	9.2 _k	19.3 _k	14.4 _k	–66.7	8.8 _k	20.4 _k	14.1 _k	–74.2
GeoAri	17.3 _k	43.0 _k	40.4 _k	–57.6	14.8 _k	35.9 _k	29.7 _k	–82.3	16.5 _k	33.5 _k	27.3 _k	–83.8	8.6 _k	7.3 _k	8.6 _k	–51.0	12.0 _k	25.2 _k	17.7 _k	–30.8
GeoHar	15.3 _k	33.5 _k	32.2 _k	–75.8	22.1 _k	42.2 _k	37.9 _k	–69.9	20.0 _k	51.1 _k	36.7 _k	–71.6	12.0 _k	24.1 _k	17.9 _k	–59.6	11.4 _k	25.1 _k	17.9 _k	–66.2
Power	23.8	48.7	46.7 _k	9.1	28.0 _k	53.9	47.2	–6.0	29.9	60.1	47.0	11.5	17.2 _k	29.1	23.3 _k	–35.4	17.8	37.0	26.5	20.7
Lehmer	23.2 _k	48.1	45.4	–13.1	27.0 _k	51.4 _k	45.2 _k	–34.1	27.3 _k	57.8 _k	44.6 _k	–49.3	15.1 _k	24.2 _k	20.5 _k	–51.0	15.8 _k	32.5 _k	24.0 _k	–44.4
Hellinger	18.1 _k	45.1	42.5 _k	–43.4	16.6 _k	39.4 _k	32.7 _k	–71.1	16.7 _k	33.9 _k	28.4 _k	–85.1	11.5 _k	18.7 _k	15.0 _k	–42.9	13.1 _k	30.6 _k	20.8 _k	–30.8
TotalVariation	7.9 _k	25.1 _k	23.3 _k	–73.7	5.0 _k	12.9 _k	11.5 _k	–92.8	6.5 _k											

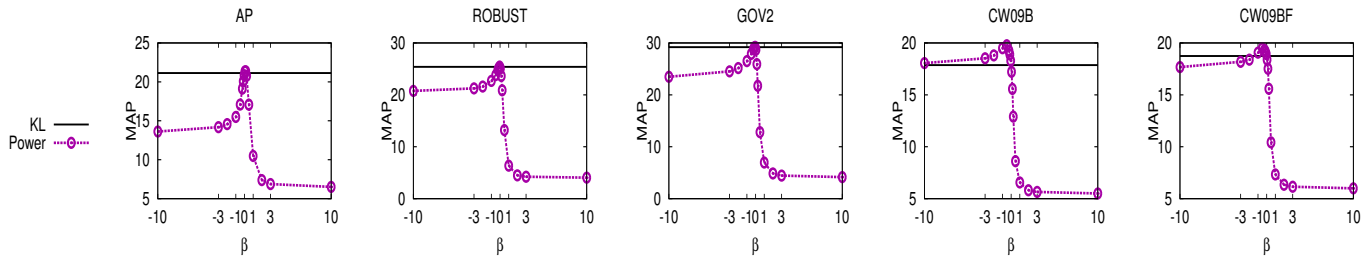


Figure 1: The effect of β on the MAP performance of Power. Short title queries and the MLE query model are used.

Table 7: Power vs. J divergence. The number of relevant comparisons out of 15 (5 datasets \times 3 evaluation measures: MAP, p@5 and NDCG20) in which the method in the row (statistically significantly) outperforms the method in the column.

	Short Query MLE			RM3			MM			Verbose Query MLE		
	KL	Power	J	KL	Power	J	KL	Power	J	KL	Power	J
KL	0 (0)	1 (0)	7 (4)	0 (0)	1 (1)	5 (2)	0 (0)	4 (0)	1 (0)	0 (0)	12 (3)	5 (2)
Power	12 (5)	0 (0)	10 (6)	12 (3)	0 (0)	6 (4)	10 (3)	0 (0)	3 (3)	2 (1)	0 (0)	5 (1)
J	7 (3)	5 (1)	0 (0)	8 (5)	6 (2)	0 (0)	13 (6)	11 (4)	0 (0)	5 (1)	8 (3)	0 (0)

5 CONCLUSIONS

Motivated by the fact that comparing query and document language models using the KL divergence is rank equivalent to using a specific weighted geometric mean, we studied alternative weighted means as well as divergence measures; specifically, we analyzed the inverse document frequency (IDF) effect of the methods. Empirical evaluation showed that KL can be often outperformed in several settings by some alternatives.

Acknowledgments We thank the reviewers for their helpful comments. This paper is based upon work supported in part by the Israel Science Foundation under grant no. 433/12.

REFERENCES

- [1] Nasreen Abdal-Jaleel, James Allan, W. Bruce Croft, Fernando Diaz, Leah Larkey, Xiaoyan Li, Marck D., and Courtney Wade. 2004. UMMASS at TREC 2004 – Novelty and HARD. In *Proc. of TREC-13*.
- [2] S. M. Ali and S. D. Silvey. 1966. A General Class of Coefficients of Divergence of One Distribution from Another. *Journal of the Royal Statistical Society. Series B (Methodological)* 28, 1 (1966), 131–142.
- [3] Javed A. Aslam and Virgiliu Pavlu. 2007. Query Hardness Estimation Using Jensen-Shannon Divergence Among Multiple Scoring Functions. In *Proc. of ECIR*. 198–209.
- [4] Gleb Beliakov, Humberto Bustince Sola, and Tomasa Calvo. 2016. *A Practical Guide to Averaging Functions*. Studies in Fuzziness and Soft Computing, Vol. 329. Springer.
- [5] Thorsten Brants, Francine Chen, and Ioannis Tsochantaridis. 2002. Topic-based document segmentation with probabilistic latent semantic analysis. In *Proc. of CIKM*. 211–218.
- [6] B.S. Bullen. 2003. *Handbook of means and their inequalities*. Springer-Science+Business Media, B.V.
- [7] Francine Chen, Ayman Farahat, and Thorsten Brants. 2004. Multiple Similarity Measures and Source-Pair Information in Story Link Detection. In *Proc. of HLT-NAACL*. 313–320.
- [8] Ruey-Cheng Chen, Chia-Jung Lee, and W. Bruce Croft. 2015. On Divergence Measures and Static Index Pruning. In *Proc. of ICTIR*. 151–160.
- [9] Stéphane Clinchant and Éric Gaussier. 2013. A Theoretical Analysis of Pseudo-Relevance Feedback Models. In *Proc. of ICTIR*. 6.

- [10] Gordon V. Cormack, Mark D. Smucker, and Charles L. A. Clarke. 2011. Efficient and effective spam filtering and re-ranking for large web datasets. *Information Retrieval* 14, 5 (2011), 441–465.
- [11] Imre Csiszar. 1967. Information-type measures of difference of probability distributions and indirect observations. *Studia Sci. Math. Hungar.* 2 (1967), 299–318.
- [12] Fernando Diaz. 2015. Condensed List Relevance Models. In *Proc. of ICTIR*. 313–316.
- [13] Hui Fang, Tao Tao, and ChengXiang Zhai. 2004. A formal study of information retrieval heuristics. In *Proc. of SIGIR*. 49–56.
- [14] Hui Fang and ChengXiang Zhai. 2005. An exploration of axiomatic approaches to information retrieval. In *Proc. of SIGIR*. 480–487.
- [15] Hussein Hazimeh and ChengXiang Zhai. 2015. Axiomatic Analysis of Smoothing Methods in Language Models for Pseudo-Relevance Feedback. In *Proc. of ICTIR*. 141–150.
- [16] Harold Jeffreys. 1939. *Theory of Probability*. Oxford University Press.
- [17] Don Johnson and Sinan Sinanovic. 2001. Symmetrizing the Kullback-Leibler Distance. *IEEE Transactions on Information Theory* (2001).
- [18] Maryam Karimzadehgan and ChengXiang Zhai. 2012. Axiomatic Analysis of Translation Language Model for Information Retrieval. In *Proc. of ECIR*. 268–280.
- [19] Oren Kurland and Lillian Lee. 2010. PageRank without hyperlinks: Structural reranking using links induced by language models. *ACM Transactions on information systems* 28, 4 (2010), 18.
- [20] John D. Lafferty and Chengxiang Zhai. 2001. Document language models, query models, and risk minimization for information retrieval. In *Proc. of SIGIR*. 111–119.
- [21] Victor Lavrenko and W. Bruce Croft. Relevance Models in Information Retrieval. In *Language Modeling for Information Retrieval*. Number 13 in Information Retrieval Book Series. Kluwer, Chapter 2, 11–56.
- [22] Lillian Lee. 1999. Measures of Distributional Similarity. In *Proc. of ACL*.
- [23] Lillian Lee. 2001. On the effectiveness of the skew divergence for statistical language analysis. In *Proc. of AISTATS*. 65–72.
- [24] Jianhua Lin. 1991. Divergence measures based on the Shannon entropy. *IEEE Trans. Information Theory* 37, 1 (1991), 145–151.
- [25] Ramesh Nallapati. 2006. *The smoothed dirichlet distribution: Understanding crossentropy ranking in information retrieval*. Ph.D. Dissertation. University of Massachusetts.
- [26] Jerzy Neyman. 1949. Contribution to the theory of the χ^2 test. In *Proc. of the Berkeley symposium on mathematical statistics and probability*, Vol. 1. University of California Press Berkeley, 239–273.
- [27] Karl Pearson. 1900. X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philosophical Magazine Series 5* 50, 302 (1900), 157–175.
- [28] Tetsuya Sakai, Toshihiko Manabe, and Makoto Koyama. 2005. Flexible pseudo-relevance feedback via selective sampling. *ACM Transactions on Asian Language Information Processing* 4, 2 (2005), 111–135.
- [29] Fei Song and W. Bruce Croft. 1999. A general language model for information retrieval. In *Proc. of SIGIR*. 279–280.
- [30] ChengXiang Zhai. 2008. Statistical Language Models for Information Retrieval: A Critical Review. *Foundations and Trends in Information Retrieval* 2, 3 (2008), 137–213.
- [31] Chengxiang Zhai and John D. Lafferty. 2001. Model-based Feedback in the Language Modeling Approach to Information Retrieval. In *Proc. of CIKM*. 403–410.
- [32] Chengxiang Zhai and John D. Lafferty. 2001. A Study of Smoothing Methods for Language Models Applied to Ad Hoc Information Retrieval. In *Proc. of SIGIR*. 334–342.
- [33] Justin Zobel and Alistair Moffat. 1998. Exploring the similarity space. *ACM SIGIR forum* 18, 1 (1998), 18–34.