

Part III

Predicting performance in recommender systems

*You can never get a cup of tea large enough or
a book long enough to suit me.*

Clive Staples Lewis

Chapter 5

Performance prediction in Information Retrieval

Information retrieval performance prediction has been mostly addressed as a query performance issue, which refers to the performance of an information retrieval system in response to a specific query. It also relates to the appropriateness of a query as an expression of the user's information needs. In general, performance prediction methods have been classified into two categories depending on the used data: pre-retrieval approaches, which make the prediction before the retrieval stage using query features, and post-retrieval approaches, which use the rankings produced by a retrieval engine. In particular, the so-called clarity score predictor – of special interest for this thesis – has been defined in terms of language models, and captures the ambiguity of a query with respect to the utilised document collection, or a specific result set.

In this chapter we provide an overview of terminology, techniques, and evaluation related to performance prediction in Information Retrieval. In Section 5.1 we introduce terminology and fundamental concepts of the performance prediction problem. In Section 5.2 we describe the different types of performance prediction approaches, which are mainly classified in the two categories mentioned above: pre-retrieval and post-retrieval approaches. Then, in Section 5.3 we provide a thorough analysis on the use of clarity score as a performance prediction technique, including examples, adaptations, and applications found in the literature. Finally, in Section 5.4 we introduce the general methodology used to evaluate performance predictors, along with the most common methods to measure their quality.

5.1 Introduction

Performance prediction has received little attention, if any, to date in the Recommender Systems field. Our research, however, finds a close and highly relevant reference in the adjacent Information Retrieval discipline, where performance prediction has gained increasing attention since the late 90's, and has become an established research topic in the field. Performance prediction finds additional motivation in personalised recommendation, inasmuch the applications they are integrated in may decide to produce recommendations or hold them back, delivering only the sufficiently reliable ones. Moreover, the ability to predict the effectiveness of individual algorithms can be envisioned as a strategy to optimise the combination of algorithms into ensemble recommenders, which currently dominate the field – rarely if ever are individual algorithms used alone in working applications, neither are they found individually in the top ranks of evaluation campaigns and competitions (Bennett and Lanning, 2007).

In Information Retrieval performance prediction has been mostly addressed as a query performance problem (Cronen-Townsend et al., 2002). Query performance refers to the performance of an information retrieval system in response to a particular query. It also relates to the appropriateness of a query as an expression of a user's information needs. Dealing effectively with poorly-performing queries is a crucial issue in Information Retrieval since it could improve the retrieval effectiveness significantly (Carmel and Yom-Tov, 2010).

In general, performance prediction techniques can be useful from different perspectives (Zhou and Croft, 2006; Yom-Tov et al., 2005a):

- From the user's perspective, it provides valuable feedback that can be used to direct a search, e.g. by rephrasing the query or suggesting alternative terms.
- From the system's perspective, it provides a means to address the problem of information retrieval consistency. The consistency of retrieval systems can be addressed by distinguishing poorly performing queries. A retrieval system may invoke different retrieval strategies depending on the query, e.g. by using query expansion or ranking functions based on the predicted difficulty of the query.
- From the system administrator's perspective, it may let identify queries related to a specific subject that are difficult for the search engine. According to such queries, the collection of documents could be extended to better answer insufficiently covered topics.
- From a distributed information retrieval's perspective, it can be used to decide which search engine (and/or database) to use, or how much weight give to different search engines when their results are combined.

Specifically, the performance prediction task in Information Retrieval is formalised based on the following three core concepts: **performance predictor**, **retrieval quality assessment**, and **predictor quality assessment**. In this context, the performance predictor is defined as a function that receives the query (and the result list D_q retrieved by the system, the set of relevant documents R_q , collection statistics \mathcal{C} , etc.), and returns a prediction of the retrieval quality for that query. Then, by means of a predictor quality assessment method, the predictive power of the performance predictor is estimated.

Based on the notation given in (Carmel and Yom-Tov, 2010), the problem of performance prediction consists of estimating a true retrieval quality metric $\mu(q)$ (retrieval quality assessment) of an information retrieval system for a given query q . Hence, a performance predictor $\hat{\mu}(q)$ has the following general form:

$$\hat{\mu}(q) \leftarrow \gamma(q, R_q, D_q, \mathcal{C}) \quad (5.1)$$

The prediction methods proposed in the literature establish different functions γ , and use a variety of available data, such as the query's terms, its properties with respect to the retrieval space (Cronen-Townsend et al., 2002), the output of the retrieval system – i.e., D_q and R_q – (Carmel et al., 2006), and the output of other systems (Aslam and Pavlu, 2007).

According to whether or not the retrieval results are used in the prediction process, such methods can be classified into pre-retrieval and post-retrieval approaches, which are described in Sections 5.2.1 and 5.2.2, respectively. Another relevant distinction is based on whether the predictors are trained or not, but this classification is less popular, and will not be considered here.

Moreover, the standard methodology to measure the effectiveness of performance prediction techniques (that is, the predictor quality assessment method) consists of comparing the rankings of several queries based on their actual precision – in terms of an evaluation metric such as MAP – with the rankings of those queries based on their performance scores, i.e., their predicted precision. In Section 5.4 we detail this methodology, along with several techniques for comparing the above rankings.

5.1.1 Notion of performance in Information Retrieval

In order to identify good performance predictors, validating or assessing their potential, we first have to define metrics of actual performance. Performance metrics and evaluation have been a core research and standardisation area for decades in the Information Retrieval field. In this section we introduce and summarise the main performance metrics and evaluation methodologies developed in the field.

The notion of performance in general, and in Information Retrieval in particular, leads itself to different interpretations, views and definitions. A number of methods

for measuring performance have been proposed and adopted (Hauff et al., 2008a; Hauff, 2010), the most prominent of which will be summarised herein; see (Baeza-Yates and Ribeiro-Neto, 2011) for an extended discussion.

As a result of several decades of research by the Information Retrieval community, a set of standard performance metrics has been established as a consensual reference for evaluating the goodness of information retrieval systems. These metrics generally require a collection of documents and a query (or alternative forms of user input such as item ratings), and assume a ground truth notion of relevance – traditional notions consider this relevance as binary, while others, more recently proposed, consider different relevance degrees.

One of the simplest and widespread performance metrics in Information Retrieval is **precision**, which is defined as the ratio of retrieved documents that are relevant for a particular query. In principle, this definition takes all the retrieved documents into account, but can also consider a given cut-off rank as the **precision at n** or **P@n**, where just the top-n ranked documents are considered. Other related and widespread metric is **recall**, which is the fraction of relevant documents retrieved by the system. These two metrics are inversely related, since increasing one generally reduces the other. For this reason, usually, they are combined into a single metric – e.g. the **F-measure**, and the **Mean Average Precision** or **MAP** –, or the values of one metric are compared at a fixed value of the other metric – e.g. the **precision-recall curve**, which is a common representation that consists of plotting a curve of precision versus recall, usually based on 11 standard recall levels (from 0.0 to 1.0 at increments of 0.1).

An inherent problem of using MAP for poorly performing queries, and in general of any query-averaged metric, is that changes in the scores of better-performing queries mask changes in the scores of poorly performing queries (Voorhees, 2005b). For instance, the MAP of a baseline system in which the effectiveness is 0.02 for a query A, and 0.40 for a query B, is the same as the MAP of a system where query A doubles its effectiveness (0.04) and query B decreases a 5% (0.38). In this context, in (Voorhees, 2005a) two metrics were proposed to measure how well information retrieval systems avoid very poor results for individual queries: the **%no measure**, which is the percentage of queries that retrieved no relevant documents in the top 10 ranked results, and the **area measure**, which is the area under the curve produced by plotting MAP(X) versus X, where X ranges over the worst quarter queries. These metrics were shown to be unstable when evaluated in small sets of 50 queries (Voorhees, 2005b). A third metric was introduced in (Voorhees, 2006): **gmap**, the geometric mean of the average precision scores of the test set of queries. This metric emphasises poorly performing queries while it minimises differences between larger scores, remaining stable in small sets of queries (e.g. 50 queries) (Voorhees, 2005b).

Nonetheless, despite the above metrics and other efforts made to obtain better measures of query performance, MAP, and more specifically the **Average Precision per query**, are still widely used and accepted. See (Carmel et al., 2006; Cronen-Townsend et al., 2002; Hauff et al., 2008b; He and Ounis, 2004; He et al., 2008; Kompaoré et al., 2007; Zhao et al., 2008; Zhou and Croft, 2006; Zhou and Croft, 2007), among others.

Almost as important as the performance metric is the query type, which can be related to the different user information needs (Broder, 2002). Most work on performance prediction has focused on the traditional ad-hoc retrieval task where query performance is measured according to topical relevance (also known as content-based queries). Some work – such as (Plachouras et al., 2003) and (Zhou and Croft, 2007) – has also addressed other types of queries such as named page finding queries, i.e., queries focused on finding the most relevant web page assuming the queries contain some form of the “name” of the page being sought (Voorhees, 2002a).

When documents are timed (e.g. a newswire system), we can also distinguish two main types of queries that have been only partially exploited in the literature (Diaz and Jones, 2004; Jones and Diaz, 2007): those queries that favour very recent documents, and those queries for which there are more relevant documents within a specific period in the past.

Finally, we note that most of the research ascribed to predict performance has been focused not on predicting the “true” performance of a query (whatever that means), but on discriminating those queries where query expansion or relevance feedback algorithms have proved to be efficient from those where these algorithms fail, such as polisemic, ambiguous, and long queries. These are typically called *bad-to-expand* queries (Cronen-Townsend et al., 2006), illustrating the implicit dependence on their final application.

5.1.2 A taxonomy of performance prediction methods

Existing prediction approaches are typically categorised into pre-retrieval methods and post-retrieval methods (Carmel and Yom-Tov, 2010). Pre-retrieval methods make the prediction before the retrieval stage, and thus only exploit the query’s terms and statistics about these terms gathered at indexing time. In contrast, post-retrieval methods use the rankings produced by a search engine, and, more specifically, the score returned for each document along with statistics about such documents and their vocabulary.

Category	Sub-category	Performance predictor (name and reference)
Pre-retrieval	Linguistics	Morphological, syntactic, semantic: (Mothe and Tanguy, 2005), (Kompaoré et al., 2007)
	Statistics	Coherency: coherence (He et al., 2008); term variance (Zhao et al., 2008) Similarity: collection query similarity (Zhao et al., 2008) Specificity: IDF-based (Plachouras et al., 2004), (He and Ounis, 2004); query scope (He and Ounis, 2004), (Macdonald et al., 2005); simplified clarity: (He and Ounis, 2004) Term relatedness: mutual information (Hauff et al., 2008a)
Post-retrieval	Clarity	Clarity (Cronen-Townsend et al., 2002), (Cronen-Townsend et al., 2006) Improved clarity (Hauff, 2010) (Hauff et al., 2008b) Jensen-Shannon Divergence (Carmel et al., 2006) Query difficulty (Amati et al., 2004)
	Robustness	Cohesion: clustering tendency (Vinay et al., 2006); spatial autocorrelation (Diaz, 2007); similarity (Kwok et al., 2004), (Grivolla et al., 2005) Document perturbation: ranking robustness (Zhou and Croft, 2006); document perturbation (Vinay et al., 2006) Query perturbation: query feedback (Zhou and Croft, 2007); autocorrelation (Diaz and Jones, 2004) (Jones and Diaz, 2007); query perturbation (Vinay et al., 2006); sub-query overlap (Yom-Tov et al., 2005a) Retrieval perturbation: (Aslam and Pavlu, 2007)
	Score analysis	Normalised Query Commitment: (Shtok et al., 2009) Standard deviation of scores: (Pérez-Iglesias and Araujo, 2009), (Cummins et al., 2011) Utility Estimation Framework: (Shtok et al., 2010) Weighted Information Gain: (Zhou and Croft, 2007)

Table 5.1. Overview of predictors presented in Section 5.2 categorised according to the taxonomy presented in (Carmel and Yom-Tov, 2010).

Pre-retrieval performance predictors do not rely on the retrieved document set, but on other information mainly extracted from the query issued by the user, such as statistics computed at indexing time (e.g. inverse term document frequencies). They have the advantage that predictions can be produced before the system's response is even started to be elaborated, which means that predictions can be taken

into account to improve the retrieval process itself. However, they have a potential handicap with regards to their accuracy on the predictions, since extra retrieval effectiveness cues available with the system's response are not exploited (Zhou, 2007). Pre-retrieval query performance has been studied from two main perspectives: based on probabilistic methods (and more generally, on collection statistics), and based on linguistic approaches. Most research on the topic has followed the former approach. Some researchers have also explored inverse document frequency (IDF) and related features as predictors, along with other collection statistics

Post-retrieval performance predictors, on the other hand, make use of the retrieved results. Broadly speaking, techniques in this category provide better prediction accuracy compared to pre-retrieval performance predictors. However, many of these techniques suffer from high computational costs. Besides, they cannot be used to improve the retrieval strategies without a post-processing step, as the output from the latter is needed to compute the predictions in the first place. In (Carmel and Yom-Tov, 2010) post-retrieval methods are classified as follows: 1) clarity based methods that measure the coherence (clarity) of the result set and its separability from the whole collection of documents; 2) robustness based methods that estimate the robustness of the result set under different types of perturbations; and 3) score analysis based methods that analyse the score distribution of results.

Table 5.1 shows a number of representative approaches on performance prediction, which will be described in the next section. These approaches are categorised according to the taxonomy and sub-categories proposed in (Carmel and Yom-Tov, 2010). In the table we can observe that the statistics category has been the most popular approach for pre-retrieval performance prediction. Several predictors have been categorised in the robustness category, probably due to its broad meaning (query, document, and retrieval perturbation). Finally, we note that recent effort from the community has been focused on the score analysis category.

5.2 Query performance predictors

In this section we explain the distinct performance predictors proposed in the literature. As mentioned before, based on whether or not retrieval results are needed to compute performance scores, predictors can be classified into two main types: pre-retrieval and post-retrieval predictors. In the following we summarise some of the approaches of each of the above types. For additional information, the reader is referred to (Carmel and Yom-Tov, 2010), (Hauff, 2010), and (Pérez Iglesias, 2012).

5.2.1 Pre-retrieval predictors

Pre-retrieval performance predictors do not rely on the retrieved document set, and exploit other collection statistics, such as the inverse document frequency (IDF). In this context, performance prediction has been studied from three main perspectives: based on linguistic methods, based on statistical methods, and based on probabilistic methods.

Linguistic methods

In (Mothe and Tanguy, 2005) and (Kompaoré et al., 2007) the authors consider 16 query features, and study their correlation with respect to average precision and recall. These features are classified into three different types according to the linguistic aspects they model:

- Morphological features:
 - **Number of words.**
 - **Average word length** in the query.
 - **Average number of morphemes per word**, obtained using the CELEX⁷ morphological database. The limit of this method is the database coverage, which leaves rare, new, and misspelled words as mono-morphemic.
 - **Average number of suffixed tokens**, obtained using the most frequent suffixes from the CELEX database (testing if each lemma in a topic is eligible for a suffix from this list).
 - **Average number of proper nouns**, obtained by POS (part-of-speech) tagger's analysis.
 - **Average number of acronyms**, detected by pattern matching.
 - **Average number of numeral values**, also detected by pattern matching.
 - **Average number of unknown tokens**, marked by a POS tagger. Most unknown words happen to be constructed words such as “mainstreaming”, “postmenopausal” and “multilingualism.”
- Syntactic features:
 - **Average number of conjunctions**, detected through POS tagging.
 - **Average number of prepositions**, also detected through POS tagging.
 - **Average number of personal pronouns**, again detected through POS tagging.
 - **Average syntactic depth**, computed from the results of a syntactic analyser. It is a straightforward measure of syntactic complexity in terms of

⁷ CELEX, English database (1993). Available at www.mpi.nl/world/celex

hierarchical depth; it simply corresponds to the maximum number of nested syntactic constituents in the query.

- **Average syntactic links span**, computed from the results of a syntactic analyser; it is the average pairwise distance (in terms of number of words) between individual syntactic links.
- Semantic features:
 - **Average polysemy value**, computed as the number of synsets in the WordNet⁸ database that a word belongs to, and averaged over all terms of the query.

In the above papers the authors investigated the correlation between these features, and precision and recall over datasets with different properties, and found that the only feature that positively correlated with the two performance metrics was the number of proper nouns. Besides, many variables did not obtain significant correlations with respect to any performance metric.

Statistical methods

Inverse document frequency is one of the most useful and widely used magnitudes in Information Retrieval. It is usually included in the information retrieval models to properly compensate how common terms are. Its formulation usually takes an ad hoc, heuristic form, even though formal definitions exist (Roelleke and Wang, 2008; Aizawa, 2003; Hiemstra, 1998). The main motivation for the inclusion of an IDF factor in a retrieval function is that terms that appear in many documents are not very useful for distinguishing a relevant document from a non-relevant one. In other words, it can be used as a measure of the specificity of terms (Jones, 1972), and thus as an indicator of their discriminatory power. In this way, IDF is commonly used as a factor in the weighting functions for terms in text documents. The general formula of IDF for a term t is the following:

$$\text{IDF}(t) = \log \frac{N}{N_t} \quad (5.2)$$

where N is the total number of documents in the system, and n_t is the number of documents in which the term t appears.

Some research work on performance prediction has studied IDF as a basis for defining predictors. He and Ounis (2004) propose a predictor based on the **standard deviation of the IDF** of the query terms. Plachouras et al. (2004) represent the quality of a query term by a modification of IDF where instead of the number of documents, the number of words in the whole collection is used (**inverse collection term**

⁸ WordNet, lexical database for the English language. Available at <http://wordnet.princeton.edu/>

frequency, or ICTF), and the query length acts as a normalising factor. These IDF-based predictors displayed moderate correlation with query performance.

Other authors have taken the similarity of the query into account. Zhao et al. (2008) compute the vector-space based query similarity with respect to the collection, considered as a large document composed of concatenation of all the documents. Then, different **collection query similarity** predictors are defined based on the SCQ values (defined below) for each query term, by summing, averaging, or taking the maximum values:

$$\text{SCQ}(t) = (1 + \log \text{TF}(t)) \cdot \text{IDF}(t) \quad (5.3)$$

The similarity of the documents returned by the query has also been explored in the field. The inter-similarity of documents containing query terms is proposed in (He et al., 2008) as a measure of **coherence**, by using the cosine similarity between every pair of documents containing each term. Additionally, two predictors based on the **pointwise mutual information** (PMI) are proposed in (Hauff et al., 2008a). The PMI of two terms is computed as follows:

$$\text{PMI}(t_1, t_2) = \log \frac{p(t_1, t_2)}{p(t_1)p(t_2)} \quad (5.4)$$

where these probabilities can be approximated by maximum likelihood estimations, that is, based on collection statistics, where $p(t_1, t_2)$ is proportional to the number of documents containing both terms, and $p(t) \propto \text{TF}(t)$. In that paper a first predictor is defined by computing the average PMI of every pair of terms in the query, whereas a second predictor is defined based on the maximum value. The predictive power of these techniques remains competitive, and is very efficient at run time.

Probabilistic methods

These methods measure characteristics of the retrieval inputs to estimate performance. He and Ounis (2004) propose a **simplified** version of the **clarity score** (see next section) in which the query model is estimated by the term frequency in the query:

$$\text{SCS} = \sum_w P_{ml}(w|q) \log_2 \frac{P_{ml}(w|q)}{P(w|\mathcal{C})} \quad (5.5)$$

$$P_{ml}(w|q) = \frac{\text{qtf}}{\text{ql}}; P(w|\mathcal{C}) = \frac{\text{TF}(w)}{|V|}$$

where **qtf** is the number of occurrences of a query term w in the query, **ql** is the query length, $\text{TF}(w)$ is the number of occurrences of a query term in the whole collection, and $|V|$ is the total number of terms in the collection.

Despite its original formulation, where the clarity score can be considered as a pre-retrieval predictor (Cronen-Townsend et al., 2002), Cronen-Townsend and colleagues use result sets to improve the computation time. For this reason, it is typically classified as a post-retrieval predictor (Zhou, 2007; Hauff et al., 2008a), and thus, we describe it with more detail in the next sections.

Kwok et al. (2004) build a query predictor using support vector regression, by training classifiers with features such as document frequencies and query term frequencies. In the conducted experiments they obtained a small correlation between predicted and actual query performances. He and Ounis (2004) propose the notion of **query scope** as a measure of the specificity of a query, which is quantified as the percentage of documents that contain at least one query term in the collection, i.e., $\log(N_Q/N)$, being N_Q the number of documents containing at least one of the query terms, and N the total number of documents in the collection. Query scope has shown to be effective in inferring query performance for short queries in ad hoc text retrieval, but very sensitive to the query length (Macdonald et al., 2005).

5.2.2 Post-retrieval predictors

Post-retrieval performance predictors make use of the retrieved results, in contrast to pre-retrieval predictions. Furthermore, computational efficiency is usually a problem for many of these techniques, which is balanced by better prediction accuracy. In the following we present the most representative approaches of each of the different sub-categories described in Section 5.1.2: clarity, robustness, and score analysis.

Clarity-based predictors

Cronen-Townsend et al. (2002) define **query clarity** as a degree of (the lack of) query ambiguity. Because of the particular importance and use of this predictor in the findings of this thesis, we shall devote a whole section (Section 5.3) for a thorough description and discussion about it. It is worth noting that the concept of query clarity has inspired a number of similar techniques. Amati et al. (2004) propose the **query difficulty** predictor to estimate query performance. In that work query performance is captured by the notion of the amount of information ($Info_{DFR}$) gained after the ranking. If there is a significant divergence in the query-term frequencies before and after the retrieval, then it is assumed that the divergence is caused by a query that is easy to respond to. $Info_{DFR}$ showed a significant correlation with average precision, but did not show any correlation between this predictor and the effectiveness of query expansion. The authors hence concluded that although the performance gains by query expansion in general increase as query difficulty decreases, very easy queries hurt the overall performance.

Adaptations of the query clarity predictor such as the one proposed in (Hauff et al., 2008b) will be discussed later in Section 5.3. Additionally, apart from the Kullback-Leibler divergence, the Jensen-Shannon Divergence on the retrieved document set and the collection also obtains a significant correlation between average precision and the distance measured (Carmel et al., 2006).

Robustness-based predictors

More recently, a related concept has been coined: **ranking robustness** (Zhou and Croft, 2006). It refers to a property of a ranked list of documents that indicates how stable a ranking is in the presence of *uncertainty* in its documents. The idea of predicting retrieval performance by measuring ranking robustness is inspired by a general observation in noisy data retrieval. The observation is that the degree of ranking robustness against noise is positively correlated with retrieval performance. This is because the authors assumed that regular documents also contain *noise*, if noise is interpreted as uncertainty. The robustness score performs better than, or at least as well as, the clarity score.

Regarding document and query perturbation, Vinay et al. (2006) propose four metrics to capture the geometry of the top retrieved documents for prediction: the **clustering tendency** as measured by the Cox-Lewis statistic, the sensitivity to **document perturbation**, the sensitivity to **query perturbation**, and the **local intrinsic dimensionality**. The most effective metric was the sensitivity to document perturbation, which is similar to the robustness score. Document perturbation, however, did not perform well for short queries, for which prediction accuracy dropped considerably when alternative state-of-the-art retrieval techniques (such as BM25 or a language modelling approach) were used instead of the TF-IDF weighting (Zhou, 2007).

Several predictors have been defined based on the concept of query perturbation. Zhou and Croft (2007) propose two performance predictors are defined based on this concept specifically oriented for Web search. First, the **Weighted Information Gain** predictor measures the amount of information gained about the quality of retrieved results (in response to a query) from an imaginary state that only an average document (represented by the whole collection) is retrieved to a posterior state that the actual search results are observed. This predictor was very efficient and showed better accuracy than clarity scores. The second predictor proposed in that work is the **Query Feedback**, which measures the degree of corruption that results from transforming Q to L (the output of the channel when the retrieval system is seen as a noisy channel, i.e., the ranked list of documents returned by the system). The authors designed a decoder that can accurately translate L back into a new query Q' , whereupon the similarity between the original query Q and the new query Q' is taken as a performance predictor, since the authors interpreted the evaluation of the quality of

the channel as the problem of predicting retrieval effectiveness. The computation of this predictor requires a higher computational cost than the previous one, being a major drawback of this technique.

Additionally, in (Diaz and Jones, 2004) and (Jones and Diaz, 2007) the authors exploited **temporal features** (time stamps) of the document retrieved by the query. They found that although temporal features are not highly correlated to performance, using them together with clarity scores improves prediction accuracy. Similarly, Diaz (2007) proposes to use the spatial autocorrelation as a metric to measure spatial similarities between documents in an embedded space, by computing the Moran's coefficient over the normalised scores of the documents. This predictor obtained good correlations results, although the author explicitly avoided collections such as question-answering and novelty related under the hypothesis that documents with high topical similarity should have correlated scores and, thus, in those collections the predictor would not work properly.

Other predictor was proposed in (Jensen et al., 2005), where visual features such as document titles and snippets are used from a surrogate document representation of retrieved documents. Such predictor was trained on a regression model with manually labelled queries to predict precision at the top 10 documents in Web search. The authors reported moderate correlation with respect to precision.

In (Yom-Tov et al., 2005a) two additional performance predictors are proposed. The first predictor builds a **histogram of the overlaps** between the results of each sub-query that agree with the full query. The second predictor is similar to the first one, but is based on a decision tree (Duda et al., 2001), which again uses overlaps between each sub-query and the full query. The authors apply these predictors to selective query expansion detecting missing content, and distributed information retrieval, where a search engine has to merge ranks obtained from different datasets. Empirical results showed that the quality of the prediction strongly depends on the query length.

The following predictors have been based on the cohesion of the retrieved documents. Kwok et al. (2004) propose predicting query performance by analysing similarities among retrieved documents. The main hypothesis of this approach is that relevant documents are similar to each other. Thus, if relevant documents are retrieved at the top ranking positions, the similarity between top documents should be high. The preliminary results, however, were inconclusive since negligible correlations were obtained. A similar approach is proposed in (Grivolla et al., 2005), where the entropy and pairwise similarity among top results are investigated. First, the entropy of the set of the K top-ranked documents for a query was computed. In this case it was assumed that the entropy should be higher when the performance for a given query is bad. Second, the mean cosine similarity between documents was proposed, using the base form of TF-IDF term weighting to define the document vec-

tors. Correlation between average precision and the proposed predictors was not consistent along the different systems used in the experiment, although the predictors could still be useful for performance prediction, especially when used in combination.

Predictors based on score analysis

Finally, the last family of post-performance predictors analyses the score distributions of the results for each query. We have to note that the Weighted Information Gain predictor (Zhou and Croft, 2007) explained above is sometimes categorised into this group. In the following we present other predictors where the retrieved scores are explicit in the predictor computation.

For instance, the **Normalised Query Commitment** (NQC) predictor (Shtok et al., 2009) measures the standard deviation of the retrieval scores, and applies a normalisation factor based on the score of the whole collection:

$$\text{NQC}(q) = \frac{\sqrt{1/|D_q| \sum_{d \in D_q} (s(d) - \mu_q)^2}}{|s(\mathcal{C})|} \quad (5.6)$$

where μ_q is the mean score of results in D_q (the retrieved set of documents for a query q). This predictor measures the divergence of results from their centroid, a “pseudo non-relevant document” that exhibits a relatively high query similarity (Carmel and Yom-Tov, 2010).

The **utility estimation framework** (UEF) was proposed in (Shtok et al., 2010) to estimate the utility of the retrieved ranking. In this framework three methods have to be specified to derive a predictor: a sampling technique for the document sets, a representativeness measure for relevance-model estimates, and a measure of similarity between ranked lists. Other authors have proposed approaches where standard deviation does not need to be computed for all the document scores in the retrieved results. Pérez-Iglesias and Araujo (2009) use a cutoff to decide how many documents are considered in the standard deviation computation. Moreover, Cummins et al. (2011) use different strategies to automatically select such cutoff.

Recently, Cummins (2012) has used Monte Carlo simulations to understand the correlations between average precision and the standard deviation of the scores in the head of a ranked list. The author found that the standard deviation of the list is positively correlated with the mean score of relevant documents, which in turn is positively correlated with average precision.

5.3 Clarity score

Cronen-Townsend et al. (2002) defined clarity score for Web retrieval as a measure of the lack of ambiguity of a particular query. More recently, it has been observed that this predictor also quantifies the diversity of the result list (Hummel et al., 2012). In this section we provide a deep analysis of this performance predictor since we shall use it along the rest of this thesis. We also describe examples and adaptations of the clarity score.

5.3.1 Definition of the clarity score

The clarity score predictor is defined as a Kullback-Leibler divergence between the query and the collection language model. It estimates the coherence of a collection with respect to a query q in the following way, given the vocabulary \mathcal{V} and a subset of the document collection R_q consisting of those documents that contain at least one query term:

$$\text{clarity}(q) = \sum_{w \in \mathcal{V}} p(w|q) \log_2 \frac{p(w|q)}{p(w|\mathcal{C})} \quad (5.7)$$

$$p(d|q) = p(q|d)p(d)$$

$$p(q|d) = \prod_{w_q \in q} p(w_q|d)$$

$$p(w|q) = \sum_{d \in R_q} p(w|d)p(d|q)$$

$$p(w|d) = \lambda p_{\text{ml}}(w|d) + (1 - \lambda)p_c(w)$$

The clarity value can thus be reduced to an estimation of the prior $p(w|\mathcal{C})$ (collection language model), and the posterior $p(w|q)$ of the query terms w (query language model) using $p(w|d)$ over the documents $d \in R_q$ and based on term frequencies and smoothing. It should be emphasised that if the set R_q is chosen as the whole collection \mathcal{C} , then this technique could be classified as a pre-retrieval performance predictor, since no information about the retrieval would be used. The importance of the size of the relevance set R_q (or number of feedback documents) has been studied in (Hauff et al., 2008b), where an adaptation of the predictor was proposed in order to automatically set the number of documents to consider.

As first published in (Cronen-Townsend et al., 2002) and (Cronen-Townsend et al., 2006), query ambiguity is defined as “the degree to which a query retrieves documents in the given collection with similar word usage.” Cronen-Townsend and

colleagues found that queries whose highly ranked documents are a mix of documents from disparate topics receive lower scores than if they result in a topically-coherent retrieved set, and reported a strong correlation between the clarity score and the performance of a query. Because of that, the clarity score method has been widely used in the area for query performance prediction.

Some applications and adaptations of the clarity score metric include query expansion (anticipating poorly performing queries that should not be expanded), improving performance in the link detection task (more specifically, in topic detection and tracking by modifying the measure of similarity of two documents) (Lavrenko et al., 2002), and document segmentation (Brants et al., 2002). More applications can be found in Section 5.3.3.

Zhou (2007) provides a complementary formulation of the clarity score by re-writing the formulation used above as follows:

$$\text{clarity}(q) = \sum_{w \in \mathcal{V}} \sum_{d \in R_q} p(w|d)p(d|q) \log \frac{\sum_{d \in R_q} p(w|d)p(d|q)}{p(w|\mathcal{C})} \quad (5.8)$$

In this way, Zhou emphasises, among other issues, the differences between the query clarity and the Weighted Information Gain predictor. Indeed, the author proposes the following generalisation of both formulations (for WIG and clarity). Specifically, the clarity formulation presented in Equations (5.7) and (5.8) is unified as follows:

$$\text{score}(q, \mathcal{C}, R) = \sum_{\xi \in T} \sum_{d \in R_q} \text{weight}(\xi, d) \log \frac{p(\xi, d)}{p(\xi, \mathcal{C})} \quad (5.9)$$

where T is a feature space, and R_q is a (ranked) document list. Besides this, $d \in R_q \subseteq \mathcal{C}$ must be comparable somehow with elements $\xi \in T$, in order to make sensible functions $\text{weight}(\xi, d)$ and $p(\xi, d)$. In this context, the query clarity as defined in (Cronen-Townsend et al., 2002) is an instantiation of Equation (5.9) where the following three aspects are considered:

- The feature space T is the whole vocabulary, consisting of single terms.
- The weight function is defined as $\text{weight}(\xi, d) = p(w|d)p(d|q)$.
- The function $p(\xi, d)$ is defined as $\sum_{d \in R_q} p(w|d)p(d|q)$, that is, it uses a document model averaged over all documents in the ranked list.

These observations help to discriminate between the underlying models used by these two predictors. In particular, for the query clarity, they also contribute to capture not so obvious divergences between a query and the collection, as we shall see in the next section.

train (0.33)	train dog (0.65)	obedience train dog (2.43)
	railroad train (0.73)	railroad train dog (0.67)
		railroad train caboose (1.46)

Table 5.2. Examples of clarity scores for related queries.

5.3.2 Interpreting clarity score in Information Retrieval

Aiming to better understand how the clarity score predictor behaves in Information Retrieval, and to what extent it is able to capture the difficulty or ambiguity of queries, in this section we summarise examples reported in the literature that let a clear interpretation of the predictor’s values.

In a seminal paper (Cronen-Townsend et al., 2002) Cronen-Townsend and colleagues present the example shown in Table 5.2, which provides the clarity scores of a number of related queries that share some of their terms. These queries are related to each other in the sense that a particular query is formed by extending other query with an additional term, starting with an initial query formed by a single term, ‘train’ in the example. According to the queries of the table, we can observe that the term ‘train’ has different meanings for the largest queries; it refers to ‘teach’ in the query ‘train dog’, to the ‘locomotive vehicle’ in the query ‘railroad train’, and can refer to any of both meanings in the query ‘railroad train dog.’ The clarity scores capture the ambiguity of the queries (due to their different meanings for the term ‘train’), independently from their length. In fact, the middle rightmost query ‘railroad train dog’ receives the lowest clarity score, corresponding to the most ambiguous query where the two considered meanings of ‘train’ are involved.

In the same paper, Cronen-Townsend and colleagues present the distribution of the language models for two queries, a clear query and a vague query (see Figure 2 in (Cronen-Townsend et al., 2002)). Each distribution is presented by plotting $p(w|q) \log_2 p(w|q)/p(w|\mathcal{C})$ against the query terms w . The authors show that the distribution of the values of this function for the clear query dominates the distribution of the values of the vague query. This makes sense since the clarity score is computed by summing the probability values in the distribution of every term in the collection. Additionally, the authors show that the clear query presents spikes in its query language model when $p(w|q)$ is plotted against the terms, and compared with the collection probability $p(w|\mathcal{C})$. Hence, some of the terms with high contribution from the query language model (i.e., with high $p(w|q)$ values) obtain low collection

probabilities ($p(w|\mathcal{C})$), thus evidencing a query that is different to the collection in its term usage (i.e., it is a non ambiguous query).

The above examples involve the (implicit) assumption known as *homogeneity assumption*, which specifies that the clarity score is higher if the documents in the considered collection are topically homogeneous. Hauff (2010) analyses the sensitivity of results with respect to that assumption. Specifically, the author computes the clarity score for three different ranked document lists: the relevant documents for a query, a non-relevant random sample, and a collection-wide random sample. The difference between the last two lists is that the second one is derived from documents judged as non-relevant, whereas the third one could contain any document in which at least one query term. Hauff shows how the clarity score is different depending on the origin of ranked document list, leading to a higher (lower) score by using relevant (non-relevant) documents for such list. However, we have to note that, as stated by Hauff, the quality in the separation of the clarity scores computed by each document list is different depending on the utilised dataset and queries.

The clarity score has been analysed in detail in Information Retrieval, mainly because its predictive power is superior to other performance predictors (in fact, it is one of the best performing post-retrieval predictors according to the overview presented in (Hauff, 2010)), but also because it provides interpretable results and high explanatory power in different IR processes, as we shall describe in the next section. Apart from that, the interest in this predictor is clear because of its probabilistic formulation and tight relationship with Language Models (Ponte and Croft, 1998).

5.3.3 Adaptations and applications of the clarity score

Cronen-Townsend and colleagues showed in (Cronen-Townsend et al., 2002) that clarity is correlated with performance, proving that the result quality is largely influenced by the amount of uncertainty involved in the inputs a system takes. In this sense, queries whose highly ranked documents belong to diverse topics receive lower scores than queries for which a topically-coherent result set is retrieved. Several authors have exploited the clarity score functionality and predictive capabilities (Buckley, 2004; Townsend et al., 2004; Dang et al., 2010), supporting its effectiveness in terms of performance prediction and high degree of adaptation. For instance, the predictor has been used for personalisation (Teevan et al., 2008) because of its proven capability of predicting ambiguity. In that paper the authors use more or less personalisation depending on the predicted ambiguity.

One of the first variants proposed in the area is the simplified clarity score proposed in (He and Ounis, 2004), presented in Section 5.2.1. In that paper He and Ouni changed the estimations of the posterior $p(w|q)$ to simple maximum likelihood estimators. Hauff et al. (2008b) proposed the Improved Clarity – called Adapted Clarity in (Hauff, 2010) –, in which the number of feedback documents

(R_q) is set automatically, and the term selection is made based on the frequency of the terms in the collection to minimise the contribution of terms with a high document frequency in the collection.

An alternative application of the clarity score is presented in (Allan and Raghavan, 2002), where the score obtained for the original set of documents returned by a query is compared against that obtained for a modified query, which was presumed to be more focused than the original one. Similarly, in (Buckley, 2004) Buckley uses the clarity score to measure the stability of the document rankings and compare it against a measure that uses the Mean Average Precision of each ranking (AnchorMap).

In (Sun and Bhowmick, 2009), Sun and Bhowmick adapted the concept of query clarity to image tagging, where a tag is visually representative if all the images annotated with that particular tag are visually similar to each other. In previous work (Sun and Datta, 2009) Sun and Datta proposed a similar concept, but in the context of blogging: a tag would receive a high clarity score if all blog posts annotated by the tag are topically cohesive.

Finally, an extension of the Kullback-Leibler divergence was proposed in (Aslam and Pavlu, 2007), where the Jensen-Shannon divergence was used instead. This distance is defined as the average of the Kullback-Leibler divergences of each distribution with respect to the average (or centroid) distribution. In this way, it is possible to compute the divergence between more than two distributions. Besides, the Jensen-Shannon divergence is symmetric, in contrast to the divergence used in the clarity score, and thus, a metric can be derived from it (Endres and Schindelin, 2003).

5.4 Evaluating performance predictors

In this section we describe the approaches proposed in the literature to evaluate the predictive power of a performance predictor. We define the different functions used to compute the quality of the performance predictors, most of them based on well known correlation coefficients between the true query performance values, and the expected or predicted performance values.

5.4.1 Task definition

Based on the notation presented in Section 5.1, in the following we present different techniques and functions to assess the effectiveness of performance predictors. Once the retrieval quality has been assessed ($\mu(q)$), and the value of the performance predictor for each query is calculated ($\hat{\mu}(q)$, using the function γ), the predictor quality is computed by using a predictor quality assessment function f^{qual} that measures the agreement between the true values of performance and the estimations, that is:

$$\text{Quality}(\gamma) = f^{qual}(\{\mu(q_1), \dots, \mu(q_n)\}, \{\hat{\mu}(q_1), \dots, \hat{\mu}(q_n)\}) \quad (5.10)$$

True quality values for each query are typically obtained by computing the per-query performance of a selected retrieval method (Cronen-Townsend et al., 2002; Hauff et al., 2008a), or by averaging the values obtained by several engines (Mothe and Tanguy, 2005), in order to avoid biases towards a particular method. As we shall see in the next section, the function f^{qual} typically represents a correlation coefficient; however, different possibilities are available and may be more appropriate depending on the prediction task.

In fact, in (Hauff et al., 2009) three estimation tasks were considered, by discriminating the output of the predictor function $\hat{\mu}$. **Query difficulty** estimation could be defined as a classification task where $\hat{\mu} \rightarrow \{0,1\}$ indicates whether the query is estimated to perform well or poorly. The standard estimation of **query performance**, nonetheless, would be defined by a function $\hat{\mu} \rightarrow \mathbb{R}$, in order to provide a ranking of queries, where the highest score denotes the best performing query. Furthermore, as stated in (Hauff et al., 2009), this function by itself does not directly estimate the performance metric μ . In order to do that we need to have normalised scores, such that the range of $\hat{\mu}$ is compatible with that of the metric, which typically requires $\hat{\mu} \rightarrow [0,1]$. In this case, we would be considering the **normalised query performance** task.

The methodology described above is general enough to be applicable to any of these three tasks, but is clearly inspired by the second one, that is, the estimation of query performance and it can be easily applied also to third one (normalised performance prediction). Because of that, we describe next a recently proposed methodology more focused on the (binary) query classification task or query difficulty prediction described in (Pérez-Iglesias and Araujo, 2010).

Let us suppose that, instead of continuous values of the performance metric μ , we are interested in estimating *as accurately as possible* the different difficulty grades of the queries, that is, $\mu \rightarrow \{1, \dots, k\}$, where k is the number of difficulty grades available. Obviously, the output of the predictor $\hat{\mu}$ also has to be grouped in one of the k classes. Typically, we would have $k = 3$, representing “Easy”, “Average”, and “Hard” queries, although a binary partition could also be acceptable. In these terms the performance prediction problem is stated as a classification problem, where the goal is to effectively predict the query class.

Furthermore, this technique lets set, at the quality computation step, whether we want to weight uniformly each of the k classes, or if we are more interested in only one of them, by building, for instance, a confusion matrix, and applying standard Machine Learning evaluation metrics to a subset of it. In the next section we describe the most popular techniques for doing this, along with a new metric introduced in (Pérez-Iglesias and Araujo, 2010) oriented to the problem of performance prediction.

5.4.2 Measuring the quality of the predictors

There are several methods for measuring the quality of the performance prediction function $\hat{\mu}$ defined in the previous section. In particular, the quality function f^{qual} may be able to capture linear relations, take into account the importance implied by the scores or the ordering given by each variable (true and estimated performance, i.e., μ and $\hat{\mu}$), and exploit the implicit partitions derived by the method.

The most commonly used quality function is correlation, which has been measured by three well-known metrics: Pearson's, Spearman's, and Kendall's correlation coefficients. **Pearson's r correlation** captures linear dependencies between the variables, whereas **Spearman's ρ** and **Kendall's τ** correlation coefficients are used in order to uncover non-linear relationships between the variables. They are generally computed as follows, although in special situations (in presence of ties, or when there are missing values in the data) alternative formulations may be used:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (5.11)$$

$$\rho = 1 - \frac{6 \sum_{i=1}^n d(x_i, y_i)^2}{n(n^2 - 1)} \quad (5.12)$$

$$\tau = 1 - \frac{4Q(x, y)}{n(n - 1)} \quad (5.13)$$

where x and y represent the two variables of interest, \bar{x} and \bar{y} denote their means, $d(x_i, y_i)$ is the difference in ranks between x_i and y_i , and $Q(x, y)$ is the minimum number of swaps needed to convert the rank ordering of x to that of y . All these coefficients return values between -1 and $+1$, where -1 denotes a perfect anti-correlation, 0 denotes statistical independence, and $+1$ denotes perfect correlation.

It can be observed that Spearman's ρ computes a Pearson's r between the ranks induced by the scores of the variables. Moreover, Kendall's τ is the number of operations required to bring one list to the order of the other list using the *bubble sort* algorithm. Besides, although Spearman's and Kendall's correlations seem more general than Pearson's since they are able to capture non-parametric relations between the variables, we have to consider that distances between the scores are ignored in the rank-based coefficients, and thus, it is typically suggested to report one correlation coefficient of each type.

It is important to note that the number of points used to compute the correlation values affects the significance of the correlation results. The confidence test for a Pearson's r correlation, modeled as the t -value of a t -distribution (assuming normality) with $N - 2$ degrees of freedom (being N the size of the sample), is defined by the following equation (Snedecor and Cochran, 1989):

<i>p</i> -value	N			Pearson's <i>r</i> value	N		
	50	100	500		50	100	500
<i>p</i> < 0.05	1.677	1.661	1.648	0.1	0.696	0.995	2.243
<i>p</i> < 0.01	2.407	2.365	2.334	0.2	1.414	2.021	<u>4.555</u>
				0.3	2.179	<u>3.113</u>	<u>7.018</u>
				0.4	<u>3.024</u>	<u>4.320</u>	<u>9.739</u>

Table 5.3. Left: minimum *t*-value for obtaining a significant value with different sample sizes (N). **Right:** *t*-value for a given Pearson's correlation value and N points. In bold when the correlation is significant for *p* < 0.05, and underlined for *p* < 0.01.

$$t = r \sqrt{\frac{N - 2}{1 - r^2}} \quad (5.14)$$

The *t*-value therefore depends on the size of the sample, and thus, the significance of a Pearson's correlation value *r* may change depending on the number of test queries. In particular, for small samples, we may eventually obtain strong but non-significant correlations; whereas for large samples, on the other hand, we may obtain significant differences, even though the strength of the correlation values may be lower. The above also applies to the correlations computed using the Spearman's coefficient, but only under the null hypothesis or large sample sizes (greater than 100) (Snedecor and Cochran, 1989; Zar, 1972). For Kendall's correlation, the confidence test can be computed using an exact algorithm when there are no ties based on a power series expansion in N^{-1} , depending again, thus, on the sample size (Best and Gipps, 1974).

Table 5.3 shows the minimum *t*-value for obtaining a significant value with different sample sizes and *p*-values, along with the *t*-value computed using Equation (5.14) for different correlation values and sample sizes. In the table we can observe that the same correlation value may be significant or not depending on the size of the sample, for instance, with 50 queries, observations are significant with *p* < 0.05 for correlation values equal or above 0.3, whereas for 100 queries it is enough to obtain Pearson's correlation values of 0.2. This observation is related to the one presented in (Hauff et al., 2009), where Hauff and colleagues compared the confidence intervals of the three correlation coefficients described before, and observed how, due to the small query set sizes, most of the predictors analysed (pre-retrieval approaches such as clarity, IDF-based, and PMI) presented no significant differences, despite having very different values. In particular, this generated a subset of the analysed predictors that were not statistically different to the best performing predictor reported, and thus, any of the predictors in subset may be used in a later application since they obtain statistically similar (strictly speaking, not statistically different) correlations.

Furthermore, in the same paper, Hauff and colleagues proposed to use the **Root Mean Squared Error** (RMSE) as a quality function. The rationale behind this is that the RMSE squared is the function being minimised when performing a linear regression, and thus, it should also be able to capture the (linear) relation between the variables. In fact, there is a close relation between the RMSE and the Pearson's r coefficient, by means of the residual sum of squares (Carmel and Yom-Tov, 2010):

$$r^2 = 1 - \frac{SS_{err}}{SS_{tot}} = 1 - \frac{\sum_{i=1}^n (x_i - y_i)^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (5.15)$$

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (x_i - y_i)^2}{n - 1}} = \sqrt{\frac{SS_{err}}{n - 1}} \quad (5.16)$$

Additional extensions to these correlation coefficients have been proposed. Most of these extensions have been focused on incorporating weights in the computation of the correlation (Melucci, 2009; Yilmaz et al., 2008). However, despite these metrics have an evident potential in the performance prediction area, to the best of our knowledge there is no work using them in order to evaluate the quality of the predictors (Pérez Iglesias, 2012).

Finally, a different family of quality functions can be considered in the query difficulty task, that is, when the performance prediction is cast as a classification problem. These techniques are based on the accuracy of the classification provided by the performance predictor, and thus, classic Machine Learning techniques could be used. In (Pérez-Iglesias and Araujo, 2010), Pérez-Iglesias and Araujo propose to use the **F-measure**:

$$F = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \quad (5.17)$$

Additionally, in the same paper, Pérez-Iglesias and Araujo introduced a new metric (**distance based error measure**, or DBEM) along with a methodology that is focused on the misclassified difficulty classes between the predictor and the true classes. With this goal in mind, the authors apply a clustering algorithm to both the performance metric values and their estimations, aimed to minimise the distance between elements in the same group, and maximise the distance between elements in different groups. Specifically, Pérez-Iglesias and Araujo used the k -means algorithm, setting the value of k to the number of relevance grades, $k = 3$ in their paper. The metric DBEM is defined as follows:

$$\text{DBEM} = \frac{\sum_i^n \text{dist}(c(x_i), c(y_i))}{\sum_i^n \max_j \text{dist}(c(x_i), c(x_j))} \quad (5.18)$$

$$\text{dist}(c_i, c_j) = \|i - j\|, 0 < i, j \leq k$$

where $c(x)$ is the function which assigns the proper class or partition to a given score x , according to the clustering algorithm. This metric captures the distance between every partition, normalised by the maximum possible distance. In this case, lower distances imply a better predictor quality.

5.5 Summary

Improvement of the predictive capabilities to infer the performance or difficulty of a query is consolidated as a major research topic in Information Retrieval, where it has been mostly applied to ad-hoc retrieval. Several performance predictors have been defined based on many different information sources, demonstrating the usefulness of such predictors in different tasks, mainly for query expansion, but also for rank fusion, distributed information retrieval, and text segmentation.

Some issues are, however, still open in the field, mostly regarding the evaluation of performance prediction. Performance prediction methods have been usually evaluated on traditional TREC document collections, which typically consist of no more than one million relatively homogenous newswire articles, and few research work has exploited these techniques with larger datasets; see, e.g. (Carmel et al., 2006; Zhou, 2007; Hauff, 2010) for some exceptions. Furthermore, reported correlation coefficient values have been typically computed using a small number of points (e.g. 50 queries for standard tracks in TREC), not always providing enough confidence to derive conclusions. And more importantly, how predictors have to be evaluated and which metric has to be used are still open research questions, that have generated some fruitful discussion in recent publications (Hauff, 2010; Pérez Iglesias, 2012), although a definitive answer has not been obtained yet.

We may presume that in the future other information retrieval applications may benefit from the framework derived by these techniques, and may develop tailored performance predictors by using purpose-designed performance metrics and evaluation methodologies, such as the recently developed concept of document difficulty in (Alvarez et al., 2012). This thesis is an example of such an application in the Recommender Systems field. More specifically, as we shall see in the next chapter, we translate the problem of performance prediction to the Recommender Systems area, where it has been barely studied. We focus our research on the query clarity predictor as a basis for the recommendation performance predictors, although additional techniques could be used, as we shall also present in Chapter 6. Finally, among the array of evaluation strategies presented above, we have decided to use correlations since it is the most common one in the literature, and provides a fair notion about the interpretability of the results.

Chapter 6

Performance prediction in recommender systems

In this chapter, we state and address the recommendation performance prediction problem, proposing and evaluating different prediction schemes. After laying out a formal frame for the problem, we start by researching the adaptation of principles and prediction techniques that have been proposed and developed in ad-hoc Information Retrieval. More specifically, we draw from the notion of query clarity as a basis for finding suitable performance predictors that provide a well grounded theoretical formalisation. In analogy to query clarity, we hypothesise that the amount of uncertainty involved in user and item data (reflecting ambiguity in user's tastes and item popularity patterns) may also correlate with the accuracy of the system's recommendations. This uncertainty can be captured as the clarity of users and/or the clarity of items by an adaptation of the query clarity formulation. This adaptation, however, is not straightforward, as we shall describe. Besides the approaches elaborating on the notion of clarity, we propose new predictors based on theories and models from Information Theory and Social Graph Theory.

In Section 6.1 we formulate the research problem we aim to address. Next, in Sections 6.2, 6.3, and 6.4 we propose several performance predictors for recommender systems, some of them based on the clarity score, information theoretical related concepts – such as entropy –, and graph-based metrics. The proposed predictors are defined upon three different spaces, namely ratings, logs, and social networks. Moreover, we also provide specific correlations of the described predictors in Section 6.5 in order to show their predictive power under different conditions along with a discussion of the results. Finally, in Section 6.6 we provide some conclusions.

6.1 Research problem

Performance prediction finds a special motivation in recommender systems. Contrary to query-based information retrieval, as far as the initiative relies on the system, a performance prediction approach may provide a basis to decide producing recommendations or holding them back, depending on the expected level of performance on a per case basis, delivering only the sufficiently reliable cases. On the other hand, recommenders based on a single algorithm are not competitive in practice, and real applications heavily rely on hybridisations and ensembles of algorithms.

The capability to foresee which algorithm can perform better in different circumstances can therefore be envisioned as a good approach to enhance the performance of the combination of algorithms by dynamically adjusting the reliance on each subsystem. Furthermore, it is well-known in the recommender systems field that the performance of individual recommendation methods is highly sensitive to different conditions, such as data sparsity, quality and reliability, which are subject to an ample dynamic variability in real settings. Hence, being able to estimate in advance which recommenders are likely to provide the best output in a particular situation opens up an important window for performance enhancement. Alternatively, estimating which users of a system are likely to receive worse recommendations allows for modifications in the recommendation algorithms to predict this situation, and react in advance.

The problem of performance prediction has been however barely addressed in the Recommender Systems field. The issue has been nonetheless mentioned in the literature – evidencing the relevance of the problem – and is in some way often implicitly addressed by means of ad hoc heuristic tweaks such as significance weighting in nearest neighbour recommenders (Herlocker et al., 1999) and confidence scores (Wang et al., 2008a), along with additional computations (mainly normalisations) which are introduced into the recommendation methods aimed to better estimate the predicted ratings.

In the recommendation context, the problem of performance prediction can be stated as follows. We define a performance predictor as a function that takes a certain input, and returns a real value that correlates with some utility dimension of a recommender system. This is an instantiation of the problem presented in Section 5.1 but in the recommendation setting. For such purpose, we first specify more precisely what the input space of predictors consists of, and how the predictor input and output relate to the data involved in recommendation. Thus, a utility predictor handles the following information:

Input variables

- The specific configuration of the recommender system. For instance, for a nearest neighbour recommender input parameters could be the neighbour map (that assigns a set of neighbours to each user) and a user similarity metric.
- Any input of the recommender, such as the active user and the active item.
- Background/context information: any known user, item, and user-item interaction data, such as user ratings, user features, item features, social network information, data timestamps, etc. We have to note that, even though the predictor will generally use this type information, we consider it as implicit input and do not include it explicitly in our notation to avoid making it needlessly cumbersome.

Output variable

- A value in \mathbb{R} .

A predictor is thus a function $\gamma: R \times \mathcal{U} \times \mathcal{J} \rightarrow \mathbb{R}$ (R being the set of all recommenders) that estimates the performance of the system, possibly using additional information available in the background. A predictor can be independent from some of these inputs, which would be then omitted in the previous notation. For instance, in this chapter we shall present predictors of the form $\gamma: \mathcal{U} \rightarrow \mathbb{R}$ and $\gamma: \mathcal{J} \rightarrow \mathbb{R}$. Additionally, a predictor may assume a specific parameterised recommender algorithm family (e.g. nearest neighbour collaborative filtering), and needs some element of its configuration as input. It may also happen that a predictor does not make any assumption on the recommender – it does not depend on it – but still the predictor works well only for certain types of recommenders. It would be syntactically possible and correct to apply the predictor with other recommenders, although it may work badly. In general, what it means for a predictor to work “well” may depend on the application, but we generally assume it can be evaluated in terms of its correlation to some utility dimension of recommendations, such as an accuracy metric (RMSE, precision, nDCG) or alternative metrics such as novelty, diversity, etc.

If a recommender system can be decomposed into its internal configuration, then a predictor can directly take as input the components of the recommender configuration. For instance, neighbourhood-based collaborative filtering recommenders can be represented in $R \equiv \mathcal{E} \times \mathcal{N} \times \mathcal{S}$, where $\mathcal{E}: \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}^+$ is a preference estimation function (based on k similarity values between the target user and her neighbours, and k neighbours’ ratings on the target item), $\mathcal{N}: \mathcal{U} \rightarrow \mathcal{P}(\mathcal{U})$ is a neighbourhood assignment map, and \mathcal{S} is a similarity metric. Upon such a model, we would have $\gamma: \mathcal{E} \times \mathcal{N} \times \mathcal{S} \times \mathcal{U} \times \mathcal{J} \rightarrow \mathbb{R}$.

We may also constrain some inputs to a relevant condition they should meet. For instance, we could limit ourselves to a neighbourhood map that considers a user v as a candidate neighbour. In that case, this map can be essentially represented by v , and then we would have $\gamma: \mathcal{E} \times \mathcal{U} \times \mathcal{S} \times \mathcal{U} \times \mathcal{I} \rightarrow \mathbb{R}$ (note that the first \mathcal{U} in the Cartesian product stands for neighbour users, and the second \mathcal{U} for target users).

It is important to note that when the predictor takes as input some of the inputs of the recommender, namely the active user and/or the active item, then the predictor's correlation with the recommender's utility must be measured on a per-input basis. For instance, if the predictor just takes users as input arguments, it should correlate with the average utility per user.

Moreover, predictors can also be used to enhance hybrid recommenders by favouring strategies that are predicted to produce better results. This can be done by relating activation switches in the recommenders to predictor values, so that one recommender or the others are activated or favoured depending on the predictor's estimation.

The way in which these activation switches are related to predictors is typically application-dependent. For instance, in ensemble recommenders consisting of a unique (Boolean) selection among a set of recommenders, the selection/discarding of recommenders can be a binary function of a predictor for each recommender. If the ensemble consists of a linear combination of recommenders, the weights in the combination can also be a function of the predictors. In neighbourhood-based collaborative filtering, activation switches can be the weights of neighbours in the prediction of user ratings. Indeed, relating predictor values to activation switches is a non-trivial problem and generally requires some research on itself.

Based on all the above mentioned issues, the general research problem we address consists of a) finding effective predictors of recommendation utility, and b) identifying and testing useful applications for the found predictors. In the remainder of this chapter we propose different predictors of recommendation utility using different types of input, namely ratings, logs, and social information. In Chapters 7 and 8 we shall exploit and evaluate such predictors in two applications: dynamic hybrid recommendation, and dynamic neighbour weighting in collaborative filtering.

6.2 Clarity for preference data: adaptations of query clarity

In this thesis, we propose different adaptations for the concept of query clarity to recommender systems. First, we deal with the definition of user clarity when rating-based preference data is available, where alternative ground models are proposed, depending on which random variables want to be considered in the computation of

the user clarity. Then, we define the concept of user clarity for log-based preference data. Additionally, for ratings we also define the concept of item clarity.

Now we propose a fairly general adaptation of query clarity, which may be instantiated into different schemes, depending on the input spaces considered. At an abstract level, we consider an adaptation that equates users in the recommendation domain to queries in the search domain, as the corresponding available representations of user needs in the respective domains. This adaptation results in the following formulation for **user clarity**:

$$\text{clarity}(u) = \sum_{x \in \mathcal{X}} p(x|u) \log_2 \frac{p(x|u)}{p(x)} \quad (6.1)$$

As we can observe, the clarity formulation strongly depends on a “vocabulary” space \mathcal{X} , which further constrains the user-conditioned model (or user model for short) $p(x|u)$, and the background probability $p(x)$. In ad-hoc information retrieval, this space is typically the space of words, and the query language model is a probability distribution over words (Cronen-Townsend et al., 2002). In recommender systems, however, we may have different interpretations, and thus, different formulations for such a probabilistic framework, as we shall show. In all cases, we will need to model and estimate two probability distributions: first, the probability that some event (depending on the current probability space \mathcal{X}) is generated by the user language model (user model); and second, the prior probability of generating that event (background model).

Under this formulation, user clarity is in fact the difference (Kullback-Leibler divergence) between a user model and a background model. The use of user and background distributions as a basis to predict recommendation performance lies on the hypothesis that a user probability model being close to the background (or collection) model is a sign of ambiguity or vagueness in the evidence of user needs, since the generative probabilities for a particular user are difficult to single out from the model of the collection as a whole. In Information Retrieval, this fact is interpreted as a query for which the relevant documents are a mix of articles about different topics (Cronen-Townsend et al., 2002).

As an additional step, we generalise the adaptation stated in Equation (6.1) to allow for different reference probability models parameterised by a generic variable θ .

$$\text{clarity}(u) = \mathbb{E}_\theta \left[\sum_{x \in \mathcal{X}} p(x|u, \theta) \log_2 \frac{p(x|u, \theta)}{p(x|\theta)} \right] \quad (6.2)$$

This generalisation will allow for the development of further varieties of the clarity scheme, and simplifies to Equation (6.1) whenever we implicitly consider a fixed θ , as we shall see next. Equivalently, the variable θ may be integrated in both user and background models by exploiting a multidimensional vocabulary space:

$$\text{clarity}(u) = \sum_{x \in \mathcal{X}, \theta \in \Theta} p(x, \theta | u) \log_2 \frac{p(x, \theta | u)}{p(x, \theta)} \quad (6.2b)$$

It is easy to see that Equations (6.2) and (6.2b) are fully equivalent, and thus allow two interpretations for the same magnitude.

As stated in (Cronen-Townsend et al., 2002), language models capture statistical aspects of the generation of language. Therefore, if we use different vocabularies, we may capture different aspects of the user. The probabilistic relations between the variables involved in Equation (6.2) also depend on the nature of the data, and the different possible generative models induced by the recorded observations of user-item interactions (the input to a recommender system). In this thesis we consider two types of interaction data records: users-rating-items (where the atomic event is a user rating an item with a value), and users “consuming” items (a user accesses an item at some time instant). The first type fits a dataset such as MovieLens and CAMRa, and the second fits well Last.fm data – the datasets on which we shall test the methods to be developed here. Across these two types, in our research we explore mainly three vocabulary spaces for \mathcal{X} : ratings, items, and time. Each of the vocabulary spaces induces its own user-specific interpretation, as we shall see. As for the optional contextual parameter θ , we shall consider here only the space of items ranging over the set of items – thus fully leveraging the triadic nature of the user-item-rating and user-item-time spaces. The scheme is however open to the exploration of further possibilities, as is the vocabulary space itself, beyond the options researched here.

In the following sections we thus explore several alternatives for rating-based and log-based data spaces (and their induced generative models).

6.2.1 Rating-based clarity

As just mentioned, in the rating space, we consider a set of user-item-rating tuples, where each user-item pair appears in a unique tuple (i.e., users only rate items once). We consider two possible vocabulary spaces: items and ratings, and two context alternatives: items (which make only sense in the rating vocabulary) and none. The resulting clarity schemes are summarised in Table 6.1, and have each their own interpretation.

The rating-based clarity model captures how differently a user uses rating values (regardless of the items the values are assigned to) with respect to the rest of users in the community. The item-based clarity takes into account which items have been rated by a user, and therefore, whether she rates (regardless of the rating value) the most rated items in the system or not. Finally, the item-and-rating-based clarity computes how likely a user would rate each item with some particular rating value, and compares that likelihood with the probability that the item is rated with some particular rating value. In this sense, the item-based user model makes the assumption that some items are more likely to be generated for some users than for others de-

User clarity	Vocabulary \mathcal{X} / Context θ	User model	Background model	Formulation
Rating-based	Ratings / None	$p(r u)$	$p_c(r)$	$\sum_r p(r u) \log_2 \frac{p(r u)}{p(r)}$
Item-based	Items / None	$p(i u)$	$p_c(i)$	$\sum_i p(i u) \log_2 \frac{p(i u)}{p(i)}$
Item-and-rating-based	Ratings / Items	$p(r u, i)$	$p_{ml}(r i)$	$\sum_{r,i} p(i)p(r u, i) \log_2 \frac{p(r u, i)}{p(r i)}$

Table 6.1. Three possible user clarity formulations, depending on the interpretation of the vocabulary and context spaces.

pending on their previous preferences. The rating-based model, on the other hand, captures the likelihood of a particular rating value being assigned by a user, which is an event not as sparse as the previous one, with a larger number of observations. Finally, the item-and-rating-based model is a combination of the two previous models into a unified model incorporating items and ratings. As we mentioned before, this could be made more explicit by considering the user model $p(r, i|u)$ in the Equation (6.2b), which would be equivalent to this model under some independence assumptions, i.e., when $p(r, i|u) = p(r|u, i)p(i)$.

Ground models for user clarity

We ground the different clarity measures defined in the previous section upon a rating-oriented probabilistic model very similar to the approaches taken in (Hofmann, 2004) and (Wang et al., 2008a). The sample space for the model is the set $\mathcal{U} \times \mathcal{I} \times \mathcal{R}$, where \mathcal{U} stands for the set of all users, \mathcal{I} is the set of all items, and \mathcal{R} is the set of all possible rating values. Hence, an observation in this sample space consists of a user assigning a rating to an item. We consider three natural random variables in this space: the user, the item, and the rating value, involved in a rating assignment by a user to an item. This gives meaning to the distributions expressed in the different versions of clarity as defined in the previous section. For instance, $p(r|i)$ represents the probability that a specific item i is rated with a value r – by a random user –, $p(i)$ is the probability that an item is rated – with any value by any user –, and so on.

The probability distributions upon which the proposed clarity models are defined can use different estimation approaches, depending on the independence assumptions one would consider, and the amount of involved information. Background models are estimated using relative frequency estimators, that is:

$$p_c(r) = \frac{|\{(u, i) \in \mathcal{U} \times \mathcal{I} | r(u, i) = r\}|}{|\{(u, i) \in \mathcal{U} \times \mathcal{I} | r(u, i) \neq \emptyset\}|} \quad (6.3)$$

$$p_c(i) = \frac{|\{u \in \mathcal{U} | r(u, i) \neq \emptyset\}|}{|\{(u, j) \in \mathcal{U} \times \mathcal{I} | r(u, j) \neq \emptyset\}|}$$

$$p_{ml}(r|i) = \frac{|\{u \in \mathcal{U} | r(u, i) = r\}|}{|\{u \in \mathcal{U} | r(u, i) \neq \emptyset\}|}$$

$$p_{ml}(r|u) = \frac{|\{i \in \mathcal{I} | r(u, i) = r\}|}{|\{i \in \mathcal{I} | r(u, i) \neq \emptyset\}|}$$

These are maximum likelihood estimations in agreement with the meaning of the random variables as defined above. Starting from these estimations, user models can be reduced to the above terms by means of different probabilistic expansions and Bayesian reformulations, which we define next for the three models introduced in the previous section.

Item based model. The $p(i|u)$ model can be simply expanded through marginalisation over ratings, but under two different assumptions: the item generated by the model only depends on the rating value, independently from the user or, on the contrary, depends on both the user and the rating. These alternatives lead to the following developments, respectively:

$$p_R(i|u) = \sum_{r \in \mathcal{R}} p_{ml}(i|r) p_{ml}(r|u) \quad (6.4)$$

$$p_{UR}(i|u) = \sum_{r \in \mathcal{R}} p(i|u, r) p_{ml}(r|u) \quad (6.5)$$

Rating based model. This model assumes that the rating value generated by the probability model depends on both the user and the item at hand. For this model, we sum over all possible items in the following way:

$$p(r|u) = \sum_{r(u, i)=r} p(r|u, i) p(i|u) \quad (6.6)$$

where the $p(i|u)$ term can be developed as in the item-based model above. The term $p(r|u, i)$ requires further development, which we define in the next model.

Item-and-rating based model. Three different models can be derived depending on how the Bayes' rule is applied. In these models, item probability is assumed to be uniform and thus it can be ignored in the computation of the expectation in Equation (6.2). In the same way as proposed in (Wang et al., 2008a), three relevance models can be defined, namely a user-based, an item-based, and a unified relevance model:

$$p_U(r|u, i) = \frac{p(u|r, i) p_{ml}(r|i)}{\sum_{r \in \mathcal{R}} p(u|r, i) p_{ml}(r|i)} \quad (6.7)$$

$$p_I(r|u, i) = \frac{p(i|u, r)p_{mI}(r|u)}{\sum_{r \in \mathcal{R}} p(i|u, r)p_{mI}(r|u)} \quad (6.8)$$

$$p_{UI}(r|u, i) = \frac{p(u, i|r)p_c(r)}{\sum_{r \in \mathcal{R}} p(u, i|r)p_c(r)} \quad (6.9)$$

The first derivation induces a user-based relevance model because it measures by $p(u|r, i)$ how probable it is that a user rates item i with a value r . The item-based relevance model is factorised proportional to an item-based probability, i.e., $p_I(r|u, i) \propto p(i|u, r)$. Finally, in the unified relevance model, we have $p_{UI}(r|u, i) \propto p(u, i|r)$. These estimations correspond respectively with the Equations 20a, 20b, and 21 from (Wang et al., 2008a); to make the thesis self-contained and facilitate the comparison between the different probability models, we present now these equations from (Wang et al., 2008a):

$$p(u|r, i) = \frac{1}{|S(r, i)|} \sum_{v \in S(r, i)} \frac{1}{h_u^{|j|}} K\left(\frac{\mathbf{u} - \mathbf{v}}{h_u}\right) \quad (6.10)$$

$$p(i|u, r) = \frac{1}{|S(r, u)|} \sum_{j \in S(r, u)} \frac{1}{h_i^{|u|}} K\left(\frac{\mathbf{i} - \mathbf{j}}{h_i}\right) \quad (6.11)$$

$$p(u, i|r) = \frac{1}{|S(r)|} \sum_{(v, j) \in S(r)} \frac{1}{h_u^{|j|}} K\left(\frac{\mathbf{u} - \mathbf{v}}{h_u}\right) \frac{1}{h_i^{|u|}} K\left(\frac{\mathbf{i} - \mathbf{j}}{h_i}\right) \quad (6.12)$$

where $K(\cdot)$ is a Parzen Kernel function (Duda et al., 2001). In this formulation, \mathbf{u} denotes the user u represented as a vector by her ratings in the space of items. Unrated items can be filled with the average rating value or with other constant value, such as 0 or the average rating in the community. Respectively, \mathbf{i} represents the item i in the user space. h_u and h_i are the bandwidth window parameter for the user and item vector, respectively; $S(\cdot)$ denotes the set of observed samples where event (\cdot) has happened. For example, $S(r, i)$ denotes the set of observed samples with event $(R = r, I = i)$. More specifically:

$$S(r, i) = \{u \in \mathcal{U} | r(u, i) = r\} \quad (6.13)$$

$$S(r, u) = \{i \in \mathcal{I} | r(u, i) = r\} \quad (6.14)$$

$$S(r) = \{(u, i) \in \mathcal{U} \times \mathcal{I} | r(u, i) = r\} \quad (6.15)$$

In the experiments, we used a Gaussian Kernel function, i.e., $K(\mathbf{x}) = e^{-x^2/2}/\sqrt{2\pi}$, and $h_i = h_u = 0.9$ as suggested in (Wang et al., 2008a).

User clarity name	User dependent model	Background model
RatUser	$p_U(r u, i); p_{UR}(i u)$	$p_c(r)$
RatItem	$p_I(r u, i); p_{UR}(i u)$	$p_c(r)$
ItemSimple	$p_R(i u)$	$p_c(i)$
ItemUser	$p_{UR}(i u)$	$p_c(i)$
IRUser	$p_U(r u, i)$	$p_{mi}(r i)$
IRItem	$p_I(r u, i)$	$p_{mi}(r i)$
IRUserItem	$p_{UI}(r u, i)$	$p_{mi}(r i)$

Table 6.2. Different user clarity models implemented.

Finally, different combinations of distribution formulations and estimations result in a fair array of alternatives. Among them, we focus on a subset that is shown in Table 6.2, which provide the most interesting combinations, in terms of experimental efficiency, of user and background distributions for each clarity model. These alternatives are further analysed in detail below (with examples) and in Section 6.5.1 where correlations obtained by each model are presented.

Qualitative observation

In order to illustrate the proposed prediction framework and give an intuitive idea of what user characteristics the predictors are capturing, we show the relevant aspects of specific users that result in clearly different predictor values, in a similar way to the examples provided in (Cronen-Townsend et al., 2002) for query clarity. We compare three user clarity models out of the seven models presented in Table 6.2: one for each formulation included in Table 6.1. In order to avoid distracting biases on the clarity scores that a too different number of ratings between users might cause, we have selected pairs of users with a similar number of ratings. This effect would be equivalent to that found in Information Retrieval between the query length and its clarity for some datasets (Hauff, 2010).

Table 6.3 shows the details of two sample users on which we will illustrate the effect of the predictors. As we may see in the table, u_2 has a higher clarity value than u_1 for the three models analysed. That is, according to our theory, u_2 is less “ambiguous” than u_1 . Figure 6.1 shows the clarity contribution in a term-by-term basis for one of the item-and-rating-based clarity models – where, in this case, terms are equivalent to a pair (rating, item) – as analysed in (Cronen-Townsend et al., 2002). In the figure, we plot $p(r|u, i) \log_2(p(r|u, i)/p(r|i))$ for the different terms in the collection, sorted in descending order of contribution to the user model, i.e.,

User	Number of ratings	ItemUser clarity	RatItem clarity	IRUserItem clarity
u_1	51	216.015	28.605	6.853
u_2	52	243.325	43.629	13.551

Table 6.3. Two example users, showing the number of ratings they have entered, and their performance prediction values for three user clarity models.

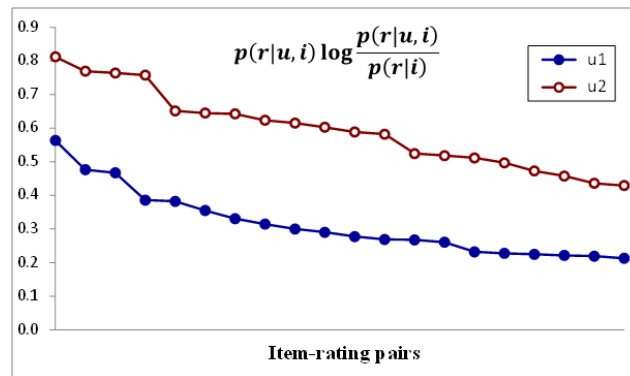


Figure 6.1. Term contributions for each user, ordered by their corresponding contribution to the user language model. IRUserItem clarity model.

$p(r|u, i)$, for each user. For the sake of clarity, only the top 20 contributions are plotted. We may see how the user with the smaller clarity value receives lower contribution values than the other user. This observation is somewhat straightforward since the clarity value, as presented in Equation (6.1), is simply the sum of all these contributions, over the set of terms conforming the vocabulary. In fact, the figures are analogous for the rest of the models, since one user always obtains higher clarity value than the other.

Let us now analyse more detailed aspects in the statistical behaviour of the users that explain their difference in clarity. The IRUserItem clarity model captures how differently a user rates an item with respect to the community. Take for instance the top item-rating pairs for users 1 and 2 in the above graphic. The top pair for u_2 is (4, “McHale’s Navy”). This means that the probability of u_2 rating this movie with 4 is much higher than the background probability (considering the whole user community) of this rating for this movie. Indeed, we may see that u_2 rated this movie with a 3, whereas the community mode rating is 1 – quite farther away from 4. This is the trend in a clear user. On the other extreme of the displayed values, the bottom term in the figure for u_1 is (2, “Donnie Brasco”), which is rated by this user with a 5, and the community mode rating for this item is 4, thus showing a very similar trend between both. This is the characteristic trend of a non-clear user.

Furthermore, if we compare the background model with the user model, we obtain more insights about how our models are discriminating distinctive from mainstream behaviour. This is depicted in Figure 6.2. In this situation, we select those terms which maximise the difference between the user and background models. Then, for this subset of the terms, we sort the vocabulary with respect to its collection probability, and then we plot the user probability model for each of the terms in the vocabulary.

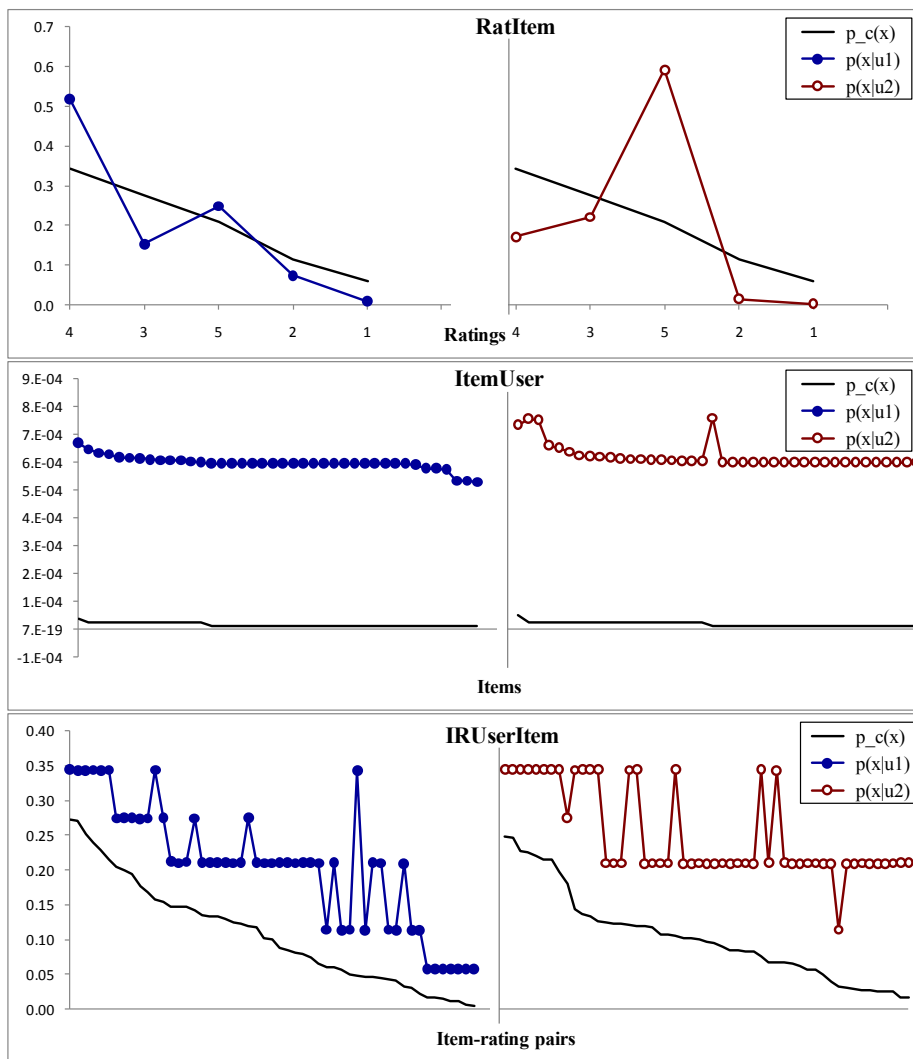


Figure 6.2. User language model sorted by collection probability.

These figures show how the most ambiguous user obtains a similar distribution to that of the background model, while the distribution of the less ambiguous user is more different. In the rating-based model this effect is clear, since the likelihood of not so popular rating values (i.e., a ‘5’) is larger for u_2 than for u_1 , and at the same time, the most popular rating value (a ‘4’) is much more likely for u_1 . The figure about the ItemUser model is less clear in this aspect, although two big spikes are observed for u_1 with respect to the collection distribution, which correspond with two unpopular movies: ‘Waiting for Guffman’ and ‘Cry, the beloved country’, both with a very low collection probability. Finally, the figure about the IRUserItem model successfully shows how u_2 has more spikes than u_1 , indicating a clear divergence from the background model; in fact, u_1 ’s distribution partially mimics that of the collection. In summary, the different models proposed are able to successfully sepa-

Item clarity	Vocabulary \mathcal{X} / Context θ	Item model	Background model	Formulation
Rating-based	Ratings / None	$p(r i)$	$p_c(r)$	$\sum_r p(r i) \log_2 \frac{p(r i)}{p(r)}$
User-based	Users / None	$p(u i)$	$p_c(u)$	$\sum_u p(u i) \log_2 \frac{p(u i)}{p(u)}$
User-and-rating-based	Ratings / Users	$p(r u, i)$	$p_{ml}(r u)$	$\sum_{r,u} p(u)p(r i, u) \log_2 \frac{p(r i, u)}{p(r u)}$

Table 6.4. Three possible item clarity formulations, depending on the interpretation of the vocabulary and context spaces.

rate information concerning the user and that from the collection, in order to infer whether a user is different or similar from the collection as a whole.

Item clarity

Alternatively to user-based predictors, we can also consider item-based predictors, where the performance prediction is made on an item-basis. Item predictors can be defined analogously as those defined previously for users, the equation for **item clarity** being as follows:

$$\text{clarity}(i) = \mathbb{E}_\theta \left[\sum_{x \in \mathcal{X}} p(x|i, \theta) \log_2 \frac{p(x|i, \theta)}{p(x|\theta)} \right] \quad (6.16)$$

The formulation of the item predictors we propose is basically equivalent to the user-based scheme but swapping users and items. That is, we have the three formulations presented in Table 6.4 where the vocabulary now may be either ratings or users, and the context variable is the user space. Based on these three formulations, and on derivations analogous to those presented before, we propose the seven item predictors defined in Table 6.5 which are further evaluated in Section 6.5.2.

In some of the instantiations of the item clarity predictor, we may observe that there are item probability models statistically equivalent to some of the user probability models, such as the $p_U(r|i, u)$ and $p_U(r|u, i)$. For this reason, we now only spec-

Item clarity name	Item dependent model	Background model
RatItem	$p_I(r i, u); p_{IR}(u i)$	$p_c(r)$
RatUser	$p_U(r i, u); p_{UR}(u i)$	$p_c(r)$
UserSimple	$p_R(u i)$	$p_c(u)$
UserItem	$p_{IR}(u i)$	$p_c(u)$
URItem	$p_I(r i, u)$	$p_{ml}(r u)$
URUser	$p_U(r i, u)$	$p_{ml}(r u)$
URItemUser	$p_{IU}(r i, u)$	$p_{ml}(r u)$

Table 6.5. Different item clarity models implemented.

ify those probability models which have not been defined before, for the rest of estimations see Equation (6.3):

$$p_R(u|i) = \sum_{r \in \mathcal{R}} p_{ml}(u|r)p_{ml}(r|i) \quad (6.17)$$

$$p_{IR}(u|i) = \sum_{r \in \mathcal{R}} p(u|i, r)p_{ml}(r|i) \quad (6.18)$$

$$p_{IU}(r|i, u) = p_{UI}(r|i, u) \quad (6.19)$$

6.2.2 Log-based clarity

In this section we adapt some of the previous models proposed for user clarity when the preference data come in the form of user-item interaction logs. Log data has a particularity we aim to exploit: the number of times a user consumes (purchased, listened, browsed, etc.) an item may be higher than one, in contrast with rating-based preferences, where the relation between a user and an item is summarised as a unique value, the rating. Moreover, the timestamp of the interactions has a stronger meaning in the implicit approach, as it informs of the very instant the user decided to use the item, rather than the time when the user decided to reflect on her quality of experience with the item (rating time). Specialised recommendation algorithms have been proposed in the literature that exploit such features in order to obtain better recommendations (Xiang et al., 2010; Lee et al., 2008). Additional alternatives for the definition of the vocabulary may be proposed, but we shall focus on these two: log co-occurrences and timestamps.

Specifically, based on Equation (6.2) and the three instantiations of \mathcal{X} and θ shown in Table 6.1, in principle only an instantiation analogous to the second one ($\mathcal{X} = \mathcal{I}$, no context – to which we shall refer as frequency-based clarity) makes sense here, as there is no rating space. However, it is possible to consider an additional space, which leads to structurally similar instantiations by taking time as the \mathcal{X} vocabulary. The similarity is only syntactic, as the meaning and implications of the resulting magnitude, to which we shall refer as time-based clarity, are quite different from rating-based clarity – in other words, ratings and time are quite different dimensions –, as we shall describe later below.

Frequency-based clarity

As mentioned above, we may define the following instantiation of the Equation (6.2) based on frequencies as follows:

$$\text{frequency-based clarity}(u) = \sum_i p(i|u) \log_2 \frac{p(i|u)}{p(i)} \quad (6.20)$$

where now the estimations of the user and background models are computed using directly the frequencies of the co-occurrences of some particular user-item interaction in the data:

$$p(i) = \frac{\text{freq}(i)}{\sum_{j \in \mathcal{J}} \text{freq}(j)}$$

$$p(i|u) = \frac{\text{freq}(i, u)}{\sum_{j \in \mathcal{J}_u} \text{freq}(j, u)}$$
(6.21)

An alternative to such estimations is to use transformations from implicit log-based to explicit ratings, such as the one proposed in (Celma, 2008). In that approach, any of the predictors based on ratings proposed in the previous section could be applied, since these transformations give the additional vocabulary space of ratings that was absent in principle in log data.

Time-based clarity

As introduced earlier, the second dimension susceptible to be exploited when log-based preference data are available is time. The time dimension is being paid increasing attention in Information Retrieval, where, for instance, it has been integrated into language models as a means to capture some temporal information needs from the user (Berberich et al., 2010), and the temporal query dynamics are being increasingly considered in the field (Kulkarni et al., 2011). In fact, temporal query features have also been used for query performance prediction, showing low or moderate correlation with query performance by themselves, although higher correlation is obtained when such features are combined with query clarity (Diaz, 2007; Diaz and Jones, 2004).

Furthermore, time has an inherent place in recommendation: recommender systems take as input (potentially long) histories of user interaction with items (Lathia, 2010; Zimdars et al., 2001; Burke, 2010). Time is an essential dimension in making sense of the data, and in explaining, analysing and interpreting the motivations behind the actions of users recorded over time. We propose to bring these ideas to recommender systems, in particular, to adapt the temporal features studied by Díaz and colleagues on a recommender system dataset. More specifically, we use the temporal Kullback-Leibler divergence described in (Diaz and Jones, 2004) as a starting point, which we generalise and elaborate upon by considering the instantiation of Equation (6.2) for a time-based space \mathcal{X} , and the space of items as a possible contextual dimension, as presented in Table 6.6. In the following, we define the specific instantiations of the temporal clarity formulations presented in this table.

User clarity	Vocabulary \mathcal{X} / Context θ	User model	Background model	Formulation
Time-based	Time / None	$p(t u)$	$p(t)$	$\sum_t p(t u) \log_2 \frac{p(t u)}{p(t)}$
Item-and-time-based	Time / Items	$p(t u, i)$	$p(t i)$	$\sum_{t,i} p(i)p(t u, i) \log_2 \frac{p(t u, i)}{p(t i)}$

Table 6.6. Two temporal user clarity formulations, depending on the interpretation of the vocabulary space.

Time based model. We denote as **TimeSimple clarity** the most direct adaptation for temporal clarity, which does not use any further extension over other dimensions. It simply computes $p(t|u)$ using smoothing (see below) and $p_c(t)$ from the collection frequencies.

Item-and-time based model. Like in the previous section, we develop conditional probabilities into sums with respect to a third variable: the items rated by the user. Here, we define two temporal clarity predictors depending on the distribution assumed for the items in the summation. If the distribution is uniform we denote such predictor as **ItemTime clarity** and $p(i) = 1/|\mathcal{I}|$. If, on the other hand, we also want to incorporate the popularity of the item for – which we have more data in this context and makes more sense than in rating data, since there the interaction between a user and an item is binary –, we include the prior item probability as $p(i) = p_c(i)$, which can be estimated considering the frequency by which i is accessed based on the interaction log.

The probabilities presented above are estimated as follows:

$$\begin{aligned}
 p_c(t) &= \frac{|\{(v, j, t) \in \mathcal{L} | v \in \mathcal{U}, j \in \mathcal{J}\}|}{|\mathcal{L}|} \\
 p_c(i) &= \frac{|\{(v, i, s) \in \mathcal{L} | v \in \mathcal{U}, s \in \mathcal{S}\}|}{|\mathcal{L}|} \\
 p_{mi}(t|i) &= \frac{|\{(v, i, t) \in \mathcal{L} | v \in \mathcal{U}\}|}{|\{(v, i, s) \in \mathcal{L} | v \in \mathcal{U}, s \in \mathcal{S}\}|} \\
 p_{mi}(t|u) &= \frac{|\{(u, j, t) \in \mathcal{L} | j \in \mathcal{J}\}|}{|\{(u, j, s) \in \mathcal{L} | j \in \mathcal{J}, s \in \mathcal{S}\}|} \\
 p_{mi}(t|u, i) &= \frac{|\{(u, i, t) \in \mathcal{L}\}|}{|\{(u, i, s) \in \mathcal{L} | s \in \mathcal{S}\}|}
 \end{aligned} \tag{6.22}$$

Note that the variable t in (u, i, t) in the above expressions denotes a timestamp in the discretised time segment (e.g. day, week) represented by t . Furthermore, these are simple estimations of the distributions; hence, it is also possible to introduce non-parametric estimations or additional expansions through similar users or items (Wang et al., 2006a; Wang et al., 2008a). Moreover, distributions can also be modeled by

other statistical theories or hypothesis (such as Bayesian inversion), and distribution fitting/modelling from time series theory could also be studied (Diaz and Jones, 2004; Wang et al., 2008b).

In particular, we have smoothed these estimations using Jelinek-Mercer as follows:

$$\begin{aligned} p(t|i) &= \lambda p_{ml}(t|i) + (1 - \lambda) p_c(t) \\ p(t|u) &= \lambda p_{ml}(t|u) + (1 - \lambda) p_c(t) \\ p(t|u, i) &= \lambda p_{ml}(t|u, i) + (1 - \lambda) p_c(t) \end{aligned} \tag{6.23}$$

6.3 Predictors based on social topology

Social information is widespread nowadays. As we surveyed in Chapter 2, recommender systems that use social information are proliferating in the research literature, as well as in the recommender system industry, because of the effectiveness they are being found to have. It seems therefore sensible to consider social information as a potentially useful input for predicting the performance of recommendation. The motivation for this approach is obvious when applied to social recommender systems, though we will also explore its potential properties in relation to non-explicitly social recommendation, in order to study whether social topologies may have an indirect effect on the results of the algorithms for different users.

With this goal in mind, we explore the use of graph-based measures as indicators of the user strength in the social network, which may in turn correlate with the ease or difficulty of users as recommendation targets. Graph-based measures developed from link-analysis theory are straightforward to interpret where they are often used to understand the structure of communities within a population (De Choudhury et al., 2010; Albert and Barabási, 2002). As a basis for user performance prediction they may thus bring an advantage in terms of explaining the predictions.

More specifically, the utilised indicators of the user strength in the network are based on the following vertex measures computed over the social network for each user, where a user is represented as a node in the graph, and the user's friends correspond with the node's neighbours:

- **Average neighbour degree:** mean number of friends of each user's friend (Kossinets and Watts, 2006).
- **Betweenness centrality:** indicator of whether a user can reach others on relative short paths (Freeman, 1977).
- **Clustering coefficient:** probability that the user's friends are friends themselves (Watts and Strogatz, 1998).

- **Degree**: number of the user's friends in the social network (Milgram, 1967).
- **Ego components size**: number of connected components remaining when the user and her friends are removed (Newman, 2003).
- **HITS**: Kleinberg, 1999) defines two complementary measures which assign recursively a weight to each vertex (user) depending on the topology of the network. In this way, they define hubs and authorities: a vertex is a hub when it links to authoritative vertices, and is an authority when it links to hub vertices. Since the social network used here (see Appendix A.1.3) is undirected, hub and authority scores are redundant and we only report one, denoted as **HITS**.
- **PageRank** score: well-known measure of connectivity relevance within a social network based on a random walk over the vertices, where a probability ($\alpha = 0.2$ in our experiments) of jumping to any other vertex is introduced (Brin and Page, 1998).
- **Two-hop neighbourhood** size: count of all the user's friends plus all the user's friend's friends (De Choudhury et al., 2010).

6.4 Other approaches

As a reference for comparison, we shall also test further predictors besides the ones proposed in the thesis, directly drawn from the literature, and not necessarily based on probabilistic formalisations, but following more loose formalisations, or heuristic approaches. As a further sanity check, we shall also examine obvious and simple functions (such as the amount of activity of a user), as a reference for the justification of elaborate approaches as proposed. Next, we present these predictors which are evaluated and compared in Section 6.5.

6.4.1 Using rating-based preference data

A fairly simple user predictor against which we would like to compare more elaborate functions is the **count** predictor, namely the number of items a user has rated at some specific moment. This predictor, as we shall see later, can be defined in the training set and in the test set, and although its rationale is the same, the output has different implications. Whereas in training this predictor is measuring how much information a recommender knows about some specific user, in test this value would be related to the amount of relevance used to obtain the performance metric. Furthermore, as observed in Chapter 4, the amount of relevance would be different depending on the evaluation methodology considered. However, we have to note that, due to statistical effects, the training count (profile size in training) and test count

(profile size in test) would probably be related if the training/test split is performed randomly.

$$\text{count}(u) = |\mathcal{J}_u| = |\{(u, \cdot)\}| \quad (6.24)$$

Two additional heuristic predictors can be defined by looking at user statistics such as the **mean** and the **standard deviation** of the user's ratings. It seems plausible that such predictors would not be equally powerful for any type of recommender: it would depend on whether these statistics are used by the recommender. For instance, one might have the intuition that the higher the standard deviation, the lower the recommendation performance as one may figure out uniform user ratings to be a somewhat easier target.

$$\text{mean}(u) = \frac{1}{|\mathcal{J}_u|} \sum_{i \in \mathcal{J}_u} r(u, i) \quad (6.25)$$

$$\text{std}(u) = \sqrt{\frac{1}{|\mathcal{J}_u|} \sum_{i \in \mathcal{J}_u} (r(u, i) - \text{mean}(u))^2} \quad (6.26)$$

Alternatively to these heuristic predictors, we have also experimented with a predictor defined upon the past observed recommender's performance. In this way, this predictor – denoted as **training performance** from now on – use a validation set (as a subset of the original training set) to evaluate the performance of each user with respect to a specific recommender; then, this value is the one returned by the predictor at test time. This approach is inspired in the Machine Learning techniques which aim to learn a feature (in this case, the user's performance) by using some training information. For this predictor, this training information is the performance computed on the validation set.

Additionally, we propose to measure the **entropy** of the user's preferences as a quantification of the uncertainty associated with a probability distribution (Cover and Thomas, 1991). We may therefore assess the uncertainty involved in the system's knowledge about a user's preferences by the entropy of the item distribution (the probability to choose an item) given the information in the user profile, using the ground models presented in Section 6.2.1. Hence, we define this predictor as follows:

$$\text{entropy}(u) = \sum_{i \in \mathcal{J}_u} p(i|u) \log_2 p(i|u) \quad (6.27)$$

Alternative measures from Information Theory could be used to define user-based predictors, like Information Gain (Bellogín, 2009), but we leave them out of this analysis because its application to Recommender Systems is neither clear nor principled and their predictive results are not optimal. Furthermore, other measures already proposed in the literature such as inverse user frequency (Breese et al., 1998)

and the analogous inverse item frequency (Bellogín, 2009), and other manipulations of the same concept, are also ignored here because they are simply transformations of the previously presented count predictor. Finally, the concept of power users (Lathia et al., 2008) could also be used as a proxy for well-performing users, but preliminary results have not shown strong predictive power.

6.4.2 Using log-based preference data

As we have observed in the previous section, recommendation performance usually has obvious predictors, obvious in the sense that they do not involve any interesting finding or insightful kind of analysis, or anything to learn from. We include in our analysis some of these obvious predictors, framed as baseline performance predictors that basically count how many interactions a user has had with the system. In this sense, these predictors are slightly different to the ones presented in the previous section, namely because in log-based datasets repetitions of items are allowed in a user's profile. In order to account for this difference, apart from **count**, **mean**, and **standard deviation** predictors, we propose to normalise the count predictor by the number of items consumed by each user, that is, we define the **average count** predictor as follows:

$$\text{average count}(u) = \frac{|\{(u, j, s) \in \mathcal{L} | j \in \mathcal{I}, s \in \mathcal{S}\}|}{|\{j \in \mathcal{L} : (u, j, s) \in \mathcal{L}, s \in \mathcal{S}\}|} \quad (6.28)$$

We also test more elaborate predictors based on the temporal dimension, such as the ones defined in (Diaz and Jones, 2004). First-order autocorrelation (or temporal self-correlation) can be considered with a reinterpretation of the random variables. Specifically, this predictor, in contrast with the temporal Kullback-Leibler divergence where the similarity with the temporal background model is assessed, captures the structure of the query time series. For instance, a uniform distribution would have an autocorrelation value of 0, whereas a query time series with strong inter-day (or whatever segment size is used to build the discrete time series) dependency will obtain a high autocorrelation value.

Thus, we define the **autocorrelation** user predictor as follows:

$$\text{autocorrelation}(u) = \frac{\sum_{t=1}^T (p(t|u) - 1/T)(p(t+1|u) - 1/T)}{\sum_{t=1}^T (p(t|u) - 1/T)^2} \quad (6.29)$$

where T is the total number of time units in the time interval. We can observe how this predictor captures the similarity between two consecutive observations.

Extensions of this predictor could use the probabilities defined in Section 6.2.2, like $p(t|u, i)$, instead of $p(t|u)$. Similarly, other predictors proposed by Díaz and Jones in (Diaz and Jones, 2004) and (Jones and Diaz, 2007) such as the kurtosis or

the burst model could be adapted to recommender systems, but we leave such extensions for future work.

6.5 Experimental results

In this section we provide correlation results where all the predictors – heuristic, social, and clarity-based – are compared against each other using an array of recommendation methods and evaluation methodologies.

6.5.1 User predictors using rating-based preference data

In this section we compare the correlations obtained for the clarity-based predictors defined in Section 6.2.1, the user entropy defined in Equation (6.27), and the baselines presented in Equations (6.24), (6.25), and (6.26) using the MovieLens 1M dataset. The λ parameter for the language model smoothing was not optimised for this task and a default value of 0.6 was used in all the models as originally used in (Cronen-Townsend et al., 2002). Here, we focus on Pearson’s correlation and P@10. Additional results are reported in Appendix A.4.1.

Table 6.7 shows the correlation values when the AR methodology is used. We can observe fairly high correlation values for recommenders pLSA, ItemPop, TFL2, and kNN, comparable to results in the query performance literature. A slightly lower correlation is found for TFL1, whereas no correlation is found for CB and IB. These results are consistent when other performance metrics are used such as nDCG, and at different cutoff points. Spearman’s correlation yields similar values. Here we also include the count predictor in test, which is obviously not a predictor in strict sense,

Predictor	Random	CB	IB	ItemPop	kNN	pLSA	TFL1	TFL2	Median	Mean
Count (training)	0.135	0.164	0.042	0.512	0.424	0.442	0.198	0.644	0.311	0.320
Count (test)	0.135	0.172	0.042	0.520	0.431	0.452	0.200	0.647	0.316	0.325
Training performance	0.024	0.176	0.258	0.429	0.296	0.357	0.215	0.485	0.277	0.280
Mean	0.019	0.067	-0.002	0.015	0.022	0.108	0.026	-0.018	0.021	0.030
Standard deviation	0.008	0.008	0.011	-0.029	-0.031	-0.032	0.011	-0.051	-0.011	-0.013
ItemSimple Clarity	0.149	0.191	0.046	0.549	0.453	0.489	0.222	0.683	0.338	0.348
ItemUser Clarity	0.134	0.166	0.048	0.493	0.416	0.428	0.215	0.634	0.316	0.317
RatUser Clarity	0.135	0.160	0.048	0.514	0.442	0.435	0.214	0.651	0.325	0.325
RatItem Clarity	0.127	0.159	0.039	0.475	0.402	0.405	0.203	0.611	0.303	0.303
IRUser Clarity	0.128	0.157	0.027	0.486	0.382	0.408	0.182	0.599	0.282	0.296
IRItem Clarity	0.122	0.165	0.034	0.446	0.352	0.386	0.188	0.551	0.270	0.281
IRUserItem Clarity	0.128	0.158	0.033	0.479	0.379	0.403	0.193	0.594	0.286	0.296
Entropy	0.121	0.168	0.025	0.492	0.389	0.483	0.140	0.589	0.279	0.301
Median	0.128	0.162	0.037	0.489	0.396	0.418	0.196	0.605		
Mean	0.112	0.145	0.033	0.413	0.338	0.367	0.166	0.511		

Table 6.7. Pearson’s correlation between rating-based user predictors and P@10 for different recommenders using the AR methodology (MovieLens dataset).

since it uses a different input than the other predictors, but we include it in our analysis as a further reference to check behaviours.

As mentioned in Chapter 5, the standard procedure in Information Retrieval for this kind of evaluation is to compute correlations between the predictor(s) and one retrieval model (like in (Cronen-Townsend et al., 2002) and (Hauff et al., 2008a)) or an average of several methods (Mothe and Tanguy, 2005). This approach may hide the correlation effect for some recommenders, as we may observe from the median and mean correlation values included in the table, which are still very large despite the fact that two of the recommenders analysed have much lower correlations. Nonetheless, these aggregated values, i.e., the mean and the median, provide competitive correlation values when compared with those in the literature.

The difference in correlation for CB and IB recommenders may be explained considering two factors: the actual recommender performance and the input sources used by the recommender. With regards to the first factor, as presented in Table 6.8, the IB algorithm performs poorly (in terms of the considered ranking quality metrics, such as precision and nDCG) in comparison to the rest of recommenders. It seems natural that a good predictor for a well performing algorithm (specifically, pLSA is the best performing recommender in this context) would hardly correlate at the same time with a poorly performing one.

This does not explain however the somewhat lower correlation with the content-based recommender, which has better performance than TFL1. The input information that this recommender and the predictors take in are very different: the latter compute probability distributions based on ratings given by users to items, while the former uses content features from items, such as directors and genres. Furthermore, the CB recommender is not coherent with the inherent probabilistic models described by the predictors, since the events modeled by each of them are different: CB would be related to the likelihood that an item is described by the same features as those items preferred by the user, whereas predictors are related to the probability that an item is rated by a user. Moreover, the predictors' ground models coherently fit in the standard collaborative framework (Wang et al., 2008a), which reinforces the suitability of the user performance predictors presented herein, at least for collaborative filtering recommenders.

It is worth noting to this respect that most clarity-based query performance predic-

Recommender	Random	CB	IB	ItemPop	kNN	pLSA	TFL1	TFL2
AR methodology	0.0025	0.0163	0.0001	0.0897	0.0307	0.1454	0.0024	0.0696
1R methodology	0.0099	0.0221	0.0074	0.0649	0.0437	0.0836	0.0221	0.0690
UIR methodology	0.0100	0.0223	0.0068	0.0406	0.0381	0.0718	0.0294	0.0524
PIR methodology	0.0101	0.0197	0.0208	0.0282	0.0265	0.0604	0.0203	0.0348

Table 6.8. Summary of recommender performance using different evaluation methodologies (evaluation metric is P@10 with the MovieLens dataset).

tion methods in Information Retrieval study their predictive power on language modeling retrieval systems (Cronen-Townsend et al., 2002; Hauff et al., 2008a; Zhou and Croft, 2007) or similar approaches (He and Ounis, 2004). This suggests that a well performing predictor should be defined upon common spaces, models, and estimation techniques as the retrieval system the performance of which is meant to be predicted.

Finally, the correlation values found by the training performance predictors, although sometimes strong, are not as high as those of the baselines predictors – such as training count – in most situations, in particular, they are always lower except for the IB and TFL1 recommenders. This highlights the importance of having a more general model for predicting the performance of a user, since these predictors in fact depend considerably on the properties of the validation (and test) partition of the data, such as the amount of sparsity, type of items evaluated and so on.

Unbiased performance prediction

In Chapter 4 we already demonstrated that some methodologies may be biased towards more popular items or sparsity constraints. We can observe in the previous table that trivial predictors such as count (either in training or in test) obtain significant (and positive) correlation, no matter the recommender. We argue whether this is because these predictors are really capturing an interesting effect or the evaluation methodology is prone to such effect. In order to overcome this problem, now we present the same correlation analysis but with the different methodologies presented in Chapter 4.

In Table 6.9 we show results with the methodology 1R. Here we can observe that most of the correlation values are lower than in the previous case; interestingly, the correlation with the Random recommender now is almost 0 for every predictor (and in particular, for the training and test profile size). This is evidence that per-

Predictor	Random	CB	IB	ItemPop	kNN	pLSA	TFL1	TFL2
Count (training)	0.061	-0.038	0.092	0.258	0.108	0.303	0.086	0.394
Count (test)	0.063	-0.033	0.091	0.266	0.115	0.312	0.089	0.398
Training performance	0.012	0.332	0.168	0.272	0.266	0.133	0.303	0.240
Mean	0.036	0.082	-0.029	0.028	0.111	0.117	0.145	0.031
Standard deviation	-0.010	0.006	0.051	-0.060	-0.116	-0.080	-0.040	-0.114
ItemSimple Clarity	0.066	-0.033	0.094	0.265	0.115	0.322	0.105	0.409
ItemUser Clarity	0.059	-0.038	0.087	0.236	0.100	0.287	0.096	0.375
RatUser Clarity	0.057	-0.054	0.083	0.245	0.130	0.285	0.086	0.372
RatItem Clarity	0.057	-0.044	0.069	0.225	0.110	0.268	0.094	0.352
IRUser Clarity	0.056	-0.020	0.053	0.250	0.069	0.280	0.077	0.364
IRItem Clarity	0.051	-0.010	0.058	0.205	0.029	0.235	0.074	0.310
IRUserItem Clarity	0.056	-0.020	0.052	0.242	0.066	0.273	0.081	0.357
Entropy	0.091	0.021	0.144	0.354	0.169	0.460	0.114	0.543

Table 6.9. Pearson’s correlation between rating-based user predictors and P@10 for different recommenders using the 1R methodology (MovieLens dataset).

Predictor	Random	CB	IB	ItemPop	kNN	pLSA	TFL1	TFL2
Count (training)	0.048	-0.012	0.237	0.162	0.115	0.140	0.022	0.235
Count (test)	0.049	-0.001	0.226	0.135	0.110	0.137	0.036	0.213
Mean	0.023	0.051	-0.035	0.009	0.108	0.075	0.155	-0.006
Standard deviation	0.015	0.032	0.023	-0.047	-0.098	-0.038	-0.061	-0.049
ItemSimple Clarity	0.055	-0.005	0.241	0.166	0.128	0.153	0.042	0.241
ItemUser Clarity	0.046	-0.009	0.232	0.142	0.109	0.133	0.028	0.216
RatUser Clarity	0.045	-0.028	0.234	0.155	0.137	0.130	0.022	0.225
RatItem Clarity	0.043	-0.025	0.212	0.136	0.119	0.117	0.033	0.203
IRUser Clarity	0.044	0.002	0.180	0.153	0.069	0.134	0.029	0.210
IRItem Clarity	0.036	0.011	0.178	0.114	0.035	0.108	0.014	0.173
IRUserItem Clarity	0.042	0.003	0.178	0.147	0.065	0.130	0.028	0.203
Entropy	0.078	0.044	0.278	0.227	0.169	0.249	0.073	0.321

Table 6.10. Pearson’s correlation between rating-based user predictors and P@10 for different recommenders using the U1R methodology (MovieLens dataset).

formance results using the AR methodology are higher for users with more items in their test, independently from the recommendation algorithm complexity (see correlations with Random recommender in Table 6.7). In the same way, the U1R (Table 6.10) and P1R (Table 6.11) methodologies also obtain negligible correlation values for the Random recommender, which confirms the suitability of these methodologies for our purposes. We also have to note that we have not applied the training performance predictor in these methodologies because their restrictions do not let to replicate the same conditions in a validation split. Furthermore, as stated in Chapter 4, both approaches aim to remove the bias towards more popular items. Here, we can observe how the correlation with respect to the ItemPop recommender is comparable to that with the Random recommender with the P1R methodology, confirming again the ability of this methodology to produce unbiased results (at least, with respect to popular items).

The main difference in the results obtained between these three methodologies (1R, U1R, and P1R) seems to be more at the recommender level rather than at the

Predictor	Random	CB	IB	ItemPop	kNN	pLSA	TFL1	TFL2
Count (training)	0.073	-0.005	0.253	0.088	0.103	0.160	-0.001	0.307
Count (test)	0.076	0.000	0.253	0.093	0.108	0.168	0.003	0.308
Mean	0.034	0.073	-0.033	0.008	0.110	0.085	0.188	-0.026
Standard deviation	-0.010	0.009	0.014	-0.058	-0.104	-0.044	-0.061	-0.051
ItemSimple Clarity	0.078	0.000	0.254	0.084	0.111	0.169	0.019	0.313
ItemUser Clarity	0.072	-0.001	0.249	0.075	0.101	0.156	0.005	0.303
RatUser Clarity	0.071	-0.016	0.252	0.086	0.128	0.148	0.003	0.297
RatItem Clarity	0.067	-0.011	0.234	0.077	0.113	0.138	0.016	0.288
IRUser Clarity	0.066	0.002	0.200	0.086	0.066	0.147	0.006	0.274
IRItem Clarity	0.059	0.010	0.192	0.061	0.037	0.123	-0.006	0.242
IRUserItem Clarity	0.066	0.003	0.200	0.082	0.065	0.145	0.006	0.272
Entropy	0.092	0.038	0.286	0.133	0.128	0.266	0.039	0.379

Table 6.11. Pearson’s correlation between rating-based user predictors and P@10 for different recommenders using the P1R methodology (MovieLens dataset).

predictor level, in the sense that the trend in predictor effectiveness is similar for each methodology but the correlations obtained for each recommender vary dramatically from one methodology to another. For instance, IB recommender obtains near zero correlations with 1R but higher (significant) values for U1R and P1R; a similar situation occurs with the TFL2 recommender, where the correlations are lower for the U1R methodology and higher for 1R and P1R. Note that the training and test sets are the same for all the methodologies except for U1R, which means that the performance predictors are entirely new for that methodology. Thus, a priori it would not be clear that such an agreement between the different methodologies should appear at the predictor level unless they are really capturing the same nuance about the user, no matter the evaluation methodology used.

It is worth noting that the correlation values of these three methodologies have been found after a careful examination of the available data, where two different trends emerged: one where the performance values were more or less uniformly distributed in the interval $[0,0.1]$ – recall that 0.1 is the maximum value for the metric $P@10$ with the 1R methodology, since there is only one relevant item – ; and a second one where a fixed value was obtained. This second trend, against which our predictors shown no correlation at all (since the performance had a zero standard deviation, and thus the correlation was impossible to calculate) is able to degrade the correlation coefficient almost to negligible values, mainly because it accounts for half of the number of points. This problem with correlation coefficients, and with Pearson’s correlation in particular, is well known in the literature of performance prediction (Hauff, 2010; Pérez Iglesias, 2012). For this reason, we have divided the performance values and computed two correlations in order to account for these two trends: the values with respect to the first trend are those presented in the previous tables, whereas the correlation with respect to the second trend was not computable because the variable had a zero standard deviation.

In summary, there seems to be no clear winner among the set of performance predictors proposed. The predictive power of each of them is clearly influenced by the actual recommender its performance aims to be predicted and the evaluation methodology in use. Nonetheless, **the proposed predictors usually obtain higher correlation values than baseline predictors** such as the mean or the standard deviation, evidencing their predictive power **independently from the evaluation methodology**. Surprisingly, the ItemSimple clarity predictor obtains very good results in most of the situations, although more complex predictors like IRUser or IRUserItem clarity obtain stronger correlations for some recommenders.

6.5.2 Item predictors using rating-based preference data

In the same way we have assessed the predictive power of user predictors, we now aim to estimate the predictive power of item predictors. However, the true perform-

	u_1	u_2	u_3		i_1	i_2	i_3	i_4
	i_1 0.8	i_2 0.6	i_3 0.9		* u_1 0.8	u_1 0.7	u_3 0.9	* u_3 0.6
	* i_2 0.7	* i_3 0.5	* i_4 0.6		* u_3 0.5	u_2 0.6	* u_1 0.6	u_1 0.5
	* i_3 0.6	* i_4 0.4	* i_1 0.5		u_2 0.3	u_3 0.1	* u_2 0.5	* u_2 0.4
	i_4 0.5	i_1 0.3	i_2 0.1					
P@2	0.5	0.5	0.5		1.0	0.0	0.5	0.5

Table 6.12. Procedure to obtain ranking for items from user rankings generated by a standard recommender. * denotes a relevant item, and the numbers are the score predicted by the recommendation method.

ance value for an item is not straightforward to compute, since the process has to produce unbiased results in the space of items (as described in Chapter 4) but with the characteristic that the item dimension is not the main input of the recommendation process, and thus, some new approach has to be put in place.

There are basically two possibilities for computing the true performance on an item: either starting from the results obtained using a standard procedure (obtain a ranking for each user by recommending items to users), then transposing users and items (generating, thus, user rankings for each item) and computing the per-ranking performance as usual; or transpose the original rating matrix in order to effectively “recommend users” for each item. This would implicitly imply a transposition of the recommendation task, which may also make sense: find the most suitable users to recommend an item – this would be the scenario, e.g. in advertisement targeting when a new product is released on the market. Here, we use the former approach since the latter does not produce consistent results in our experiments, probably because the recommendation problem is not completely symmetric and, thus, this method is not able to properly capture the recommender’s performance for each item. On the other hand, non-personalised recommenders (such as recommendation by item popularity) cannot be applied in the symmetric formulation: since the same item ranking is built for all users, the user ranking for an item would be a global tie on all users. Table 6.12 shows an example of how we may transpose users and items from an item ranking for three users. We show that the precision for all the users is the same, whereas for the items is completely diverse, ranging from zero to perfect precision.

In our experiments, we have tested the different methodologies already presented along with a modified version of the U1R evaluation methodology (user-uniform U1R, or uuU1R). The rationale for the uuU1R design goes as follows: in the U1R methodology we force the same number of ratings (or, equivalently, users) for the items in the test set, however, users are freely assigned to each item. Now, when we transpose users and items this situation may produce a new problem, since there

Predictor	Random	CB	ItemPop	kNN	pLSA
Count (training)	0.414	0.060	-0.151	-0.021	-0.269
Count (test)	-----	-----	-----	-----	-----
Mean	0.602	0.125	0.096	0.040	-0.038
Standard deviation	-0.313	0.025	-0.006	-0.003	0.075
UserSimple Clarity	0.467	0.080	-0.120	-0.015	-0.240
UserItem Clarity	0.419	0.064	-0.145	-0.018	-0.261
RatItem Clarity	0.440	0.075	-0.127	-0.015	-0.230
RatUser Clarity	0.451	0.085	-0.103	-0.004	-0.201
URItem Clarity	0.396	0.053	-0.174	-0.026	-0.289
URUser Clarity	0.408	0.072	-0.132	-0.004	-0.243
URItemUser Clarity	0.409	0.061	-0.161	-0.021	-0.277
Entropy	0.381	-0.001	-0.216	-0.055	-0.442

Table 6.13. Pearson’s correlation for rating-based item predictors and precision using the uuU1R methodology (MovieLens dataset).

could be users assigned to more items which would bias the ranking’s performance towards items contained in the test set of heavy raters. Therefore, if we impose a uniform distribution also on the user’s dimension, this bias should decrease. We refer to the reader to Appendix A.3 for more details.

However, despite these efforts, we have not found a reliable methodology to evaluate the item performance. We present in Table 6.13 the results using the uuU1R methodology and the predictors defined in Table 6.5 for the precision metric. Recall that, since we transpose users and items from the generated rankings, to obtain a similar measure of $P@10$ we only use the top 10 items from each original ranking and then compute precision over the whole ranking for each item. We may observe in the table that the correlations with the Random recommender are very strong, questioning the validity of such results. Besides, the entropy predictor obtains stronger correlation than clarity-based in this case, and most of them (except for URItem) show little difference to training count. Note that it is not possible to compute a correlation with the test count predictor since that predictor has a constant value with zero standard deviation (see Equation (5.11) for more details on Pearson’s correlation) since every item has the same number of ratings in the test set in the uuU1R methodology.

As a conclusion, we have found that **a proper evaluation of item performance is not obvious**, mainly because the task of suggesting users to items is not completely symmetric with respect to the standard task of recommendation. We have devised different methodologies to estimate the recommendatoin performance of an item, however the difficulty lies mainly in forming consistent lists of “recommended” users for items, a difficulty which is not conceptual (ranking target users to whom an item may be recommended does make sense as a task in many scenarios), but technical (obtaining balanced result lists that allow for undistorted performance measurements).

6.5.3 User predictors using log-based preference data

In this section we analyse the correlation obtained between the predictors defined in Sections 6.2.2 and 6.4.2 and five recommenders using the 1R methodology on two versions of the Last.fm dataset – one where a temporal partition is performed and another where the partition is randomly made (more details about the splits in Appendix A.1.2). No smoothing was used in the language models since preliminary tests obtained better results with lower values of λ . Besides, for comparison purposes, we also include one of the clarity models proposed for rating-based preference data using the transformation proposed in Section 6.2.2 to use such predictors with log data along with the frequency-based clarity proposed in Equation (6.20). Like in the previous section, Pearson’s correlation with the P@10 evaluation metric is reported; for additional metrics, see Appendix A.4.2.

First, we can observe in Table 6.14 (temporal split) that ItemPriorTime clarity obtains strong correlation values, especially for the ItemPop and kNN recommenders. It is interesting to compare the correlations between this predictor and the ItemTime clarity, which are much lower. This is probably because the ItemPriorTime clarity predictor, as opposed to ItemTime clarity, incorporates a component that measures the item popularity, i.e., $p(i)$. The TimeSimple and the frequency-based clarity predictors, on the other hand, obtain strong correlation but negative values for all the recommenders except the ItemPop for the TimeSimple predictor. Furthermore, the ItemSimple clarity (a predictor based on explicit information) obtains negligible correlations except for the ItemPop and kNN recommenders.

Table 6.15, on the other hand, shows the results when a random split is used. We have to note that such split does not preserve the temporal continuity of the user’s preferences, and thus, any recommender or technique which makes use of temporal features is not guaranteed to succeed. Here, we can observe that TimeSimple predictor obtains strong correlations for all the recommenders except for the Random

Predictor	Random	CB	ItemPop	kNN	pLSA
Average Count	0.027	0.138	0.069	-0.013	0.191
Count	0.046	0.118	-0.058	0.131	0.139
Mean	-0.079	-0.361	0.054	-0.110	-0.376
Standard deviation	-0.050	-0.158	0.082	-0.132	-0.177
Autocorrelation	0.004	0.139	-0.066	-0.105	0.100
TimeSimple Clarity	-0.091	-0.342	0.093	-0.317	-0.354
ItemTime Clarity	0.037	0.078	0.038	0.258	0.064
ItemPriorTime Clarity	0.057	0.154	0.189	0.307	0.154
Frequency-based Clarity	-0.049	-0.410	-0.221	-0.291	-0.376
ItemSimple Clarity	0.027	0.047	-0.107	0.221	0.029

Table 6.14. Pearson’s correlation between log-based predictors and P@10 for different recommenders using 1R methodology (Last.fm temporal dataset).

Predictor	Random	CB	ItemPop	kNN	pLSA
Average Count	-0.023	-0.068	-0.170	-0.018	-0.087
Count	-0.012	-0.236	-0.242	-0.086	-0.198
Mean	0.036	0.182	0.100	0.047	0.118
Standard deviation	-0.009	0.089	0.079	0.092	0.082
Autocorrelation	0.045	-0.069	-0.089	-0.012	-0.055
TimeSimple Clarity	0.031	0.274	0.314	0.169	0.240
ItemTime Clarity	0.021	-0.145	0.004	0.025	-0.053
ItemPriorTime Clarity	0.011	-0.057	0.176	0.145	0.083
Frequency-based Clarity	0.025	0.018	-0.287	-0.182	-0.220
ItemSimple Clarity	0.020	-0.247	-0.163	-0.068	-0.186

Table 6.15. Pearson’s correlation between log-based predictors and P@10 for different recommenders using 1R methodology (Last.fm five-fold dataset).

technique. Like before, ItemPriorTime has a high correlation with the ItemPop recommender. In contrast with the previous situation, the ItemSimple clarity obtains strong but negative correlations for the personalised recommenders. Besides, the frequency-based clarity has negative correlations for all the recommenders except CB, a consistent situation with the results obtained with the temporal split.

Hence, we may conclude that **log-based and time-aware predictors successfully predict the performance of the recommendation algorithms**, although in some situations the sign of the prediction is negative. Moreover, frequency-based, ItemSimple, and TimeSimple clarity obtain consistently strong correlations both in a temporal split and in a random split of the data, evidencing their predictive power.

6.5.4 User predictors using social-based preference data

In this section we study the correlation between the predictors described in Section 6.3 and several recommenders using the two versions of the CAMRa dataset: social and collaborative. In this case, we also consider social filtering recommenders in order to analyse whether these predictors are sensitive to the source of information used by the recommender, and thus, whether they obtain stronger correlations with social filtering recommenders. Besides, one clarity-based predictor (ItemSimple) and the baseline rating predictors presented in Section 6.4.1 are also included in the analysis for comparison purposes. Additionally, for the HITS and PageRank graph metrics in this experiment we use the implementation developed in the JUNG library (O’Madadhain et al., 2003).

Table 6.16 shows correlation values obtained when using the AR methodology in the social version of the dataset. Here, we can observe that most of the correlation values obtained for the social predictors are negative, representing that the lower the predictor output, the better the performance, which may seem a little counter-intuitive, at least for the social filtering recommenders (Personal and PureSocial).

Among the social-based predictors, degree and two-hop neighbourhood size obtain better correlations than the rest.

A similar situation is presented in Table 6.17, where the collaborative-social version of the dataset is used. Again, most of the correlations with the social-based predictors are negative, and degree and two-hop neighbourhood size obtain higher correlations (in absolute value). Interestingly, in this situation strong correlations are obtained with the user-based recommender (kNN), in particular with degree and the average neighbour degree predictors. Nonetheless, these correlations are lower than those obtained for the ItemSimple predictor with the collaborative filtering recommenders. At the same time, this predictor always obtains worse correlations (in absolute value) than the social-based predictors for the social filtering recommenders, as expected.

Additionally, note that the number of points used in the correlation computation is different in each version of the dataset, namely: in the collaborative-social version

Predictor	Random	ItemPop	kNN	pLSA	Personal	PureSocial
Count (training)	0.032	0.122	0.113	0.031	0.062	0.111
Count (test)	0.158	0.252	0.382	0.167	0.235	0.174
Mean	-0.066	0.033	-0.012	0.023	-0.057	-0.051
Standard deviation	0.034	0.054	-0.020	0.115	0.128	0.183
Avg neighbour degree	-0.062	-0.089	-0.013	0.011	-0.074	-0.106
Betweenness centrality	-0.031	-0.016	0.027	-0.038	-0.012	-0.079
Clustering coefficient	0.049	-0.084	-0.023	0.048	-0.027	-0.035
Degree	-0.038	-0.046	0.015	-0.059	-0.147	-0.133
Ego components size	-0.058	0.005	0.004	-0.046	-0.056	-0.020
HITS	-0.021	-0.043	0.005	0.061	0.038	0.000
PageRank	-0.022	-0.025	-0.023	-0.039	-0.102	-0.037
Two-hop neighbourhood	-0.080	-0.082	0.004	-0.054	-0.123	-0.136
ItemSimple Clarity	0.030	0.157	0.130	0.050	0.072	0.126

Table 6.16. Pearson's correlation between social-based predictors and P@10 for different recommenders using AR methodology (CAMRa Social).

Predictor	Random	ItemPop	kNN	pLSA	Personal	PureSocial
Count (training)	0.012	0.098	0.203	0.107	0.058	0.111
Count (test)	0.096	0.207	0.389	0.179	0.232	0.170
Mean	-0.067	0.000	-0.126	-0.024	-0.051	-0.050
Standard deviation	0.082	0.014	-0.029	0.016	0.129	0.182
Avg neighbour degree	0.071	-0.008	0.152	0.046	-0.073	-0.104
Betweenness centrality	-0.007	-0.008	0.010	-0.005	-0.012	-0.078
Clustering coefficient	0.006	-0.022	0.152	0.076	-0.032	-0.035
Degree	0.032	0.018	0.164	0.006	-0.143	-0.134
Ego components size	0.026	0.044	0.133	0.002	-0.053	-0.022
HITS	-0.011	-0.034	-0.001	0.061	0.038	0.001
PageRank	-0.002	0.021	0.118	0.014	-0.099	-0.040
Two-hop neighbourhood	0.059	-0.015	0.130	0.012	-0.121	-0.135
ItemSimple Clarity	0.010	0.120	0.211	0.129	0.070	0.126

Table 6.17. Pearson's correlation between social-based predictors and P@10 for different recommenders using AR methodology (CAMRa Collaborative).

the number of users contained in the test set is twice the number available in the social version (see Appendix A.1.3), which means that significant correlations can be achieved with lower values (as described in Chapter 5).

In the results described above, we can observe how, like in the previous sections, the size of the user profile in test (predictor count in test) obtains significant correlations. This trend, however, is almost neutralised in the collaborative-social dataset with respect to the Random recommender. Thus, as before, we would attempt to use the 1R methodology with each dataset in order to obtain unbiased correlations towards users with more ratings in test. However, due to the lack of coverage of Personal and PureSocial recommenders, this methodology do not obtain sensible results (for instance, the value of precision at 10 is almost invariably 0.10, that is, the maximum possible value when only one relevant document – as assumed in the 1R methodology – is retrieved in the top 10, mainly because the recommender is not able to retrieve most of the not relevant items). This lack of coverage is natural for these recommenders since they can only suggest items rated by users in the active user’s social network (see Appendix A.2 for details on the implementation of the algorithms).

In conclusion, most of the social performance predictors proposed obtain significant correlations, however, **correlations with the social filtering methods are not so strong as we would expect**. Nonetheless, the **ItemSimple clarity does obtain significant correlations** with respect to most of the recommenders, highlighting the importance and validity of this predictor even when the main input of some recommenders (social network) is so different to that of the predictor (ratings).

6.5.5 Discussion

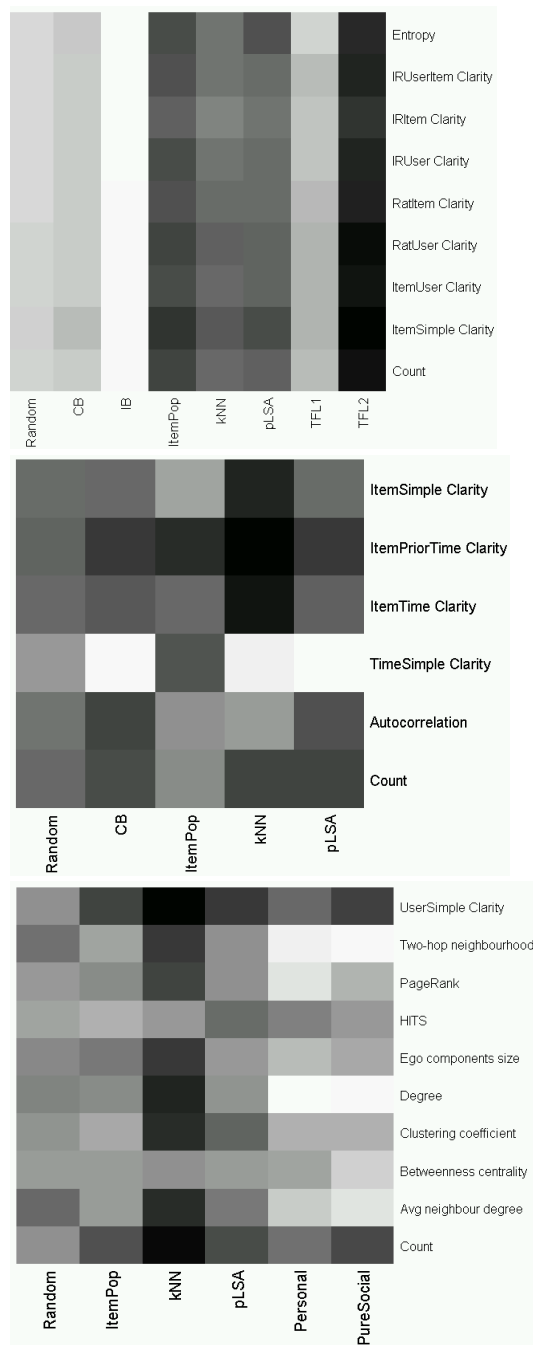
The reported experiments confirm that it is possible to predict a recommender’s performance and obtain strong correlations in this regard. The results show that, in general, the proposed predictors (mostly based on Kullback-Leibler divergences over different language models and other concepts from Social Graphs and Information Theory such as entropy) obtain significant correlations in the three spaces considered: ratings, logs, and social networks. More importantly, these correlations are stronger than those obtained by more simple predictors, such as the profile size of a user, the standard deviation of her ratings, and the user’s performance using a validation split. Specifically, for each recommendation input considered we have observed the following:

- Clarity-based predictors are very powerful for rating-based preferences, in particular, the ItemSimple, IRUser, and IRUserItem clarity predictors obtain strong correlations for most of the recommendation methods.

- The use of the item space as a contextual variable shows strong correlation values when the AR methodology is used, but these correlations decrease when we use unbiased methodologies, which may indicate that this new dimension is in fact capturing the item popularity and, thus, when the popularity bias is neutralised such predictors show less predictive power. We find a similar situation with the item clarity and the user space used as the contextual dimension.
- Temporal and log-based versions of the clarity predictor show higher prediction power than the rest of predictors.
- Social-based predictors are not the ones with the strongest correlation regarding the social filtering recommenders in this experiment, but the correlation found is significant and they could serve as a complement to other predictors based on a different input such as the rating-based.
- The ItemSimple clarity predictor consistently obtains strong correlation values in most of the datasets where we have analysed it. This is an evidence of the theoretical power of the user clarity to capture the uncertainty in user's tastes, even when the recommender's input is different (social filtering recommenders) or when we apply some transformation to the data (frequency-based clarity with transformation from implicit to explicit).
- As described in the Appendix A.4, most of the correlations presented in this chapter are stable when other evaluation metrics and correlation coefficients are used.

In the Recommender Systems field there are, however, additional problems due to subtle differences with respect to the common settings and experimental assumptions in Information Retrieval. Since we aim to predict the performance of a recommender, we have to be sure that we are using an unbiased performance metric, and its subsequent evaluation methodology. As we analysed in Chapter 4 there are at least two biases in the evaluation of recommender systems which may distort the results: data sparsity and item popularity. Thus, in this chapter we have computed correlations between the output of the predictors and the evaluation metrics using different evaluation methodologies, in order to analyse how sensitive the different proposed predictors are to these biases. Interestingly, although the correlations may change drastically when different evaluation methodologies are considered, most of the performance predictors still obtain good correlations. In particular, this result evidences that our proposed predictor are not so prone to the analysed biases like other simple predictors.

Finally, in Figure 6.3 we summarise the correlations found for the proposed predictors in each dimension – ratings, logs, and social. We have selected the most representative evaluation methodology (AR for rating and social data, and 1R for log



Correlation between rating-based user predictors and recommenders using AR methodology in MovieLens (Table 6.7).

Correlation between log-based user predictors and recommenders using 1R methodology in Last.fm (Table 6.14).

Correlation between social-based user predictors and recommenders using AR methodology in CAMRa collaborative (Table 6.17).

Figure 6.3. Heatmap of the correlation values between a subset of predictors and recommenders, using the most representative methodologies for the three considered spaces.

data) and a subset of the evaluated predictors and recommenders from each experiment, where the same information presented in Table 6.7, Table 6.14, and Table 6.17 (except for the average and median correlation values) is depicted in a more visual form. In particular, we may observe that predictors in MovieLens seem to be more redundant since the correlations are too similar.

From the figure we may also observe that in Last.fm and CAMRa datasets such redundancy is much lower and the predictors are quite different. Moreover, the first column and row (from the bottom) represent the recommender and predictor baselines, which serve as references from where the correlations should be analysed. In the three cases we can observe that most of the predictors obtain larger (darker) values than the count predictor. In the first case (rating-based predictors), however, it is clear that the correlation depends more on the recommender and less on the actual predictor.

6.6 Conclusions

We have proposed adaptations of query performance techniques from ad-hoc Information Retrieval to define performance predictors in Recommender Systems. Taking inspiration in the query predictor known as query clarity, we have defined and elaborated in the Recommender Systems domain several predictive models according to different formulations and assumptions. Furthermore, we propose performance predictors from theories and models of Information Theory, Social Graph Theory, and Information Retrieval based on three types of preference data: rating-based, log-based, and social-based.

We find several effective schemes with a high predictive power for recommender systems performance. We have proposed different ways for the adaptation of the query clarity predictor to recommender systems depending on the equivalences between the involved spaces. The clarity formulation is powerful because of its theoretical soundness, which is suitable to different domain-oriented adaptations. Hence, for rating-based preferences we use different expansions which take into account the rating values and the items rated by the user. For log-based preferences we exploit the co-occurrences of the items in the user profile and, more importantly, the temporal dimension, which allows for more principled functions such as the temporal Kullback-Leibler divergence or the user's autocorrelation. Finally, for social-based preferences we exploit the user's social network and different graph metrics are used apart from the user clarity based on the ratings. The results, as summarised in the previous section, are in general positive and provide evidences that the proposed functions are able to indeed predict the performance of user or items in recommender systems.

Furthermore, by analysing the behaviour of trivial predictors (such as the count of ratings in training and test) we have been able to uncover noisy biases or sensitivity to irrelevant variables in the way performance is measured. Irrelevant and uninteresting in the sense that it is not clear that the variations due to these variables really reflect actual differences in quality. As a result, we have used unbiased evaluation

methodologies where non trivial predictors still obtain positive results with respect to performance correlation.

As a side-effect, our study introduces an interesting revision of the gray sheep user concept. A simplistic interpretation of the gray sheep intuition would suggest that users with a too unusual behavior are a difficult target for recommendations. It appears however in our study that, on the contrary, users who somewhat distinguish themselves from the main trends in the community are easier to give well-performing recommendations. This suggests that perhaps the right characterisation of a gray sheep user might be one who has scarce overlap with other users. On the other hand, the fact that a clear user distinguishes herself from the aggregate trends does not mean that she does not have a sufficiently strong neighbourhood of similar users. In particular, this seems to indicate that users who follow mainstream trends are more difficult to be suggested successful items by a recommender system (at least, by a personalised one). In Information Retrieval, one can observe a similar trend: more ambiguous (mixture of topics) queries perform worse than higher-coherence queries (Cronen-Townsend et al., 2002).

In the future we plan to explore further performance predictors. Specifically, we are interested in incorporating explicit recommender dependence into the predictors, so as to better exploit the information managed by the recommender, allowing to the predictor a smoother adaptation to the recommender performance, and increasing the final correlation between them. Additionally, we are also interested in exploring alternative item-based predictors apart from those defined in this chapter, and, eventually, using other information sources such as log-based preference data and even the social network of the users who rated a particular item.

