

Table 2: Instantiations of the click model used in our experiments.

<i>perfect model</i>					
$R(d)$	0	1	2	3	4
$P(c R(d))$	0.00	0.20	0.40	0.80	1.00
$P(s R(d))$	0.00	0.00	0.00	0.00	0.00
<i>realistic model</i>					
$R(d)$	0	1	2	3	4
$P(c R(d))$	0.05	0.10	0.20	0.40	0.80
$P(s R(d))$	0.00	0.20	0.40	0.60	0.80

els) to study the behavior of comparison methods under *perfect* and *realistic* conditions. The perfect click model is used to examine the behavior of comparison methods in an ideal setting. The stop probabilities are 0 for all relevance levels, i.e., the user continues to examine documents regardless of the relevance of previously encountered documents. Therefore, all documents in the top 10 of each result list are examined. The click probabilities are defined such that highly relevant documents ($R(d) = 4$) are always clicked, and that non-relevant documents ($R(d) = 0$) are never clicked. Click probabilities for relevance levels between these extremes decrease monotonically. The realistic click model was designed to approximate noisier click behavior. Here, highly relevant documents are most likely to be clicked ($P(c|R(d) = 4) = 0.8$), and click probabilities decrease for lower relevance levels. Similarly, stop probabilities for this model are highest after seeing a highly relevant document and decrease for lower probabilities. Besides the different levels of noise, the two models differ in the overall expected number of clicks, which is lower for the realistic click model because the click probabilities are consistently lower and the stop probabilities are higher. Table 2 lists the probabilities defined for the two models.

The click models instantiated here are appropriate for our evaluation of comparison methods, as they satisfy the assumptions made by all these methods. In preliminary experiments, we found that the level of noise introduced by the realistic click model is of the same order of magnitude as observed in previous work [21]. In addition, we observed qualitatively similar results for all other instantiations we tested. Finer details, such as click behavior that takes relations between documents into account are not captured. Therefore, extending our evaluation to include real-life experiments is important future work.

5.2 Data

All experiments make use of the MSLR-WEB30k Microsoft learning to rank (MSLR) data set.⁴ This data set consists of approximately 30,000 queries with 136 precomputed document features and relevance judgments. Relevance judgments are provided on a 5-point scale. Unless stated otherwise, our experiments use the training set of fold 1. This set contains 18,919 queries, with an average of 9.6 judged documents per query.

To generate rankers to compare in our experiments, we make use of the features provided for the documents of the data set. Specifically, we take each feature to be a ranker, and formulate the task of the comparison method as identifying the ranking features that are expected to give the best performance. In our first set of experiments, we exhaustively compare all distinct ranker pairs that can be generated in this way. For the second and third sets, we select a small subset that appears most suitable for more detailed analysis.

⁴Downloaded from <http://research.microsoft.com/en-us/projects/mslr/default.aspx>.

The experiments are designed to evaluate the comparison methods’ ability to identify the better of two rankers based on (simulated) user clicks. To set the ground truth for which ranker is better, we use NDCG based on the explicit relevance judgments provided with the data set, following previous work [21]. We use the complete set of documents provided with the data set to determine NDCG for the ground truth, i.e., no cut-off is used. Note that comparisons between rankers using implicit feedback depend not only on the true difference in NDCG between two rankers but also on the magnitude of the difference per query and the number of affected queries.

5.3 Runs

All our experiments evaluate the same four interleaved comparison methods. As baselines, we use three existing methods (cf., §2):

- **balanced interleave**, as described in [23].
- **team draft**, as introduced in [23]; as in the original, clicks on top results are ignored when the documents are placed in identical order by both rankers.
- **document constraint**, introduced in [10], as described in §2.

In addition, we evaluate the experimental condition:

- **marginalized probabilities**: our model using the probabilistic interleave process (cf. §4.1) and the probabilistic comparisons defined in Eq. 4 (cf. §4.3).

Below we describe the three sets of experiments designed to assess interleaved comparison methods in terms of accuracy, convergence behavior, and robustness to noise.

Accuracy. The first experiment assesses the accuracy of interleaved comparison methods when evaluated on a large set of ranker pairs. We construct a ranker for each of the 136 individual features provided with the MSLR data set. We then evaluate the interleaved comparison methods on all 9,180 possible distinct pairs that can be obtained from these rankers. We evaluate the methods on this large set of ranker pairs using the perfect click model and a fixed subset of the data consisting of 1,000 queries. We then compare the outcome predicted by the interleaved comparison methods to the true NDCG of the rankers, computed on the same set of queries. We report on *accuracy* after 1,000 queries: the percentage of pairs for which a comparison method correctly identified the better ranker.

Convergence. Based on the findings of the first experiment, we select a small subset of ranker pairs with which to investigate the methods’ behavior in more detail. These experiments were designed to follow the setup used in [21] as closely as possible, with the difference that we use simulated interactions with the perfect click model instead of real-life search engine traffic, as discussed in §5.1. For each pair of rankers, each comparison method was run on randomly sampled query sets of size 1 to 10,000; queries for which an interleaved list did not obtain any clicks were ignored. Each experiment was repeated 1,000 times. We evaluate each method by computing, for each sampled query set size, the fraction of these 1,000 repetitions for which it detected a preference for the correct ranker. We compute binomial confidence intervals on this fraction using Wilson score intervals at the 95% level [3].

Noise. The third set of experiments assesses the robustness of interleaved comparison methods to noise in click feedback. The setup is identical to the previous one except that the *realistic* click model is used instead.