
On Using Simultaneous Perturbation Stochastic Approximation for IR Measures, and the Empirical Optimality of LambdaRank

Yisong Yue¹
Dept. of Computer Science
Cornell University
Ithaca, NY 14850
yyue@cs.cornell.edu

Christopher J. C. Burges
Microsoft Research
Microsoft Corporation
Redmond, WA 98052
cburges@microsoft.com

Abstract

One shortfall of existing machine learning (ML) methods when applied to information retrieval (IR) is the inability to directly optimize for typical IR performance measures. This is in part due to the discrete nature, and thus non-differentiability, of these measures. When cast as an optimization problem, many methods require computing the gradient. In this paper, we explore conditions where the gradient might be numerically estimated. We use Simultaneous Perturbation Stochastic Approximation as our gradient approximation method. We also examine the empirical optimality of LambdaRank, which has performed very well in practice.

1 Introduction

In recent years, the problem of learning to rank has gained importance in the fields of machine learning (ML) and information retrieval (IR). Ranking functions trained using ML techniques are currently in use by commercial web search engines. Due to the growth of interest, the amount of training data available for learning has likewise grown. One shortfall of existing methods is the inability to directly optimize for IR performance measures, such as mean average precision and normalized discounted cumulative gain (NDCG) [4].

Gradient descent is a common and effective method for directly optimizing an objective function within some search space. When casting learning to rank as an optimization problem, one can consider the search space to be the space of possible parameter values for some ranking function. Unfortunately, one cannot easily use IR measures as the objective function since, in general, they are not differentiable with respect to the ranking function parameters. As a result, ML methods typically optimize for a surrogate objective function (which is differentiable, and often convex) or use an approximate gradient. A more detailed version of this paper with additional results is available as a tech report [13].

Given the availability of larger training sets, we investigate whether the NDCG gradient might be numerically approximated. Numerical approximation requires measuring the change in the objective function over a small interval in the search space. This can become extremely expensive when dealing with high dimensional search spaces. We therefore use Simultaneous Perturbation Stochastic Approximation (SPSA) as our gradient approximation method, since it is very efficient and requires only two function evaluations for gradient approximation. We find that NDCG does become significantly smoother with additional

¹This work was conducted while the first author was an intern at Microsoft Research.

training data, but not enough to effectively perform gradient descent. However, we do anticipate that datasets of sufficient size might become available in the foreseeable future.

We also examine the potential optimality of LambdaRank. LambdaRank is a gradient descent method which uses an approximation to the NDCG “gradient”, and has performed very well in practice. Our experiments show that LambdaRank does in fact find a local optimum with respect to NDCG, even though the gradient used is a heuristic approximation.

2 Common Performance Measures for Information Retrieval

Performance measures used for information retrieval tasks are typically defined over rankings of documents. Relevance labels can be either binary (0 for non-relevant, 1 for relevant) or multilevel (0, 1, 2, ...). Binary measures include Mean Average Precision and Mean Reciprocal Rank. Normalized Discounted Cumulative Gain (NDCG) [4, 9] is a cumulative, multilevel measure that is usually truncated at a particular rank level. For a given query q_i , NDCG is computed as

$$\text{NDCG}_i \equiv N_i \sum_{j=1}^T \frac{2^{l_i(j)} - 1}{\log(1 + j)}, \quad (1)$$

where $l_i(j)$ is the label of the j th document in the ranking for q_i . The normalization constant N_i is chosen so that the perfect ranking would result in $\text{NDCG}_i = 1$, and T is the ranking truncation level at which NDCG is computed. NDCG is well suited for applications to Web search since it is multilevel and the truncation level can be set to model user behavior. Thus we will focus on NDCG in this paper.

2.1 Direct Optimization

Tuning model parameters to maximize performance is often viewed as an optimization problem in parameter space. In this setting, given a collection of training examples, we are concerned with optimizing NDCG with respect to the parameters of some ranking function.

Existing ranking functions usually score each document’s relevance independently of other documents. The ranking is then computed by sorting the scores. Variations to these ranking functions’ parameters will change the scores, but not necessarily the ranking. The measures discussed above are all computed over the rank positions of the documents. Therefore, the above measures have gradients that are zero wherever they are defined: that is, viewed as functions of the model score, typical IR measures are either flat or discontinuous everywhere.

However, what is optimized is usually the IR measure averaged over all queries in the training set. Given enough data, we might hope that the corresponding function becomes smooth enough for empirical gradients to be computable. This paper explores conditions under which an empirical NDCG gradient might exist and whether SPSA can be used to efficiently perform stochastic gradient descent. We focus on neural nets as our function class, which were also considered in [1, 2].

3 Related Work

Previous approaches to directly optimizing IR measures either used grid search, coordinate ascent, or steepest ascent using finite difference approximation methods (see Section 4 for discussion on stochastic approximation methods). Metzler & Croft [8] used a Markov Random Field ranking method and showed that MAP is empirically concave when using a parameter space with two degrees of freedom. In this study, we consider parameter spaces with much larger degrees of freedom.

Direct optimization becomes difficult with large datasets and parameter spaces with many degrees of freedom. Most other approaches choose instead to optimize an alternative smooth objective function. Perhaps the most straightforward approach is learning to predict the relevance level of individual documents using either regression or multiclass classification

(e.g., [7]). Another popular approach learns using the pairwise preferences between documents of different relevance levels (e.g., [5, 2, 3]). While these methods perform reasonably well in practice, they do not optimize for IR measures directly and offer no performance guarantees. Some recent studies focus on minimizing upper bounds of IR performance loss [14, 6, 12]. These methods do offer partial performance guarantees.

Another important class of approaches uses approximations to conceptualize a gradient for IR performance measures (despite these measures being non-differentiable in general). Of these, we examine LambdaRank by Burges et al. [1], as it performs very well in practice and, like this study, uses neural nets for its function class. We discover that LambdaRank appears to find a local optimum for NDCG.

3.1 LambdaRank

LambdaRank is a general gradient descent optimization framework that only requires the gradient to be defined, rather than the objective function. In the case of learning to rank, we focus on document pairs (i, j) of different relevance classes (i more relevant than j). The derivative of such pair with respect to the neural net’s output scores is defined as

$$\lambda_{ij} = N \left(\frac{1}{1 + e^{s_i - s_j}} \right) \left| (2^{\ell_i} - 2^{\ell_j}) \left(\frac{1}{\log(1 + r_i)} - \frac{1}{\log(1 + r_j)} \right) \right|,$$

where s_i is the output score, ℓ_i is the relevance label, and r_i is the (sorted) rank position of document i . The normalization factor N is identical to the one in (1). Let D_i^+ and D_i^- denote the set of documents with higher and lower relevance classes than i , respectively. The total partial derivative with respect to document i ’s output score is

$$\lambda_i \equiv \sum_{j \in D_i^-} \lambda_{ij} - \sum_{j \in D_i^+} \lambda_{ji}. \quad (2)$$

For each document, the LambdaRank gradient computes the NDCG gain from swapping rank positions with every other document (of a different relevance class) discounted by a function of the score difference between the two documents. This discount function is actually the RankNet gradient [2] (see (3) in Section 5.1 below). The objective function of LambdaRank can in principle be left undefined, since only the gradient is required to perform gradient descent, although for a given sorted order of the documents, the objective function is simply a weighted version of the RankNet objective function [2].

4 Stochastic Approximation

Assuming an objective function $L : \mathcal{R}^d \rightarrow \mathcal{R}$, an optimum w^* of L satisfies the property that the gradient vanishes (i.e., $\partial L(w^*)/\partial w = 0$). In cases when the gradient is not directly computable and evaluations of L are noisy, stochastic approximation techniques are often used to approximate the gradient.

Given an approximation $\hat{g}(w)$ to the true gradient, L can be iteratively optimized by stepwise gradient descent. The most common stochastic approximation technique, Finite Difference Stochastic Approximation (FDSA), approximates each partial derivative separately as

$$\hat{g}_k(w_k)_i = \frac{L(w_k + c_k e_i) - L(w_k - c_k e_i)}{2c_k},$$

where $c_k \in \mathcal{R}$ is the approximation step size and $e_i \in \mathcal{R}^d$ is the unit vector along the i th axis. This method requires $2d$ function evaluations at each iteration, which can be prohibitively expensive if L is non-trivial to compute (e.g., the forward propagation required to score documents and sort required to compute NDCG).

4.1 SPSA

We now describe Simultaneous Perturbation Stochastic Approximation (SPSA), which was first proposed by Spall [10, 11]. SPSA is an efficient method for stochastic gradient approximation. In contrast to FDSA, which performs $2d$ function evaluations per iteration, the simplest form of SPSA requires only 2.

As the name suggests, a simultaneous perturbation vector $\Delta_k \in \mathcal{R}^d$ is used in each iteration. Given Δ_k , the gradient approximation is computed as

$$\hat{g}_k(w_k) = \begin{bmatrix} 1/\Delta_{k1} \\ \dots \\ 1/\Delta_{kd} \end{bmatrix} \cdot \frac{L(w_k + c_k \Delta_k) - L(w_k - c_k \Delta_k)}{2c_k}.$$

Following the conditions stated in [10], Δ_k is a vector of d mutually we choose each Δ_{kl} to be independently symmetrically Bernoulli distributed (+1 or -1 with equal probability). When the objective function is extremely noisy or non-linear, multiple SPSA gradients can be computed and averaged together at each iteration.

The correctness results regarding SPSA are available in [10]. We ommit any detailed discussion due to space constraints. The results in [10] essentially conclude that SPSA produces an unbiased estimate of the true gradient and that the accumulated sampling errors will cancel out over time.

If evaluations of L is the computation bottleneck for gradient approximation, then each iteration of SPSA will be d times faster than FDSA. SPSA will thus converge faster than FDSA if it requires less than d times as many iterations. Spall [10] empirically demonstrated such results for relatively low dimensional data. In this paper we will empirically evaluate the convergence rate of SPSA vs. FDSA on a large Web search dataset.

5 Experiments

We performed experiments on two datasets: an artificial dataset and a real Web search dataset. Our experiments used neural nets trained with SPSA, FDSA and LambdaRank. Our experiments were designed to investigate three questions: (A) whether SPSA converges faster than FDSA for Web search data, (B) whether NDCG becomes empirically smooth given enough data, and therefore become trainable using SPSA, and if so then (C) whether SPSA can achieve results competitive with LambdaRank.

We tried a variety of learning rates for all our methods. We used a validation set to choose the best model for evaluation on the test set. We fixed the size of the hidden layer to 10 nodes for all two layer nets. Certain parameters were set as recommended by [10, 11]. We denote an SPSA variant using F function evaluations per iteration as SPSA:F. The basic SPSA algorithm is thus named SPSA:2.

We used both an artificial dataset as well as a “real” Web dataset generated from a commercial search engine. The Artificial dataset comprises of artificially generated data with 50 dimensions. The Web dataset is drawn from a commercial search engine and has 367 dimensions. More experiment and dataset details can be found in our technical report [13].

5.1 SPSA vs. FDSA

We empirically evaluated the convergence speed of SPSA vs. FDSA on the Web dataset for minimizing pairwise cross entropy as a “sanity check”. Pairwise cross entropy is the objective used for RankNet training [2]. We chose this metric since the objective function is differentiable, and so the non-smoothness of the cost function is removed as a possible factor, enabling us to cleanly compare FDSA and SPSA. For this experiment, we chose for our function class single layer neural networks.

The Web data contains 367 features, causing FDSA to perform 734 objective function evaluations per iteration. We evaluated three SPSA variants. SPSA:2 performs one gradient approximation (2 function evaluations) per iteration. SPSA:4 and SPSA:8 perform two and four gradient approximations, respectively. We report the mean pairwise cross entropy values number of function evaluations over ten runs of each method.

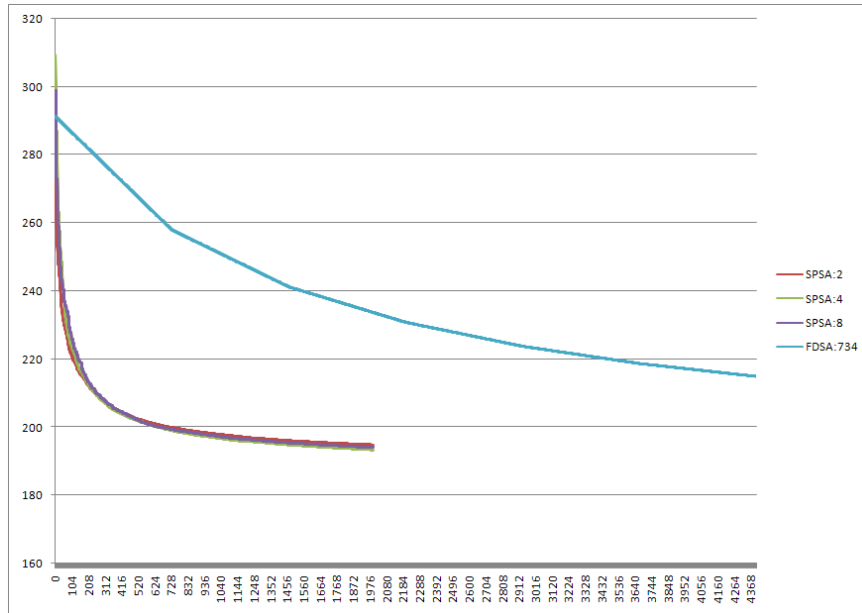


Figure 1: Cross Entropy w.r.t. Number of Function Evaluations on Web Training Data

For any pair of documents (i, j) with document i having higher relevance than j , the cross entropy gradient is computed as

$$\frac{\partial C_{ij}}{\partial s_i} = -\frac{\partial C_{ij}}{\partial s_j} = \frac{-1}{1 + e^{s_i - s_j}}. \quad (3)$$

where s_i is the neural network output score of document i . The total partial derivative for s_i can be written in the same form as (2). We empirically verified that the resulting FDSA gradient is virtually identical to the closed form solution. This allows us to use the gradient formulation in place of the FDSA gradient, and is much faster computationally.

The performance difference of SPSA vs FDSA is quite striking. Figure 1 shows the pairwise cross entropy value of FDSA and three variants of SPSA plotted against number of function evaluations. We see that all variants of SPSA converge significantly faster than FDSA. In this setting, performing a function evaluation requires forward propagating all the training examples to then performing a pairwise difference to compute the cross entropy. Interestingly, all variants of SPSA achieved roughly equivalent convergence rates.

5.2 Smoothness of NDCG

SPSA assumes the objective function (in this case NDCG) is thrice differentiable [10]. While the NDCG is either flat or discontinuous everywhere for a single query, it may become empirically smooth when averaged over a sufficient number of queries, just as any smooth function may be approximated as a linear combination of step functions. We can investigate this by taking a trained LambdaRank net and comparing the change mean NDCG of the training set as one weight of the net is varied while fixing the rest. We performed this comparison on both single and two layer nets. In all our comparisons, we varied each weight by a percentage of its trained weight.

We first compared the NDCG change when varying the weights of a single layer net. We report this comparison of the top five weights by absolute value. Figure 2 shows this result for a 100 query subsample, a 1,000 query subsample, and the entire 10,000 query Artificial training set. We observe that the NDCG function curve consistently becomes smoother as the query sample increases.

We also compared the NDCG change when varying the weights of a two layer net. We report this comparison for all ten weights connecting the hidden units to the output unit.

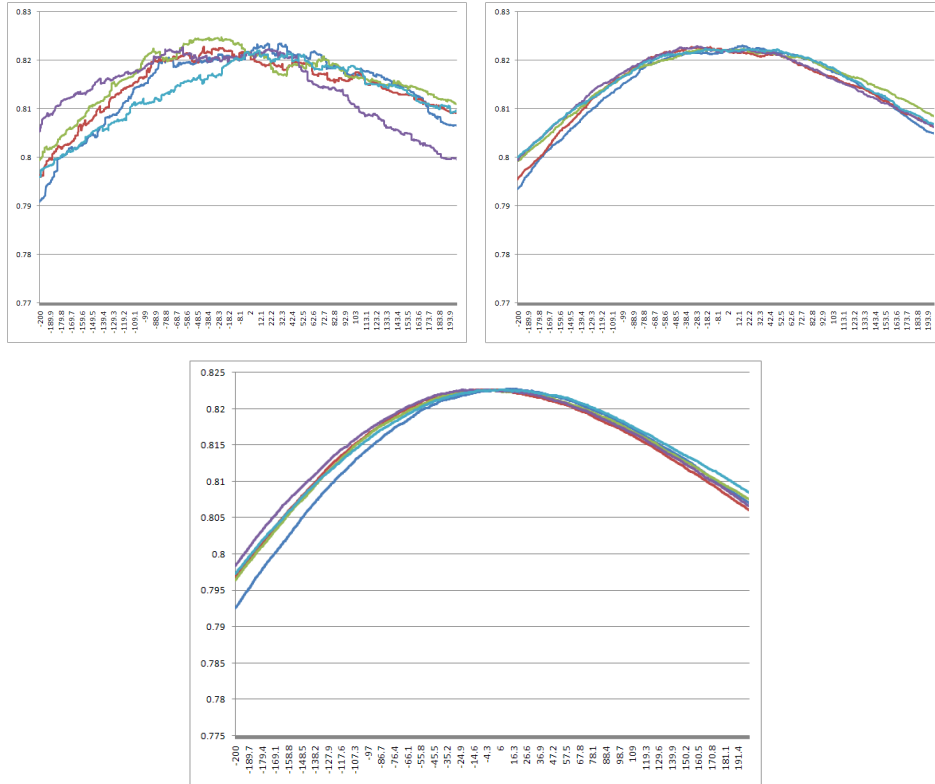


Figure 2: NDCG w.r.t. Shifting Top 5 Weights of Single Layer Net on 100, 1000 and 10000 Queries of Artificial Data

Figure 3 show the results for the entire Web training sets. Again, we observe that the NDCG function curves are relatively smooth on a macro scale. The inherent discontinuity of NDCG is observable only at a local scale.

However, the curves still contain numerous smaller local optima even when averaged over 10,000 queries. Figure 3 also shows a blown up section of the 2 layer net. We see that the discontinuities are very noticeable at this scale. In order for SPSA to work well, the scale of the discontinuities must be smaller than the gradient approximation step size. Not surprisingly, our results for SPSA are significantly worse than both LambdaRank and RankNet. Table 1 shows the comparison between LambdaRank and the best performing SPSA result.

Nonetheless, these results are encouraging. They suggest the feasibility of collecting a sufficient number of queries whereby the NDCG “gradient” becomes smooth enough for effective approximation. More generally, they suggest that many objective functions previously considered infeasible to optimize directly might yield computable gradient approximations.

Method	Train	Valid	Test
LambdaRank	0.721	0.713	0.707
SPSA:4	0.690	0.682	0.677

Table 1: NDCG@10 for Linear Nets on Web Data

5.3 Empirical Optimality of LambdaRank

The evaluation of the smoothness of NDCG also yields another interesting observation. The local optimum found using LambdaRank (2) also corresponds very closely to a local optimum of NDCG. Perhaps we should not find this too surprising, since the true NDCG “gradient” should reflect the instantaneous change in NDCG as the scores of the documents vary. The

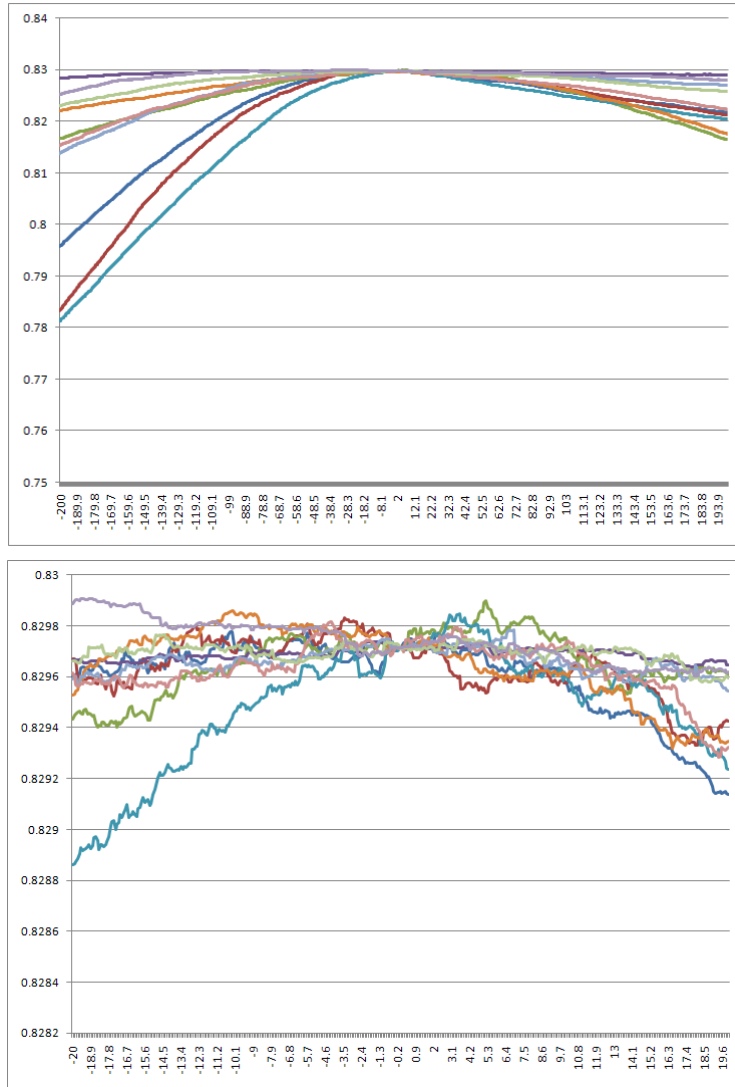


Figure 3: NDCG w.r.t. Shifting Hidden-Output Weights of Two Layer Net on Web Data, and a blown up section

LambdaRank gradient can be interpreted as a smoothed or convolved version of the change in NDCG, where the smoothing is a function of the distance between two documents' scores. It should be noted that LambdaRank is not simply a smoothed version of FDSA; it will contain contributions from documents that have very different scores, if the current ranker puts them in the wrong order. In fact, our prior experiments using FDSA (with respect to the model score) as the LambdaRank gradient yielded inferior results to the published LambdaRank gradient [1]. Roughly speaking, LambdaRank incorporates smoothing into its gradient approximation whereas SPSA requires empirical smoothing (averaging over enough data). This result suggests that it will be non-trivial to improve on LambdaRank performance using (two layer) neural networks as the function class.

6 Conclusions & Future Work

We've presented evidence demonstrating a trend towards smoothness of NDCG as the dataset size grows. While we cannot exactly characterize this smoothness, we find it reasonable to expect, in the foreseeable future, having enough data for methods such as SPSA

effective. The scale of the inherent discontinuities need only be smaller than the gradient approximation step size. Given the current training data available, we find that SPSA does not compare well with LambdaRank.

We also showed empirically that LambdaRank finds a local optimum for NDCG, despite using a (smooth) approximation of the NDCG gradient. Given these results, it appears difficult to improve on LambdaRank NDCG performance using (two layer) neural networks as the ranking function class.

These results also beg the question of whether LambdaRank has additional theoretical properties. One such question to ask is: given some distribution over the space of examples (queries and relevance labels), does a local optimum with respect to the LambdaRank gradient imply anything about the true gradient of the expected NDCG over that distribution?

We finally note that SPSA is a very general optimization framework. The objective function need only satisfy (or approximately satisfy) the condition that its third derivatives exist, and be bounded, in order for SPSA to work in practice. While optimizing for NDCG is a well-studied problem, for many IR optimization problems SPSA may work reasonably well.

Acknowledgements

The authors would like to thank John Platt for his helpful comments as well as for first pointing us to SPSA.

References

- [1] C. Burges, R. Ragno, and Q. Le. Learning to rank with non-smooth cost functions. In *Proceedings of NIPS'06*, 2006.
- [2] C. Burges, T. Sheked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender. Learning to rank using gradient descent. In *Proceedings of ICML'05*, Bonn, Germany, August 2005.
- [3] A. Herschtal and B. Raskutti. Optimising area under the roc curve using gradient descent. In *Proceedings of ICML'04*, 2004.
- [4] K. Jarvelin and J. Kekalainen. IR evaluation methods for retrieving highly relevant documents. In *Proceedings of SIGIR'00*, 2000.
- [5] T. Joachims. A support vector method for multivariate performance measures. In *Proceedings of ICML'05*, 2005.
- [6] Q. Le and A. Smola. Direct optimization of ranking measures. arXiv:0704.3359, 2007.
- [7] P. Li, C. Burges, and Q. Wu. Learning to rank using classification and gradient boosting. Technical report, Microsoft Research, 2007.
- [8] D. Metzler and B. Croft. A markov random field for term dependencies. In *Proceedings of SIGIR'05*, 2005.
- [9] S. Robertson and H. Zaragoza. On rank-based effectiveness measures and optimisation. Technical report, Microsoft Research, 2006.
- [10] J. Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE Transactions on Automatic Control*, 37:332–341, 1992.
- [11] J. Spall. Implementation of the simultaneous perturbation algorithm for stochastic approximation. *IEEE Transactions on Aerospace and Electronic Systems*, 34:817–823, 1998.
- [12] J. Xu and H. Li. A boosting algorithm for information retrieval. In *Proceedings of SIGIR'07*, 2007.
- [13] Y. Yue and C. Burges. On using simultaneous perturbation stochastic approximation for learning to rank; and, the empirical optimality of lambdarank. Technical Report MSR-TR-2007-115, Microsoft Research, 2007.
- [14] Y. Yue, T. Finley, F. Radlinski, and T. Joachims. A support vector method for optimizing average precision. In *Proceedings of SIGIR'07*, 2007.