

# BanditRank: Learning to Rank Using Contextual Bandits

Phanideep Gampa\*

Indian Institute of Technology (BHU) Varanasi  
gampa.phanideep.mat15@iitbhu.ac.in

Sumio Fujita

Yahoo Japan Corporation  
sufujita@yahoo-corp.jp

## ABSTRACT

We propose an extensible deep learning method that uses reinforcement learning to train neural networks for offline ranking in information retrieval (IR). We call our method BanditRank as it treats ranking as a contextual bandit problem. In the domain of learning to rank for IR, current deep learning models are trained on objective functions different from the measures they are evaluated on. Since most evaluation measures are discrete quantities, they cannot be leveraged by directly using gradient descent algorithms without an approximation. BanditRank bridges this gap by directly optimizing a task-specific measure, such as mean average precision (MAP), using gradient descent. Specifically, a contextual bandit whose action is to rank input documents is trained using a policy gradient algorithm to directly maximize the reward. The reward can be a single measure, such as MAP, or a combination of several measures. The notion of ranking is also inherent in BanditRank, similar to the current *listwise* approaches. To evaluate the effectiveness of BanditRank, we conducted a series of experiments on datasets related to three different tasks, i.e., web search, community, and factoid question answering. We found that it performs better than state-of-the-art methods when applied on the question answering datasets. On the web search dataset, we found that BanditRank performed better than four strong listwise baselines including LambdaMART, AdaRank, ListNet and Coordinate Ascent.

## CCS CONCEPTS

• **Information systems** → **Learning to rank**; *Novelty in information retrieval.*

## KEYWORDS

Information Retrieval, Learning to Rank, Question Answering, Web Search, Contextual bandits, Policy Gradient, REINFORCE

## ACM Reference Format:

Phanideep Gampa and Sumio Fujita. 2019. BanditRank: Learning to Rank Using Contextual Bandits. In *ACM*, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

\*Work conducted while the first author was in research internship at Yahoo! JAPAN Research

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

*Arxiv*, *ACM*

© 2019 Association for Computing Machinery.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Learning to rank is an important sub-field of information retrieval (IR), which involves designing models that rank documents corresponding to a query in order of their relevance. Considering the type of learning approach used, all ranking models can be classified into three categories, i.e., pointwise, pairwise, and listwise. The ranking models are either trained on indirect objective functions, such as classification related functions, or direct objective functions related to the evaluation measures. Direct optimization of IR measures has been a long standing challenge in the learning-to-rank domain. If we only consider bounded IR measures such as MAP, a theoretical justification is provided regarding the superiority of direct optimization techniques [31]. That study states that if an algorithm can directly optimize an IR measure on the training data, the ranking function learned with the algorithm will be one of the best ranking functions one can obtain in terms of expected test performance with respect to the same IR measure. Several algorithms have been developed that use direct optimization, and they can be grouped into three categories. The algorithms in the first category try to optimize the surrogate objective functions, which are either upper bounds of IR measures [7, 51, 55] or smooth approximations of IR measures [12, 41]. The algorithms in the second category smoothly approximate the true gradient of the evaluation measures, similar to LambdaRank [4, 5, 10, 54]. The algorithms in the third category directly optimize evaluation measures in the form of rewards without any approximation using reinforcement learning such as MDPRank [48, 58]. However, except for some algorithms like LambdaRank, most of the algorithms in all the categories are only suitable for models with less parameters [10] making it difficult to use deep neural networks, which are quite effective.

Deep learning [21] models have been proven to be effective with state-of-the-art results in many machine learning applications such as speech recognition, computer vision, and natural language processing, which leads to the introduction of neural networks in IR. Neural networks have been used for functions such as automatic feature extraction and comparison and aggregation of local relevance [13, 16, 17, 26, 42]. But, the neural networks are generally trained on objective functions such as cross entropy, which is not related to the evaluation measures. They do not have information about the measures that they are going to be evaluated on, i.e., the objective functions indirectly optimize the evaluation measures. Since most evaluation measures such as MAP, mean reciprocal rank (MRR), and normalized discounted cumulative gain (nDCG) are not differentiable, they cannot be used as the objective functions for training the neural networks.

For leveraging the efficacy of neural networks and superiority of direct optimization, we propose an extensible deep learning method called BanditRank. BanditRank formulates ranking as a contextual bandit problem and trains neural networks using the policy gradient algorithm [38], for directly maximizing the target measures.

Contextual bandit is a type of reinforcement learning algorithm used in decision-making scenarios in which an action has to be taken by an agent depending on the provided context. The exact details of the formulation are provided in Section 3. BanditRank follows the listwise approach by treating a query and the corresponding candidate documents as a single instance for training. BanditRank is extensible in the sense that it provides a methodology for training neural networks using reinforcement learning for ranking. Therefore, it can be used with any text-matching architecture for feature extraction and jointly trained or it can use the features extracted from a pre-trained model. For example, the LETOR 4.0 dataset [30] provides 46-dimensional feature vectors corresponding to each query-document pair that can be leveraged directly for training. Since BanditRank is a deep learning method, it is also extensible with respect to the choice of architectures that can be used. We focused on offline ranking tasks in which external relevance labels are provided for training.

Empirically, we prove that BanditRank is superior when compared to other strong baselines. We conducted a series of experiments on three datasets in the domains of question answering and web search. The major contributions of this paper are summarized as follows:

- For training neural networks by directly optimizing evaluation measures using gradient descent algorithms, we formulate the ranking problem as a contextual bandit and introduce a new deep learning method called *BanditRank* for ranking.
- To the best of our knowledge, BanditRank is the first listwise deep learning method that uses reinforcement learning to train neural networks for offline ranking purposes. We enabled this by introducing a hybrid training objective in order to solve the exploration problem when the number of possible actions is large.
- BanditRank provided state-of-the-art results when applied on both InsuranceQA [11] and WikiQA [52] datasets outperforming the previous best method at the time of writing of this paper.
- In the web-search task, when applied on the benchmark MQ2007 [30] dataset using only the provided 46-dimensional features, BanditRank achieved better results than the state-of-the-art learning-to-rank algorithm LambdaMART [4] and clearly outperformed other listwise baselines such as CoordinateAscent [24], ListNet [6], and AdaRank [51].

The remainder of the paper is structured as follows. In the next section, we briefly discuss related studies. In Section 3, we give the formulation of ranking as a contextual bandit. In Section 4, we provide the details of the model architecture used for our experiments. We explain the experiments we conducted along with a comparative study of rewards in Section 5. We conclude the paper in Section 6.

## 2 RELATED WORK

BanditRank is similar to BanditSum [9], which was proposed earlier for extractive summarization tasks in NLP. BanditSum introduces a theoretically grounded method based on contextual bandit formalism for training neural-network-based summarizers with reinforcement learning. We have adapted the formulation of ranking

as a contextual bandit from that of BanditSum. Adaptation of the contextual bandit framework to the ranking problem is not straightforward at all, for example, a naive application of BanditSum suffers from inadequate exploration when the number of actions is very large which is prevalent in ranking tasks. Thus we propose the use of hybrid loss for leveraging the feedback from a supervised loss function as explained in Section 3.4. Reinforcement learning was used for directly optimizing measures such as BLEU [27] and ROUGE [23] in different tasks of natural language processing such as summarization and sequence prediction [1, 22, 28, 33].

In the domain of learning-to-rank for IR, MDPRank [48] uses reinforcement learning for ranking by formulating ranking as a sequential decision process. Since the sequential models are affected by the order of the decisions, they may be biased towards selecting documents with low relevance level at the beginning. MDPRank is not suitable for training neural networks because a model with only 46 weight parameters requires more than 10000 epochs for convergence. In contrast, BanditRank is suitable for deep architectures, and all the best results of BanditRank were achieved in less than 30 epochs of training. Another issue with the setting of MDPRank is that the number of possible rankings for a query  $q$  with  $n_q$  number of candidate documents is  $n_q!$ , which is quite large making exploration more difficult. In contrast, BanditRank has more flexibility and freedom to explore the search space, as it samples a fixed number of documents  $M$  without replacement based on the affinity scores during training, reducing the search space to  $n_q P_M \ll n_q!$  for small  $M$ . This is because BanditRank uses listwise approach by treating all the candidate documents corresponding to a query as a single state. The policy gradient algorithm was also used to train the generator of IRGAN [44], but the rewards for the generator depend on the scoring function learned by the discriminator. The training of IRGAN is similar to that of SeqGAN [53], which is based on the idea of using the policy gradient algorithm for tackling the generator differentiation problem due to the discrete outputs produced by the generator. Both Bandits [18, 20, 32] and MDPs [56] were used to model the interactive process between a search engine and user with the user providing implicit relevance feedback. An overview of approaches that use reinforcement learning for different IR tasks such as query reformulation, recommendation and session search can be found in a previous paper [58]. BanditRank’s action is similar to the formulation of ListNet [6], which is based on the permutation of the input documents. However, both approaches differ with respect to the training method and structure of the probability model used.

## 3 BANDITRANK FORMULATION

We formulate ranking as a contextual bandit trained using policy gradient reinforcement learning. A bandit is a decision-making algorithm in which an agent repeatedly chooses one out of several actions and receives a reward based on this choice. The goal of the agent is to maximize the cumulative reward it achieves by learning the actions that yield good rewards. The term *agent* is generally used to refer to an entity or model that interacts with the environment. Contextual Bandit is a variant of the bandit problem that conditions its action on the context or state of the environment and observes the reward for the chosen action only. It forms a

subclass of Markov decision processes with the length of each episode being one. Formally, assume there is an environment with context space  $X$  and action space  $A$ . The agent interacts with the environment in a series of time steps. At each time step  $t$ , the agent observes a context  $x_t \in X$ , chooses an action  $a_t \in A$ , and observes a reward for that action  $r(a_t)$ . The goal of the agent is to maximize the cumulative rewards it achieves over a certain period.

Now, we can formulate the ranking problem as a contextual bandit with the environment being the dataset of queries and documents. The set of query-document pairs corresponding to a single query is treated as a context, and each permutation of the candidate documents is treated as a different action. Formally, given a query  $q$  and its candidate documents  $d = \{d_1, d_2, \dots, d_{n_q}\}$ , each context is the set  $c$  given by  $c = \{(q, d_1), (q, d_2), \dots, (q, d_{n_q})\}$ . Where  $n_q$  is the number of candidate documents of  $q$ , and the cardinality of  $c$  is given by  $n_c = n_q$ . Given  $c$ , the action of the agent is given by the permutation  $a_c = (d_{k_1}, d_{k_2}, \dots, d_{k_{n_q}})$  of the candidate documents, where  $k_t \in \{1, 2, \dots, n_q\}$  and  $k_t \neq k_{t'}$  for  $t \neq t'$ . The reward is given by a scalar function  $R(a_c, g_c)$  that takes action  $a_c$  and the ground-truth permutation  $g_c$  corresponding to  $c$  as the input. The  $g_c$  is nothing but the candidate documents sorted in descending order according to their relevance levels. The notation  $R$  is a scalar reward function defined using a combination of measures such as MAP and MRR.

The action taken by the agent is determined by its *policy*. In the current formulation, a policy is a neural network  $p_\theta(\cdot|c)$  parameterized by  $\theta$ . For each input  $c$ ,  $p_\theta(\cdot|c)$  encodes a probability distribution over permutations of the candidate documents. The goal is to find  $\theta$  that cause the network to assign high probability to the permutations, which can yield good rewards induced by  $R$ . This can be achieved by maximizing the following objective function with respect to  $\theta$ :

$$J(\theta) = E[R(a_c, g_c)], \quad (1)$$

where the expectation is taken over  $c$  paired with  $g_c$  and  $a_c$  generated according to  $p_\theta(\cdot|c)$ . The above objective function is a standard objective function used in the reinforcement-learning domain, which maximizes the expected reward. The negative of the expectation can be treated as the loss function.

### 3.1 Structure of Policy $p_\theta(\cdot|c)$

The exact action of the agent depends on the chosen structure of  $p_\theta(\cdot|c)$ . We follow the approach used for extractive summarization [9] because of its simplicity and effectiveness. With this approach,  $p_\theta(\cdot|c)$  is decomposed into a deterministic function  $\pi_\theta$ , which contains all the network's parameters, and  $\mu$ , a probability distribution induced by the output of  $\pi_\theta$  defined as

$$p_\theta(\cdot|c) = \mu(\cdot|\pi_\theta(c)) \quad (2)$$

Provided a  $c$  corresponding to a  $q$ , the network  $\pi_\theta$  outputs a real valued vector of document affinities within the range  $[0, 1]$ . The length of the vector is equal to the number of candidate documents  $n_c$  corresponding to  $q$ , i.e.,  $\pi_\theta(c) \in \mathbb{R}^{n_c}$ . The affinity score of a document  $d_i$  given by  $\pi_\theta(c)_i$  represents the network's propensity to keep the document at the top position in the output permutation. Specifically, the interpretation of the affinity scores is highly dependent upon the type of reward signal used. For example, if

only Precision@1 is used as the reward signal, the focus of the network would mainly be on the permutations that contain a relevant document at the first position.

Provided the above document affinities  $\pi_\theta(c)$ ,  $\mu$  implements a process of repeated sampling without replacement by repeatedly normalizing the set of affinities of documents not yet selected. In total,  $M$  unique documents are sampled yielding an ordered subset of the candidate documents. For exploring the action space, a small probability  $\epsilon$  of sampling uniformly from all remaining documents is included at each step of the sampling. This is similar to the  $\epsilon$ -greedy technique generally used in the reinforcement learning problems for exploration. According to the prescribed definition of  $\mu$ , the probability  $p_\theta(a_c|c)$  of producing a permutation  $a_c$  corresponding to  $c$  according to (2) is given by

$$p_\theta(a_c|c) = \prod_{i=1}^M \left( \frac{\epsilon}{n_c - i + 1} + \frac{(1 - \epsilon)\pi_\theta(c)_{k_i}}{z(c) - \sum_{l=1}^{i-1} \pi_\theta(c)_{k_l}} \right), \quad (3)$$

where  $k_t$  is the index to the  $t$ -th document in  $a_c$ ,  $d_{k_t}$  and  $z(c) = \sum_{m=1}^{n_c} \pi_\theta(c)_m$ . We define  $M = \min(n_c, M')$ , where  $M'$  is an integer hyper parameter that depends on the environment or dataset. This sampling method with exploration is followed only during training time. At test time, we output all the candidate documents sorted in descending order according to their affinity scores.

### 3.2 Policy Gradient Reinforcement Learning

The gradient of the objective function (1) cannot be calculated directly as  $a_c$  is discretely sampled while calculating  $R(a_c, g_c)$ . This is a common situation in most reinforcement-learning tasks. However, the gradient of the objective function can be calculated after a reformulation of the expectation term according to the REINFORCE algorithm [38, 49]. It tells us that the gradient of that function can be calculated using the following equation:

$$\nabla_\theta J(\theta) = E[\nabla_\theta \log p_\theta(a_c|c) R(a_c, g_c)], \quad (4)$$

where the expectation is over the same variables as (1).

Given a context-true permutation pair  $(c, g_c)$  sampled from the dataset or environment  $D(c, g_c)$ , the gradient can be derived using the following reformulation of the expectation in (1):

$$\begin{aligned} \nabla_\theta J(\theta) &= \nabla_\theta E[R(a_c, g_c)] \\ &= \nabla_\theta E_{(c, g_c) \sim D(c, g_c), a_c \sim p_\theta(\cdot|c)}[R(a_c, g_c)] \\ &= \sum_{(c, g_c) \sim D(c, g_c)} \sum_i \nabla_\theta p_\theta(a_c^i|c) R(a_c^i, g_c) \quad (5) \\ &= \sum_{(c, g_c) \sim D(c, g_c)} \sum_i p_\theta(a_c^i|c) \nabla_\theta \log p_\theta(a_c^i|c) R(a_c^i, g_c) \quad (6) \\ &= E_{(c, g_c) \sim D(c, g_c), a_c \sim p_\theta(\cdot|c)}[\nabla_\theta \log p_\theta(a_c|c) R(a_c, g_c)] \\ &= E[\nabla_\theta \log p_\theta(a_c|c) R(a_c, g_c)] \end{aligned}$$

Step (5) follows from the definition of expectation for discrete quantities and the linearity of the gradient operator. Step (6) is the reformulation of the gradient term, which is an important step in the derivation initially given by the REINFORCE [49] algorithm. The expectation in (4) is empirically calculated by first sampling a context-true permutation pair  $(c, g_c)$ , sampling  $B$  permutations  $a_c^1, a_c^2, \dots, a_c^B$  from  $p_\theta(\cdot|c)$  using the sampling method mentioned

in Section 3.1, and finally taking the average. Empirically, the inner expectation of (4) is given by

$$\nabla_{\theta} J_c(\theta) \approx \frac{\sum_{i=1}^B \nabla_{\theta} \log p_{\theta}(a_c^i | c) R(a_c^i, g_c)}{B} \quad (7)$$

The number  $B$  is also an integer hyperparameter that mainly depends on the dataset. Given the expression for  $p_{\theta}(a_c | c)$  (3), the gradient (7) can be calculated by any automatic differentiation package. As mentioned in Section 3.1, we sample  $M = \min(n_c, M')$  number of documents from the candidate documents during training time. Therefore, we take reward feedback from an  $M$ -length ordered subset. Since we cannot efficiently explore the whole action space for large  $M$  as the number of possible actions or permutations would then become  $n_c P_M^1$ , we choose  $M$  based on the average number of relevant documents per query in the dataset. The  $B$  determines the exact number of actions we explore for each context during each epoch, which can be seen from the approximation in (7). In our experiments, we obtained very good results even-though  $B$  was set to be small, i.e.,  $B = 20$  or  $B = 30$ .

The gradient estimate in (7) is prone to have high variance [38]. Moreover, all target measures, such as MAP, MRR, and nDCG, are always non negative, which increases the probability of every sampled permutation according to the objective function. We would prefer the probability of a bad permutation in terms of reward should be decreased. We use a baseline function, which is subtracted from all rewards. This decreases the variance of the estimate by acting as an advantage function, and it ensures that the permutations with low rewards receive negative rewards. If chosen appropriately, the advantage function can significantly reduce the variance of the estimate [38] without biasing the estimate. Using a baseline  $r_{base}$ , the sample-based estimate (7) becomes

$$\nabla_{\theta} J_c(\theta) \approx \frac{\sum_{i=1}^B \nabla_{\theta} \log p_{\theta}(a_c^i | c) [R(a_c^i, g_c) - r_{base}]}{B} \quad (8)$$

For choosing the baseline function, we follow the terminology of self-critical reinforcement learning, in which the test time performance of the current model is used as the baseline [9, 28, 33, 35]. Therefore, while calculating the gradient estimate (8) after sampling the context-true permutation pair  $(c, g_c)$ , we greedily generate a permutation using the current model similar to the test time action.

$$a_c^{greedy} = \arg \max_{a_c} p_{\theta}(a_c | c) \quad (9)$$

The baseline for a  $c$  is then calculated by setting  $r_{base} = R(a_c^{greedy}, g_c)$ . Therefore, all the permutations with reward greater than the *greedy* permutation receive positive rewards and other permutations receive negative rewards. The baseline is also intuitive in the way that it is different for different contexts.

### 3.3 Reward Function $R$

As mentioned earlier, the reward function can be a single target measure or a combination of several measures. For the question answering datasets, the following reward function was used:

$$R(a_c, g_c) = \frac{AP(a_c, g_c) + RR(a_c, g_c)}{2} \quad (10)$$

<sup>1</sup>Permutation  $nPr$  is an increasing function of  $r$ .

For the web search dataset, the following reward function was used:

$$R'(a_c, g_c) = \frac{AP(a_c, g_c) + nDCG@10(a_c, g_c)}{2}, \quad (11)$$

where the measures average precision (AP), reciprocal rank (RR), and nDCG@10 are traditional IR measures. In the experiments section, we also provide a simple comparison of different reward functions on the web search dataset.

### 3.4 Hybrid Training Objective

As mentioned in the Section 3.2, the problem of exploring the action space when  $M$  is large can be tackled using a hybrid loss, which is a combination of the reinforcement learning loss and a standard supervised learning loss such as binary cross entropy. The supervised loss can guide the training initially when the exploration by the model is in the starting stages for large  $M$ . Even though the number of actions explored with the model at each epoch given by  $B$  is small for large  $M$ , i.e.,  $B \ll n_c P_M$ , a supervised signal can help the model by compensating the loss incurred due to the inefficient exploration. The hybrid loss function is given as follows:

$$L_{hybrid} = \gamma L_{rl} + (1 - \gamma) L_{sl} \quad (12)$$

where  $L_{rl}$  is the loss given by the reinforcement-learning algorithm, which is the negative of (1), and  $L_{sl}$  is a supervised loss such as binary cross entropy. The notation  $\gamma$  is a scaling factor accounting for the difference in magnitude between  $L_{rl}$  and  $L_{sl}$ . It is a hyperparameter lying between 0 and 1. We found the hybrid loss to be effective in the case of the web search dataset where the average number of relevant documents per query was equal to 10.3. Since we use binary cross entropy as the supervised loss, the hybrid training objective is a blend of the *pointwise* objective function  $L_{sl}$  and a *listwise* objective function  $L_{rl}$ . The hybrid training objective still has direct control over the target measures weighted by  $\gamma$ . Similar hybrid loss was used in the domain of NLP in some papers, e.g., [28, 50].

## 4 MODEL ARCHITECTURE

In this section, we discuss the neural architecture  $\pi_{\theta}$  we used in our experiments. For demonstrating the extensibility of BanditRank, we considered two scenarios for the experiments and show that BanditRank performs well in both the scenarios.

### Scenario 1

In this scenario, we decomposed  $\pi_{\theta}$  into two neural network architectures  $f_{\theta_1}$  and  $b_{\theta_2}$ :  $f_{\theta_1}$  for extracting the feature vectors from raw texts of query or documents and  $b_{\theta_2}$  for a bandit network that yields the document affinities. For  $f_{\theta_1}$ , since any text-matching neural network [3, 40, 45] that can provide a single feature vector corresponding to each query-document pair was suitable, we have used the architecture similar to the recently proposed Multi Cast Attention Networks (MCAN) [40]. The  $b_{\theta_2}$  architecture was chosen to be a simple feed forward neural network with an output sigmoid unit for yielding the document affinities corresponding to a query. While training, the two architectures were treated as a single architecture by positioning the bandit network on top of the text-matching network. Formally, provided with a context

$c = \{(q, d_1), (q, d_2), \dots, (q, d_{n_q})\}$ , we passed it through both networks to obtain the document affinities  $\pi_\theta(d)$  as given below:

$$\begin{aligned} f_{\theta_1}(c) &= c_1, c_2, \dots, c_{n_q} \\ b_{\theta_2}(c_1, c_2, \dots, c_{n_q}) &= \pi_\theta(d) \end{aligned}$$

where  $c_i$  is a feature vector corresponding to the query-document pair  $(q, d_i)$ .

## Scenario 2

In this scenario, we have only used  $b_\theta$  for training. The feature vectors corresponding to the query and document texts were extracted using a pre-trained neural language model such as BERT [8], which is a state-of-the-art unsupervised language model in NLP. In web search datasets, such as MQ2007 [30], 46-dimensional feature vectors composed of basic features including BM25, term frequency (TF), and LMIR, are provided for each query-document pair. We directly use those vectors while conducting experiment on MQ2007 dataset. Formally, provided with the query-document feature vectors  $c_1, \dots, c_{n_q}$  corresponding to a context  $c = \{(q, d_1), \dots, (q, d_{n_q})\}$ , we have obtained the document affinities  $\pi_\theta(d)$  as given below :

$$b_\theta(c_1, c_2, \dots, c_{n_q}) = \pi_\theta(d)$$

These scenarios are intended for separating the functionality of the text-matching and bandit networks. Therefore, the inputs and outputs are not necessarily in the above prescribed format.

The results obtained in both the scenarios indicate that the bandit network is not entirely dependent on a text-matching network for providing good results. The exact details of the scenarios and the neural network architectures used for each dataset are provided in the experiments section <sup>2</sup>.

## 5 EXPERIMENTS

We conducted our experiments on three different datasets in the domains of question answering and web search. For the question answer ring task, we tested BanditRank on InsuranceQA [11], which is a community question answering dataset (closed domain), and on WikiQA [52], which is a well studied factoid question answering dataset (open domain). For the web search task, we conducted our experiments on the benchmark MQ2007 [30] dataset. Since the baselines and evaluation measures are different for the above three datasets, we divide this section into three subsections each dealing with a specific dataset. For each dataset, we provide the details of the architecture used, implementation details, details of the baselines, and obtained results. We adopted Scenario 1 for the InsuranceQA dataset and Scenario 2 for the MQ2007 dataset and compared both scenarios on the WikiQA dataset. Finally, we conducted a comparative study on different reward choices for training on the MQ2007 dataset. We choose the exploration probability mentioned in Section 3.1 as  $\epsilon = 0.1$  for all the experiments.

### 5.1 Experiment 1: InsuranceQA

**5.1.1 Dataset and Evaluation Measures.** InsuranceQA [11] is a well studied community question answering dataset <sup>3</sup> with questions submitted by real users and answers composed by professionals

with good domain knowledge. In this task, BanditRank was expected to select the correct answer among a pool of candidate answers with negative answers randomly sampled from the whole dataset. The dataset consists of two test datasets for evaluation. The evaluation measure for this task was Precision@1. The statistics of the dataset along with the average number of relevant answers per question are given in Table 1.

**Table 1: Statistics of InsuranceQA dataset**

Split	# of questions	# of correct answers	# of avg-rel-per-q
train	12887	18540	1.43
dev	1000	1454	1.45
test1	1800	2616	1.45
test2	1800	2593	1.44

**5.1.2 Model and Implementation Details.** We used two architectures for Scenario 1. For the text matching part, we used the architecture of the recently proposed Multi Cast Attention Network (MCAN) [40]. Specifically, we used the same architecture until the mean max pooling operation layer of the MCAN, which provides a fixed dimensional feature vector of each question or answer sentence. We used the sum function (SM) as the compression function for casting the attention. For the bandit network, we modified the prediction layer of MCAN with a sigmoid unit in the output layer and kept the two highway layers [37] intact. We conducted the experiments by replacing the highway networks with simple neural networks with the same dimensions, but the best results were obtained when highway layers were used. Highway networks [37] are gated nonlinear transform layers that control the information flow similar to the LSTM layers [15] for sequential tasks.

We randomly initialized the embedding matrix with 100 dimensional vectors sampled from standard normal distribution and fixed them during training. The hidden size of the LSTM layers was set to 300. The dimensions of the two highway prediction layers were set to 200 with ReLU being the activation function. Once we obtained the feature vectors corresponding to the question and answer sentences from the text-matching network as  $h_q, h_a \in \mathbb{R}^{600}$ , we passed the following vector  $h_{qa} = [h_q; h_a; h_q \odot h_a; h_q - h_a] \in \mathbb{R}^{2400}$  to the bandit network, where  $\odot$  is the pointwise multiplication operation and  $[\cdot; \cdot]$  is a vector concatenation operator. The vectors  $h_{qa}$  correspond to the query-document feature vectors  $c_i$  mentioned in Section 4. Before passing  $h_{qa}$  to the highway layers of bandit network, we used a single feed forward layer with a ReLU activation function for projecting  $h_{qa}$  into a 200-dimensional space. A dropout of 0.2 was applied to all layers except the embedding layer. The sequences were padded to their batch-wise maximums. We optimized the model using the Adam optimizer [19] with the beta parameters set to (0.9, 0.999) and a weight decay of  $1e^{-6}$  was used for regularization. We used the hybrid training objective defined in Eq. (12) with  $\gamma$  tuned over the set of values [0.5, 0.75, 1]. An initial learning rate of  $5e^{-5}$  was used for BanditRank with  $\gamma = 1$  and  $1e^{-4}$  for BanditRank with other  $\gamma$  values. We set  $M', B$  to  $M' = 5$  and  $B = 20$  for calculating the gradient in Eq. (8). For BanditRank with  $\gamma = 0.75$ ,  $M'$  was set to 20. We used the reward function defined as Eq. (10) during training.

<sup>2</sup>We will make our implementation publicly available when the paper is accepted.

<sup>3</sup><https://github.com/shuzi/insuranceqa>

**Table 2: Precision@1 for InsuranceQA dataset. Best results are in bold and second best are underlined.**

	test-1	test-2
IR model [2]	0.551	0.508
QA-CNN [36]	0.6133	0.5689
LambdaCNN [36, 57]	0.6294	0.6006
IRGAN [44]	0.6444	0.6111
CNN with GESD [11]	0.653	0.61
Attentive LSTM [39]	0.69	0.648
IARNN-Occam [43]	0.689	0.651
IARNN-Gate [43]	0.701	0.628
Comp-Agg(MULT) [45]	0.752	0.734
Comp-Agg(SUBMULT+NN) [45]	0.756	0.723
BanditRank( $\gamma = 1$ )	<u>0.8494</u>	<u>0.8283</u>
BanditRank( $\gamma = 0.75$ )	<b>0.8572</b>	<b>0.8522</b>

**5.1.3 Baselines and Results.** We compared the performance of BanditRank against all current methods that achieved significant results on this dataset. The competitive baselines are the IR model [2], CNN with GESD (from the authors who created the InsuranceQA dataset) [11], Attentive LSTM [39], IARNN-Occam [43], IARNN-Gate [43], QA-CNN [36], LambdaCNN [36, 57], IRGAN [44], and the method with the previous best P@1 measure, Comp-Agg [45]. A description of all the baselines can be found in previous studies [44, 45]. As the testing splits were the same for all methods, we report the P@1 measures directly from those studies.

The results in Table 2 indicate the superiority of BanditRank over all other methods. BanditRank with  $\gamma = 1$  achieved the second best P@1 measure, this is equivalent to training the model only with the reinforcement loss. The results further improved using a hybrid loss with  $\gamma = 0.75$  weight to the reinforcement loss  $L_{rl}$ . Therefore, providing some weight to the supervised loss improved the performance of BanditRank. BanditRank exhibited significant improvement in P@1 measure by 13.3% on the test-1 dataset and 16.1% on the test-2 dataset when compared to the previous best method.

## 5.2 Experiment 2: WikiQA

**5.2.1 Dataset and Evaluation Measures.** WikiQA [52] is a well-known open domain question answering dataset in contrast to InsuranceQA, which is a closed domain question answering dataset. The dataset was constructed from real queries on Bing and Wikipedia. In this task, the models were expected to rank the candidate answers according to the question. The evaluation measures for this task were MAP and MRR. The statistics of the dataset are given in Table 3.

**Table 3: Statistics of WikiQA dataset**

Split	# of questions	# of correct answers	# of avg-rel-per-q
train	873	1040	1.19
dev	126	140	1.11
test	243	293	1.20

**5.2.2 Model and Implementation details.** We considered both the scenarios given in Section 4 for WikiQA dataset. For Scenario 1, we used the same setting of MCAN as InsuranceQA. The only difference was the type of embedding used. We initialized the embedding matrix with 300-dimensional GloVe embeddings [29] and fixed them during training. We used the reward function defined as Eq. (10) during training. A dropout of 0.2 was applied to all layers except the embedding layer. The sequences were padded to their batch-wise maximums. We optimized the model using the Adam optimizer [19] with the beta parameters set to (0.9, 0.999), and a weight decay of  $1e^{-6}$  was used for regularization. We used the hybrid training objective defined in Eq. (12) with  $\gamma$  tuned over the set of values [0.25, 0.5, 0.75, 1]. We provide the results of the two best performing models with respect to  $\gamma$ . An initial learning rate of  $5e^{-5}$  was used for BanditRank with  $\gamma = 1$  and  $1e^{-4}$  for BanditRank with other  $\gamma$  values. We set  $M', B$  to  $M' = 3$  and  $B = 20$  for calculating the gradient in Eq. (8). The  $M'$  was chosen according to the average number of relevant queries, which is much less for the WikiQA dataset.

For Scenario 2, we extracted word-level<sup>4</sup> features from a pre-trained<sup>5</sup> BERT [8] language model, which takes a question-answer sentence pair  $(q, a)$  as the input. There are two versions of the BERT model available, BERT-base and BERT-large. We conducted our experiments on features extracted with both versions. We used the concatenation of word-level features obtained from the last four layers of the BERT language model for training. BERT-base produced 3072-dimensional feature vectors for each word in the input sentence while BERT-large produced 4096-dimensional feature vectors after the concatenation. These contextual word embeddings were passed through a two-layer bidirectional LSTM layer followed by mean pooling for obtaining a fixed dimensional representation<sup>6</sup> of  $(q, a)$ . These vectors correspond to the  $c_i$  vectors mentioned in Section 4. Regarding the architecture of the bandit network, we chose a feed forward network with a single hidden layer followed by a sigmoid unit at the output layer. Tanh was used as the activation function. For the BanditRank method trained on features extracted from BERT-base, we set the dimensions of the LSTM layer to 768. For BERT-large, we set the dimensions of the LSTM layer to 1024. For the feed forward layer, we set the dimensions of the hidden layer to 256 for both type of features. A dropout of 0.4 was applied to all layers. We optimized BanditRank using the Adam optimizer [19] with the beta parameters set to (0, 0.999) and used a weight decay of  $1e^{-6}$  for regularization. We used the hybrid training objective defined in Eq. (12) with  $\gamma$  tuned over the set of values [0.5, 0.75, 1]. We provide the results of the best performing model with respect to  $\gamma$ . An initial learning rate of  $8e^{-5}$  was used for BanditRank with  $\gamma = 1$  and  $1e^{-4}$  for the BanditRank method with other  $\gamma$  values. We set  $M', B$  to  $M' = 5$  and  $B = 20$  for calculating the gradient in Eq. (8). We used the reward function defined with Eq. (10) during training.

<sup>4</sup>Please note that, even-though only word level features are extracted from BERT, these features encode contextual information of both the input sentences similar to that of a text matching network.

<sup>5</sup><https://github.com/huggingface/pytorch-pretrained-BERT>

<sup>6</sup>Although LSTM layers are used before the bandit network, we include this in Scenario 2 as the true text matching was carried out with the BERT model, LSTM layers were only used to obtain a fixed dimensional representation of the question-answer pair.

**Table 4: Test-set performance on WikiQA dataset. Best results are in bold and second best are underlined.**

	MAP	MRR
CNN-Cnt [52]	0.652	0.665
QA-CNN [36]	0.689	0.696
NASM [25]	0.689	0.707
Wang et al [47]	0.706	0.723
He and Lin [14]	0.709	0.723
NCE-CNN [34]	0.701	0.718
BIMPM [46]	0.718	0.731
Comp-Agg [45]	0.743	0.755
Comp-Clip [3]	<u>0.754</u>	<u>0.764</u>
Scenario 1		
BanditRank( $\gamma = 0.75$ )	0.6663	0.673
BanditRank( $\gamma = 1$ )	0.7043	0.716
Scenario 2		
BanditRank-BERT-base ( $\gamma = 1$ )	0.7437	0.7589
BanditRank-BERT-large ( $\gamma = 1$ )	<b>0.7649</b>	<b>0.7807</b>

**5.2.3 Baselines and Results.** We compared the performance of BanditRank against all other previous methods on this dataset. The baselines were CNN-Cnt [52], QA-CNN [36], NASM [25], Wang et al [47], He and Lin [14], NCE-CNN [34], BIMPM [46], Comp-Agg [45], and the state-of-the-art method Comp-Clip [3]. Since the testing split is same for all methods, we report the highest measures directly from the respective papers.

The results given in Table 4 indicate the superiority of BanditRank over all other methods. BanditRank trained using the features extracted from BERT-large produced the best results. Interestingly in Scenario 1, we observed performance degradation when hybrid loss was used. Although, this degradation may depend on many factors, one possible explanation can be that the model can explore with the help of  $L_{r,l}$  efficiently since the average number of relevant documents is much less. The results in Scenario 2 indicate that, provided with good features, training a text-matching network along with the bandit network is not necessary for achieving good results.

### 5.3 Experiment 3: MQ2007

**5.3.1 Dataset and Evaluation Measures.** For the web search task, we used the benchmark Million Query Track 2007 (MQ2007) [30] dataset <sup>7</sup>. In this task, BanditRank was expected to rank the documents corresponding to a query according to their relevance. Unlike the previous tasks in which the relevance was binary, this task consisted of multiple-levels of relevance with  $\{0, 1, 2\}$ . This dataset provides 46-dimensional feature vectors corresponding to each query-document pair. Moreover, the average number of relevant documents per query was very large compared to the previous tasks. The statistics are given in Table 5. Out of the total 1692 queries, the number of queries with at least one relevant document was only 1455.

The loss for the above mentioned 237 query instances would be zero as the reward generated would be zero. Therefore, we conduct experiments on the dataset after removing the queries

<sup>7</sup><https://www.microsoft.com/en-us/research/project/letor-learning-rank-information-retrieval/>

**Table 5: Statistics of MQ2007 dataset**

	MQ2007
# of queries	1692
# of q-with-rel	1455
# of documents	65,323
# of avg-rel-per-q	10.3
# of features	46

with no relevant documents as they would not help BanditRank during training. We carried out 60-20-20 splitting for the train-val-test datasets after the dataset was cleaned. The baselines were also trained and evaluated on the same splits. As per evaluation measures, we report the measures of MAP, MRR, precision, and nDCG at positions 1, 3, 10. We also conducted significance tests using both paired t-test and Wilcoxon signed rank test.

**5.3.2 Model and Implementation details.** We used Scenario 2 for this task. For the bandit network, we used a feed forward layer with three highway network layers followed by an output layer with a sigmoid unit. The dimensions of the highway layer were set to 92. An input projection layer with an ReLU activation function was used to project the input vectors into 92 dimensions. The provided 46-dimensional feature vectors correspond to vectors  $c_i$  mentioned in Section 4.

We used the reward function defined in Eq. (11). We chose nDCG@10 instead of reciprocal rank (RR) for this task as the number of relevant documents was large. We optimized the model using the Adam optimizer [19] with the beta parameters set to  $(0, 0.999)$ , and a weight decay of  $1e^{-6}$  was used for regularization. We used the hybrid training objective defined in Eq. (12) with  $\gamma$  tuned over the set of values  $[0.25, 0.5, 0.75, 1]$ . The best results were obtained for BanditRank with  $\gamma = 0.5$ . A dropout of 0.4 was applied to all layers. An initial learning rate of  $7e^{-5}$  was used for all models. We set  $M', B$  to  $M' = 40$  and  $B = 30$  for calculating the gradient in Eq. (8). A high value was chosen for  $M'$  because the number of queries with at least 30 relevant documents was 99, which is a significant number. Moreover, the average number of relevant documents per query was large, i.e., 10.3.

**5.3.3 Baselines and Results.** We compared BanditRank with four strong listwise baselines. The baselines were AdaRank [51], ListNet [6], Coordinate Ascent [24] and the state-of-the-art listwise ranking method LambdaMART [4]. All baselines were implemented using the RankLib <sup>8</sup> software. As mentioned in Section 3.4, hybrid training objective with  $\gamma = 0.5$  resulted in the best performance as BanditRank with  $\gamma = 1$  cannot efficiently explore the relatively large action space in this task.

The results given in Table 6 show that BanditRank clearly outperformed AdaRank, ListNet, and Coordinate Ascent. When compared with the stronger baseline LambdaMART, except for the measures P@10, nDCG@3, and nDCG@10, BanditRank achieved minimum of 1% improvement in all other measures. Except for the measure nDCG@10, the improvement shown by BanditRank on all other measures is statistically significant according to the paired t-test and wilcoxon signed rank test. In the next section, we discuss the

<sup>8</sup><https://sourceforge.net/projects/lemur/>

**Table 6: Results of MQ2007 dataset. Best results are in bold. Statistically significant differences compared to best model according to paired t-test is denoted as \* and wilcoxon signed rank test is denoted as + (p-value < 0.05).**

	P@1	P@3	P@10	MAP
ListNet	0.446**	0.409**	0.366**	0.452**
AdaRank	0.474**	0.434**	0.379**	0.471**
Coordinate Ascent	0.474**	0.435**	0.382**	0.474**
LambdaMART	0.477**	0.444*	0.390**	0.477**
BanditRank( $\gamma = 1$ )	0.460	0.432	0.382	0.468
BanditRank( $\gamma = 0.5$ )	<b>0.498</b>	<b>0.457</b>	<b>0.393</b>	<b>0.483</b>
	nDCG@1	nDCG@3	nDCG@10	MRR
ListNet	0.391**	0.392**	0.435	0.556**
AdaRank	0.432**	0.426**	0.457	0.577**
Coordinate Ascent	0.418**	0.420**	0.449	0.574**
LambdaMART	0.431**	0.434**	0.470*	0.582**
BanditRank( $\gamma = 1$ )	0.412	0.413	0.454	0.572
BanditRank( $\gamma = 0.5$ )	<b>0.447</b>	<b>0.437</b>	<b>0.473</b>	<b>0.597</b>

behavior of different reward functions when trained on the MQ2007 dataset.

#### 5.4 Comparison of Reward Functions

The reward function plays a significant role in the training of an agent in reinforcement learning. Gradual feedback through rewards is often required for training a good agent. For comparing the behavior of reward functions during training, we conducted experiments using different reward functions on the MQ2007 dataset with the same architecture as the previous section. Since we wanted to compare the behavior of the reward function, we chose  $\gamma = 1$  for all the experiments. The following reward functions were used:

$$R_1(a_c, g_c) = AP(a_c, g_c)$$

$$R_2(a_c, g_c) = nDCG@10(a_c, g_c)$$

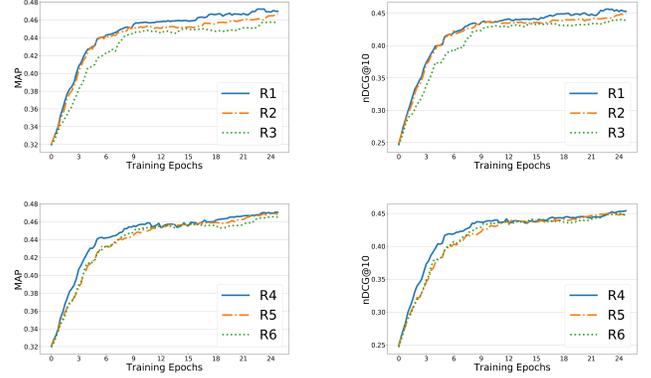
$$R_3(a_c, g_c) = DCG@5(a_c, g_c)$$

$$R_4(a_c, g_c) = \frac{[AP + nDCG@10](a_c, g_c)}{2}$$

$$R_5(a_c, g_c) = \frac{[AP + RR](a_c, g_c)}{2}$$

$$R_6(a_c, g_c) = \frac{[AP + P@3 + P@5 + nDCG@3 + nDCG@5](a_c, g_c)}{5}$$

The first three reward functions are the direct evaluation measures and the last three are a combination of several evaluation measures. Figure 1 plots the test set performance of the models during training epochs. We can observe that  $R_1$  achieved good measures right from the start when compared with  $R_2$  and  $R_3$ . The performance of  $R_6$  when compared with  $R_5$  and  $R_4$  shows that using many evaluation measures will not necessarily improve the measures of MAP and nDCG@10. During the initial epochs,  $R_4$  was clearly a better performer than  $R_5$  and  $R_6$ . After a similar performance by all three models in the next few epochs,  $R_4$  achieved better measures of MAP and nDCG@10 in the final stages.



**Figure 1: Test set performance of BanditRank trained using different reward functions on MQ2007 dataset. Top row corresponds to functions  $R_1, R_2$  and  $R_3$  and bottom row corresponds to functions  $R_4, R_5$  and  $R_6$ .**

## 6 CONCLUSION

We proposed an extensible listwise deep learning method *BanditRank* for ranking. It can directly optimize the evaluation measures using the policy gradient algorithm. Experimental results indicate the superiority of BanditRank over other methods on the tested datasets. Future work can involve modifying the structure of the policy network discussed in Section 3.1 for efficiently addressing the issue of exploration when the number of actions is large. For example, we could use adaptive exploration strategies instead of simple  $\epsilon$ -greedy strategy, for exploring the action space. We can define new reward functions for handling queries with no relevant documents. For example, we can penalize the model if any of the document affinity scores for such queries is greater than 0.5. There is also a possibility of defining reward functions as the weighted average of different measures with trainable weights for better feedback. Regarding the theoretical aspects, we can compare the directness of BanditRank to other algorithms such as LambdaRank.

## REFERENCES

- [1] Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2016. An actor-critic algorithm for sequence prediction. *arXiv preprint arXiv:1607.07086* (2016).
- [2] Michael Bendersky, Donald Metzler, and W Bruce Croft. 2010. Learning concept importance using a weighted dependence model. In *Proceedings of the third ACM international conference on Web search and data mining*. ACM, 31–40.
- [3] Weijie Bian, Si Li, Zhao Yang, Guang Chen, and Zhiqing Lin. 2017. A compare-aggregate model with dynamic-clip attention for answer selection. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. ACM, 1987–1990.
- [4] Christopher JC Burges. 2010. From ranknet to lambdarank to lambdamart: An overview. *Learning* 11, 23–581 (2010), 81.
- [5] Christopher J Burges, Robert Ragno, and Quoc V Le. 2007. Learning to rank with nonsmooth cost functions. In *Advances in neural information processing systems*. 193–200.
- [6] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*. ACM, 129–136.
- [7] Olivier Chapelle, Quoc Le, and Alex Smola. 2007. Large margin optimization of ranking measures. In *NIPS workshop: Machine learning for Web search*.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

- [9] Yue Dong, Yikang Shen, Eric Crawford, Herke van Hoof, and Jackie Chi Kit Cheung. 2018. Banditsum: Extractive summarization as a contextual bandit. *arXiv preprint arXiv:1809.09672* (2018).
- [10] Pinar Donmez, Krysta M Svore, and Christopher JC Burges. 2009. On the local optimality of LambdaRank. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 460–467.
- [11] Minwei Feng, Bing Xiang, Michael R Glass, Lidan Wang, and Bowen Zhou. 2015. Applying deep learning to answer selection: A study and an open task. In *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*. IEEE, 813–820.
- [12] John Guiver and Edward Snelson. 2008. Learning to rank with softrank and gaussian processes. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 259–266.
- [13] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. A deep relevance matching model for ad-hoc retrieval. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 55–64.
- [14] Hua He and Jimmy Lin. 2016. Pairwise word interaction modeling with deep neural networks for semantic similarity measurement. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 937–948.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [16] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. 2014. Convolutional neural network architectures for matching natural language sentences. In *Advances in neural information processing systems*. 2042–2050.
- [17] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*. ACM, 2333–2338.
- [18] Sumeet Katariya, Branislav Kveton, Csaba Szepesvari, and Zheng Wen. 2016. DCM bandits: Learning to rank with multiple clicks. In *International Conference on Machine Learning*. 1215–1224.
- [19] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [20] Branislav Kveton, Csaba Szepesvari, Zheng Wen, and Azin Ashkan. 2015. Cascading bandits: Learning to rank in the cascade model. In *International Conference on Machine Learning*. 767–776.
- [21] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436.
- [22] Gyoung Ho Lee and Kong Joo Lee. 2017. Automatic Text Summarization Using Reinforcement Learning with Embedding Features. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. 193–197.
- [23] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. *Text Summarization Branches Out* (2004).
- [24] Donald Metzler and W Bruce Croft. 2007. Linear feature-based models for information retrieval. *Information Retrieval* 10, 3 (2007), 257–274.
- [25] Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural variational inference for text processing. In *International conference on machine learning*. 1727–1736.
- [26] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Jingfang Xu, and Xueqi Cheng. 2017. DeepRank: A new deep architecture for relevance ranking in information retrieval. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. ACM, 257–266.
- [27] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 311–318.
- [28] Romain Paulus, Caiming Xiong, and Richard Socher. 2017. A deep reinforced model for abstractive summarization. *arXiv preprint arXiv:1705.04304* (2017).
- [29] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
- [30] Tao Qin and Tie-Yan Liu. 2013. Introducing LETOR 4.0 datasets. *arXiv preprint arXiv:1306.2597* (2013).
- [31] Tao Qin, Tie-Yan Liu, and Hang Li. 2010. A general approximation framework for direct optimization of information retrieval measures. *Information retrieval* 13, 4 (2010), 375–397.
- [32] Filip Radlinski, Robert Kleinberg, and Thorsten Joachims. 2008. Learning diverse rankings with multi-armed bandits. In *Proceedings of the 25th international conference on Machine learning*. ACM, 784–791.
- [33] Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732* (2015).
- [34] Jinfeng Rao, Hua He, and Jimmy Lin. 2016. Noise-contrastive estimation for answer selection with deep neural networks. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 1913–1916.
- [35] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7008–7024.
- [36] Cicero dos Santos, Ming Tan, Bing Xiang, and Bowen Zhou. 2016. Attentive pooling networks. *arXiv preprint arXiv:1602.03609* (2016).
- [37] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. 2015. Highway networks. *arXiv preprint arXiv:1505.00387* (2015).
- [38] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. 2000. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*. 1057–1063.
- [39] Ming Tan, Cicero Dos Santos, Bing Xiang, and Bowen Zhou. 2016. Improved representation learning for question answer matching. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vol. 1. 464–473.
- [40] Yi Tay, Luu Anh Tuan, and Siu Cheung Hui. 2018. Multi-cast attention networks for retrieval-based question answering and response prediction. *arXiv preprint arXiv:1806.00778* (2018).
- [41] Michael Taylor, John Guiver, Stephen Robertson, and Tom Minka. 2008. Softrank: optimizing non-smooth rank metrics. In *Proceedings of the 2008 International Conference on Web Search and Data Mining*. ACM, 77–86.
- [42] Shengxian Wan, Yanyan Lan, Jun Xu, Jiafeng Guo, Liang Pang, and Xueqi Cheng. 2016. Match-srnn: Modeling the recursive matching structure with spatial rnn. *arXiv preprint arXiv:1604.04378* (2016).
- [43] Bingning Wang, Kang Liu, and Jun Zhao. 2016. Inner attention based recurrent neural networks for answer selection. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vol. 1. 1288–1297.
- [44] Jun Wang, Lantao Yu, Weinan Zhang, Yu Gong, Yinghui Xu, Benyou Wang, Peng Zhang, and Dell Zhang. 2017. Irgan: A minimax game for unifying generative and discriminative information retrieval models. In *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM, 515–524.
- [45] Shuohang Wang and Jing Jiang. 2016. A compare-aggregate model for matching text sequences. *arXiv preprint arXiv:1611.01747* (2016).
- [46] Zhiguo Wang, Wael Hamza, and Radu Florian. 2017. Bilateral multi-perspective matching for natural language sentences. *arXiv preprint arXiv:1702.03814* (2017).
- [47] Zhiguo Wang, Haitao Mi, and Abraham Ittycheriah. 2016. Sentence similarity learning by lexical decomposition and composition. *arXiv preprint arXiv:1602.07019* (2016).
- [48] Zeng Wei, Jun Xu, Yanyan Lan, Jiafeng Guo, and Xueqi Cheng. 2017. Reinforcement learning to rank with Markov decision process. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 945–948.
- [49] Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* 8, 3-4 (1992), 229–256.
- [50] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [51] Jun Xu and Hang Li. 2007. Adarank: a boosting algorithm for information retrieval. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 391–398.
- [52] Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. Wikiqa: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 2013–2018.
- [53] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- [54] Fajie Yuan, Guibing Guo, Joemon M Jose, Long Chen, Haitao Yu, and Weinan Zhang. 2016. Lambdafm: learning optimal ranking with factorization machines using lambda surrogates. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 227–236.
- [55] Yisong Yue, Thomas Finley, Filip Radlinski, and Thorsten Joachims. 2007. A support vector method for optimizing average precision. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 271–278.
- [56] Wei Zeng, Jun Xu, Yanyan Lan, Jiafeng Guo, and Xueqi Cheng. 2018. Multi Page Search with Reinforcement Learning to Rank. In *Proceedings of the 2018 ACM SIGIR International Conference on Theory of Information Retrieval*. ACM, 175–178.
- [57] Weinan Zhang, Tianqi Chen, Jun Wang, and Yong Yu. 2013. Optimizing top-n collaborative filtering via dynamic negative item sampling. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*. ACM, 785–788.
- [58] Xiangyu Zhao, Long Xia, Jiliang Tang, and Dawei Yin. 2018. Reinforcement Learning for Online Information Seeking. *arXiv preprint arXiv:1812.07127* (2018).