
LEARNING TO RANK FOR PUSH NOTIFICATIONS USING PAIRWISE EXPECTED REGRET

Yuguang Yue*, Yuanpu Xie, Huasen Wu, Haofeng Jia, Shaodan Zhai, Wenzhe Shi, Jonathan J Hunt*
{yuguangy, yxie, huasenw, hjia, szhai, wshi, jjh}@twitter.com

ABSTRACT

Listwise ranking losses have been widely studied in recommender systems. However, new paradigms of content consumption present new challenges for ranking methods. In this work we contribute an analysis of learning to rank for personalized mobile push notifications and discuss the unique challenges this presents compared to traditional ranking problems. To address these challenges, we introduce a novel ranking loss based on weighting the pairwise loss between candidates by the expected regret incurred for misordering the pair. We demonstrate that the proposed method can outperform prior methods both in a simulated environment and in a production experiment on a major social network.

1 INTRODUCTION

The majority of internet users now access the internet via a mobile device (Handley, 2019). As a result, mobile devices have become increasingly important for content consumption. Push notifications allow a mobile user to receive timely messages and relevant information from an app they have installed, even while the app is in the background.

Because push notifications interrupt the user flow, users have a low tolerance for irrelevant content and expect to receive notifications only for content that is timely and important to them. This makes ranking and decision making for push notifications crucial to the user experience and a demanding problem.

Learning to rank is an important topic in information retrieval for designing machine learning ranking systems to surface relevant content to users from a large corpus (Liu, 2011). These ranking models are typically learned using logged user feedback on previous content. Listwise ranking is an approach to modeling the ranking problem with a loss which more closely approximates the utility of a ranking over documents to a user. For example, in a traditional “10 blue links” search paradigm (Chen et al., 2012), the user is only shown the top 10 results. Therefore, the ordering of documents after the 10th will have no impact on the utility of the ranking, a property a listwise ranking loss can exploit. Many metrics and listwise ranking losses implicitly focus on this search paradigm.

In this work, we present the ranking problem for push notifications and its unique challenges on a the Twitter platform. While some specifics may be unique to the platform, the problem constraints will be similar in other content-based push settings.

We observe that push notification differ from many classical ranking problems in several key ways. 1. Unlike other ranking problems such as search, the user will observe only one notification from each set of candidates. 2. User behavior is highly personalized since there is no explicit context from the user to indicate their information need. 3. User responses are highly non-stationary, a notification that is timely and relevant now, may be irrelevant if sent later. 4. The candidates available to be sent to the user changes rapidly and is highly personalized, so the approach must generalize to rank documents never seen before. These characteristics combine to make counterfactual estimates challenging. See section 3.1 for more details.

We derive a ranking loss which approximately minimizes the expected regret of the ranking in the push notification setting by weighting pairwise losses. We test this approach in both a simulation based on the production system and a production experiment on real users. We show that the

*These authors contributed equally.

proposed loss can outperform prior approaches and results in significant gains over baselines in a production AB test.

2 RELATED WORK

The ranking problem is to sort a candidate set of documents by relevance to the user. There has been a significant body of work studying ranking losses and metrics known as learning to rank (Oosterhuis et al., 2020). A number of metrics have proposed to estimate performance of a ranking. For example, the widely used Normalized Discounted Cumulative Gain (Järvelin and Kekäläinen, 2002), Expected Reciprocal Rank (Chapelle et al., 2009), and Average Relevance Position (Zhu, 2004). These ranking metrics are non-differentiable and therefore not practical to optimize directly. Their surrogate functions can be optimized via gradient based algorithms (Cao et al., 2007; Xia et al., 2008).

Pairwise losses are a class of ranking losses that compare positive and negative pairs to minimize the number of inversions (Chen et al., 2009). In practice, Pasumarthi et al. (2019) empirically demonstrates that *listwise* losses typically outperform *pairwise* losses on a range of tasks, because of *pairwise* losses do not take ranking information into account.

Many listwise losses can be decomposed into weighted pairwise losses. Well known ranking losses include RankNet (Burgess et al., 2005), LambdaRank (Burgess et al., 2006) which improves on RankNet by adding weights to the gradient based on the cost of inversions, a tree-based LambdaMART (Wu et al., 2010).

To take advantage of the ranking information, Usunier et al. (2009) proposed the ordered weighted pairwise classification (OWPC) algorithm. However, as with many *listwise* losses (e.g. *ListNet* (Xia et al., 2008) and *SVM_{map}* (Yue et al., 2007)), OWPC requires knowledge of the rank of positive examples to calculate the loss, which is prohibitively expensive when the candidate set has millions of items. To make this idea more scalable Weston et al. (2013) proposed the *K-Order Statistics Loss* (K-OS) loss which is an extended loss function of the weighted approximately ranked pairwise loss (Weston et al., 2011); K-OS loss takes care of both the estimation of rank and reweighting the pairwise loss. Weighted Margin-Rank Batch loss proposes a similar method for the batch training regime (Liu and Natarajan, 2017).

In addition to the learning to rank approaches to push notifications, two recent papers (Yancey and Settles, 2020; Xu et al., 2020) proposed contextual bandit based algorithms for push notification and personalized recommendation problems, where the arm of the bandit is interpreted as the item to be sent. However, both algorithms require repeated arms, which are not available in our setting where content creation happens rapidly.

3 PUSH NOTIFICATION RANKING PROBLEM

Changes in technology have resulted in new modes of users consuming content. In particular, on mobile devices users can choose to receive “push notifications” to be alerted to relevant new content in a timely manner. We describe the basic setup of the recommender based notification system at Twitter.

Algorithm 1 shows pseudocode of the notification system. \mathcal{U} denotes the set of users and \mathcal{X} the space of features of documents. Periodically a pass is made through the set of users, candidate documents for each user are obtained and a scoring function, parametrized by θ , $f_\theta(u, x) : \mathcal{U} \times \mathcal{X} \rightarrow \mathbb{R}$ is used to score the candidates. The highest scoring document is sent to the user and the user response (which is binary, either the user opens the notification $y = 1$ or ignores it $y = 0$) is logged and used to train subsequent ranking models.

The objective of the ranking system is to surface documents relevant to the user, as measured by the user choosing to open the notifications rather than dismiss them.

Algorithm 1 Pseudo-code of a single pass of the ranking based push notification system. At each timestep a set of candidates for a user are obtained (note this list will be distinct for each user and each pass through the loop due to new content being created), ranked and the highest scoring candidate sent to the user.

```

Accept scoring function  $f_\theta(u, x)$ 
for  $u \in \mathcal{U}$  do
  Obtain candidate documents  $\{x_1, \dots, x_n\}$  for user  $u$  available this time.
  Find the highest scoring document  $x = \arg \max_{x \in \{x_1, \dots, x_n\}} f_\theta(u, x)$ 
  Send the corresponding document to the user and receive a response  $y \in \{0, 1\}$ 
  Log  $u, x, y$ 
end for

```

3.1 CHARACTERISTICS OF PUSH NOTIFICATIONS

Push notifications have characteristics that make ranking challenging and distinct from many existing ranking problems. Firstly, of all the candidates ranked, only a single candidate can be sent to the user and receive feedback due to the display limitations of push notifications. This is important when considering what is the appropriate ranking loss below, and distinct from many other ranking problems.

Secondly, document relevance is highly personalized. Unlike some other content consumption routes, users receive push notifications without actively interacting with an app, so there is limited context (e.g. in a search engine if several users search for the same keywords, their responses could be grouped together). This high degree of personalization means that it is not valid to average user responses to the same document.

Thirdly, user responses are non-stationary. That is a document may be relevant to the user now, but not relevant a short time in the future. For example, users make use of Twitter to obtain information on breaking news (Petrovic et al., 2013; Osborne and Dredze, 2014). A document sent to a user now may be opened, the same document sent later may be irrelevant and dismissed later.

Finally, new documents are created at high volume. Therefore, ranking approaches must generalize to documents that have not been seen before. Each time a set of candidates is ranked, both this set and many individual candidates in the set may have never been seen by the system before.

These properties combine to make it challenging to estimate counterfactual outcomes. In contrast, in a search engine where the user views multiple results, after correcting for position bias (Craswell et al., 2008), it is then possible with some assumptions to estimate the performance of a ranking with a different ordering of the same documents (Agarwal et al., 2019). The personalization and non-stationarity problems combine to mean that we can't average multiple user responses together to get better relevance estimates than binary labels in contrast to typical contextual search problems (Joachims, 2002).

Like most ranking problems, we could also frame this as a contextual bandit problem as (see section 2). However, because of the rapid creation of new documents, standard bandit approaches would not be relevant since each arm of the bandit would only be seen once and we need to generalize to unseen arms. A contextual bandit in our problem setup would reduce to a pointwise loss. For this reason, in this work we have focused treating the problem as a ranking problem.

3.2 PUSH NOTIFICATION RANKING UTILITY

In this work we deal only with deterministic policies. Because in push notifications, only the top ranked document is sent to the user and all other documents have no effect on the outcome, we define a policy $\pi : \mathcal{U} \times \{\mathcal{X}_1, \dots, \mathcal{X}_n\}$ as a mapping from a user u and unordered candidate set of documents $\{x_1, \dots, x_n\}$ to one document in that set x . Note that the ordering over documents in the set is arbitrary. The chosen document $x = \pi(u, \{x_1, \dots, x_n\})$ results in a reward $r(u, x) = y$ based on the response of the user $y \in \{0, 1\}$, who may either open the document $y = 1$ or dismiss it $y = 0$.

The reward is stochastic, as there are many unobserved factors which affect the outcome (e.g. a relevant notification sent when the user is busy may be dismissed). In simulation we model each document as having a latent likelihood of being opened if it is sent to the user $\tilde{y} = p(y = 1|u, x)$ which defines a Bernoulli distribution. In the simulation, where we have access to the latent $p(y = 1|u, x)$ for each document, we also define a deterministic reward function $r_{sim}(u, x) = p(y = 1|u, x)$. Note that $r_{sim} = \mathbb{E}_{p(y|u, x)} r(u, x)$. This latent variable is never used in training (that would be unrealistic), but allows us to evaluate the ranking algorithms with lower variance in simulation.

The regret (which we can only compute exactly in simulation when we have access to the latent properties of the document) incurred by a ranking policy π for a particular user u and candidate set of documents x_1, \dots, x_n is given as:

$$g(u, \{x_1, \dots, x_n\}, \pi) = \max_{x_j \in \{x_1, \dots, x_n\}} r_{sim}(u, x_j) - r_{sim}(u, \pi(u, \{x_1, \dots, x_n\})). \quad (1)$$

The way we actually parametrized our ranking policy for all methods here is by a scoring function $f_\theta : \mathcal{U} \times \mathcal{X} \rightarrow \mathbb{R}$ that maps from user features and document features to a real-valued score, f_θ is parameterized by θ . The highest scoring document is selected $\pi(u, \{x_1, \dots, x_n\}) = \arg \max_{x \in \{x_1, \dots, x_n\}} f_\theta(u, x)$.

3.3 USER TYPES

Users at Twitter are categorized into 7 different user types based on their usage of the platform. The behavior of users to push notifications varies by user type (e.g. a user who uses Twitter every day is more likely to open a notification than an occasional user). We made use of the user type in our ranking method described below.

4 METHOD

Here we describe both the prior methods for ranking that we use as baselines and introduce a novel loss designed for ranking push notifications.

We write all the losses over a single candidate set. We split the set of candidates into candidates with positive labels (opened) \mathcal{X}_{pos} and candidates with negative labels (dismissed notifications) \mathcal{X}_{neg} . We assume the final output of the model $f_\theta(u, x)$ has no non-linearity and we write any final non-linearity in the loss itself.

4.1 EXISTING RANKING LOSSES

4.1.1 POINTWISE

A straightforward approach, widely used in practice (e.g. [Liu et al., 2017](#)), is to treat the ranking problem as a classification problem on the outcome y and use the cross-entropy loss to train the score function with empirical risk minimization to predict the likelihood that a document will be opened. The loss is:

$$\ell_{ce}(u, \theta) = \sum_{x_{pos} \in \mathcal{X}_{pos}} -\log(\sigma(f_\theta(u, x_{pos}))) - \sum_{x_{neg} \in \mathcal{X}_{neg}} \log(1 - \sigma(f_\theta(u, x_{neg}))) \quad (2)$$

where a sigmoid function σ is used to bound the prediction $f_\theta(u, x)$ to $(0, 1)$.

The well-known weakness of a pointwise loss is that it is a poor proxy of the true objective of the ranking function ([Liu, 2011](#)) because of the mismatch between the objective of predicting the outcome for each document, and the ranking problem in which only the score relative to other documents in the candidate set are important. This is particularly acute in the push notification setting where only a single item will be sent to the user and therefore only the item which is ranked highest matters for the outcome.

4.1.2 PAIRWISE

Pairwise losses focus on the relative ordering of two documents (with different outcomes). The pairwise hinge loss is

$$\ell_{pair}(u, \theta) = \sum_{x_{pos} \in \mathcal{X}_{pos}} \sum_{x_{neg} \in \mathcal{X}_{neg}} \max(0, 1 - (f_{\theta}(u, x_{pos}) - f_{\theta}(u, x_{neg}))) \quad (3)$$

The hinge loss “pushes” the score function to score documents with a positive label higher than negative documents. There are other pairwise losses that could be used, such as a pairwise logistic loss, but in initial experiments we found a hinge loss performed best.

4.1.3 LISTWISE RANKING LOSSES

Listwise ranking losses try to approximate the utility of the ordering of the set more directly. The closest prior example to the loss we introduce below is the *K-Order Statistics Loss* (K-OS) AUC loss (Weston et al., 2013) which incorporates ranking information into a pairwise loss.

K-OS first computes an ordering over the positive examples based on the current scores of the model being optimized $(x_{pos}^1, \dots, x_{pos}^{|\mathcal{X}_{pos}|})$.

$$\ell_{K-OS-AUC}(u, \theta) = \frac{1}{Z} \sum_{i=1}^{|\mathcal{X}_{pos}|} W\left(\frac{i}{|\mathcal{X}_{pos}|}\right) \sum_{x_{neg}} \max(0, 1 - (f_{\theta}(u, x_{pos}^i) - f_{\theta}(u, x_{neg}))) \quad (4)$$

where $Z = \sum_i W\left(\frac{i}{|\mathcal{X}_{pos}|}\right)$. W weights pairs by their rank. If $W(j) = C$ for all j and C is a positive constant, K-OS loss just reduces to the pairwise loss; if $W(i) > W(j)$ when $i < j$, it will penalize errors on high ranked positive items more. In our setup, $W\left(\frac{i}{|\mathcal{X}_{pos}|}\right) = 1$ for $i = 1$ and 0 otherwise to reflect that only the highest ranked document is important to the outcome in push notifications. We also tested a “weight capping” approach for comparison with the loss introduced (see equation 8) below where $W\left(\frac{i}{|\mathcal{X}_{pos}|}\right) = k < 1$ for $i > 1$, so that lower ranked positive examples still contribute to the loss, just at a reduced weight. We treated k as a hyperparameter and include $k = 0$ (where it reduces to the original K-OS loss) in the hyperparameter search space. In Weston et al. (2013) equation 4 was estimated by sampling due the prohibitively large number of documents. In our case, each candidate set is of moderate size, so we compute the loss exactly.

4.2 PSEUDO-CANDIDATE SETS

Applying pairwise or listwise losses to push notifications is not straightforward due to the previously discussed requirement that only a single document from each candidate set can be sent to the user. Since only one item in the set is labelled, we must find some way to group items for pairwise or listwise comparisons.

The approach we used for all the experiments outlined below was to load a batch of candidates from the training log, and group examples by user type. This approach was used for all pairwise and listwise losses in the experiments. Appendix A outlines some alternative approaches we considered.

4.3 EXPECTED REGRET

Here we introduce the ranking loss that is the primary contribution of this work. The motivation of the *expected regret* (ER) loss is to weight pairwise losses by the regret that is expected to be incurred if only this pair was misordered in the ranking. The loss is over pseudo-candidate sets.

$$\ell_{er}(u, \theta) = \sum_{x_{pos}} \sum_{x_{neg}} w_{er}(x_{pos}, x_{neg}) \times \max(0, 1 - (f_{\theta}(u, x_{pos}) - f_{\theta}(u, x_{neg}))) \quad (5)$$

The weighting of each pairwise loss $w_{er}(x_{pos}, x_{neg})$ is the crucial aspect. We introduce the approach by assuming that for each document x_{pos} we have three additional pieces of information: the probability of the user opening the document $\tilde{y} = p(y = 1 | x_{pos}, u)$ (note $\tilde{y} \in [0, 1]$ is not binary),

the cumulative distribution function for \tilde{y} for this candidate set (we denote as $F(\tilde{y})$) and the number of candidates in the candidate set n . Knowledge of \tilde{y} is unrealistic, and we will remove later. However, $F(\tilde{y})$ can be estimated from logged data and n is known.

We compute w_{er} for a pair by using the information to compute the expected regret in CTR (Click Through Rate) incurred if the model misorders only this pair in a candidate set. Since we send only one candidate, the regret incurred for misordering anything but the top ranked candidate is 0 so this can be decomposed into two parts: the probability x_{pos} is the top ranking candidate and the expected regret is x_{neg} is sent rather than x_{pos} .

The probability that a candidate with a CTR of \tilde{y}_{pos} is the top ranked candidate in a candidate set of size n is:

$$p_{top}(\tilde{y}_{pos}) = (1 - F(\tilde{y}_{pos}))^{n-1} \quad (6)$$

If the document x_{pos} should have been the top-ranked document and we misorder and send the user x_{neg} instead, then we will incur a regret in expected CTR of $\tilde{y}_{pos} - \tilde{y}_{neg}$. Therefore the expected regret of misranking a pair is

$$w'_{er} = p_{top}(\tilde{y}_{pos}) \times (\tilde{y}_{pos} - \tilde{y}_{neg}) \quad (7)$$

$$w_{er} = \max(w'_{er}, k) \quad (8)$$

and we bound the weights so they cannot be negative or completely ignore any examples.

4.3.1 REMOVING THE CTR ASSUMPTION

The ER loss proposed above requires using \tilde{y} when weighting the pairs. Obviously, assuming knowledge of the latent value of a document is unworkable, if we knew that we would know the perfect ranking. Here we remove that assumption by observing even if the pointwise losses don't provide good ranking performance, for the purposes of weighting pairs, we only need to estimate \tilde{y} approximately. The estimated CTR is used only for weighting, the labels are used for generating pairs so errors in the CTR estimate will result in weighting a pair incorrectly but won't result in the loss ranking a negative document above a positive one.

To avoid doubling the training time, by requiring first training a pointwise loss and then a ER loss, we used the same model and train with both losses simultaneously. Combining pairwise and pairwise losses has previously been successful with logistic regression (Sculley, 2010; Li et al., 2015). However, it was crucial to ensure the losses are approximately compatible between the pairwise and pointwise loss. We used an ℓ_2 pointwise loss and used the label values 1, -1 for the ℓ_2 loss so that for any pair x_{pos}, x_{neg} there is a solution that minimizes both the pairwise and pointwise loss. We add these losses together with a hyperparameter α to control their relative contribution to the final loss.

Algorithm 2 puts everything together to show how these losses can be combined to simultaneously learn online both a CTR estimate, and using this estimate to weight the pairwise loss (on the same model).

Algorithm 2 Pseudo-code of the (unbatched) ER loss.

Initialize θ

for training iterations **do**

 Sample pair of examples $(u_{pos}, x_{pos}), (u_{neg}, x_{neg})$

 Estimate CTRs $\tilde{y}'_{pos} = f_{\theta}(u_{pos}, x_{pos}), \tilde{y}'_{neg} = f_{\theta}(u_{neg}, x_{neg})$.

 Compute w_{er} using equation 7 and the estimated values $\tilde{y}'_{pos}, \tilde{y}'_{neg}$.

 Compute ER loss ℓ_{er} using equation 5

 Sample pointwise sample u_i, x_i, y_i .

 Compute pointwise loss ℓ_2 .

 Update θ to minimize loss

$$\ell = \ell_{er} + \alpha \ell_2 \quad (9)$$

end for

5 EXPERIMENTS

We compared the Expected Regret loss with the pointwise loss, pairwise loss, K-OS-AUC loss. We treated the capping weight k as a hyperparameter for K-OS loss. For all losses the same neural network architecture, dataset, and feature pre-processing were used, and models were trained using stochastic gradient descent with Adam optimizer (Kingma and Ba, 2014) in simulation and plain SGD in the production experiment.

5.1 SIMULATION

We tested Expected Regret loss (along with baselines) in an production experiment (next section) and a simulated experiment using synthetic data which we describe here.

Simulations of recommender systems have gained popularity as systems to facilitate experimental setups beyond simple pointwise losses (Rohde et al., 2018; Ie et al., 2019a;b; Mladenov et al., 2021). We used the recsim framework (Ie et al., 2019a) to construct the simulation, and we fit key parameters of the simulation to production data in order to make the simulation more realistic.

We used an extremely simple model of users as being fully described by their user type, sampled from a Categorical distribution matching the production system. Conditioned on the user type a candidate set consisting of 60 documents is sampled at each interaction. Each document has an underlying probability of being opened $\tilde{y} = p(y = 1|x, u)$ which is drawn independently and identically from a Beta distribution fit (using maximum likelihood estimation) to the candidate distribution of the production system. For fitting the distributions, the production ranking system estimates were taken as ground truth. Each document has a corresponding set of features x , which in our case is a 5 degree projection of $p(y = 1|x, u)$ with independent Gaussian noise added to each dimension. Although a very simple model of features, this captures the essence that features provide a model with a noisy estimate of the underlying CTR of the document. A model is used to score the documents and select which one is sent to the user, and a label is generated by sampling from the Bernoulli distribution defined by the latent CTR \tilde{y} associated with the document. For evaluation we used the latent probability of open to compute the regret without any noise due to the label sampling (see section 3.2).

5.2 DATA BIAS

The ranking models are trained on feedback from users. A well-known problem in ranking is the data bias (Wang et al., 2016; Levine et al., 2020) issue this introduces, since feedback will only be received on documents that are sent to a user, typically documents which rank highly under an existing ranker. We don't propose to solve this common problem here. Ideally, we might explore by sending the users documents chosen uniformly at random over all possible candidates. However, this would result in a poor user experience, so a common trade-off is ϵ -greedy exploration, where for a small ϵ fraction of requests documents are sent at random, and for the remainder the top ranked document is sent.

In simulation we test the ranking losses with both "unbiased" data, which simulates sending documents to a user chosen at random and a biased production dataset. We simulated the production data bias by ranking documents according to a pointwise ranker (trained on an second set of simulated data) with ϵ -greedy exploration with ϵ of 0.14, approximately the same as the production data.

5.3 ONLINE PRODUCTION EXPERIMENTS

For the production experiment all models were trained offline on the same logged production data which was obtained using ϵ -greedy exploration $\epsilon \approx 0.14$ ¹ and the production ranking model. The models were tested by serving a subset of Twitter users for push notifications. The experiment ran for approximately two weeks. Users were kept in the same "bucket" throughout the experiment, receiving notifications from only one of the models. The hyperparameters of the model can be found in Appendix B.

¹This is not the true value of ϵ , rather non-exploration notifications are downsampled so this is a description of the data.

Model	unbiased	gain	Model	biased	gain
Pointwise	0.07624 ± 0.00006	0.0%	Pointwise	0.07417 ± 0.00005	0.0%
Pairwise	0.07623 ± 0.00006	0.01%	Pairwise	0.07412 ± 0.00004	0.07%
K-OS-AUC	0.07679 ± 0.00016	-0.72%	K-OS-AUC	0.08139 ± 0.00119	-9.73%
Expect Regret	0.07602 ± 0.00003	0.28%	Expect Regret	0.07429 ± 0.00008	-0.16%

Table 1: **The regret of different ranking losses on simulated data.** The units of regret are arbitrary, and lower is better. Errorbars show SEM from 10 runs. Gains show the % improvement in the regret compared to the pointwise loss. The left panel is evaluated on unbiased data and the right panel is on biased data.

6 RESULTS

The results for the simulated data are in Table 1. The key finding is that on unbiased data expected regret outperforms all the alternative approaches. However, on biased data expected regret performs slightly worse than a simple pointwise or pairwise loss (but still outperforms the baseline listwise ranking loss K-OS). The poorer performance on biased data is probably due to the expected regret method weighting high-ranking examples heavily. The data bias already results in over-representing high-ranking examples and this may be resulting in an extreme overweighting on high ranked examples.

We also compared the methods in a large-scale AB testing experiment on Twitter. Due to the work incurred in production experiments and its poor performance in simulation, we did not include the K-OS loss and its variant. We did compare against pointwise and pairwise ranking models.

Table 2 shows the results of the production experiment. Despite the issues of bias in simulation, expected regret clearly outperformed other methods in the real-world simulation. It resulted in both a statistically significant improvement in the open rate and users interacted more with the opened documents, indicating the documents being recommended were of more interest to them.

One reason the data bias issue may affect the production experiment less is that the simulated data has symmetric Gaussian noise, which means the performance of a model is very sensitive to precise weighting over the data, whereas this is not true of the production experiments.

Model	Open	Interact
Pointwise	0.0%	0.0%
Pairwise	0.04%	0.46%
Expected regret	0.33%	2.46%

Table 2: **Online experiment results (bold number for differences $p < 0.01$).** The table shows % gain in open and interaction rate compared with the baseline pointwise model (the pointwise model performance is 0% by definition). The expected regret model resulted in a higher open rate compared to both a pointwise and pairwise model. Although not directly optimized for, the interact rate, which indicates a user subsequently interacted with a document after opening was also significantly improved. This indicates the improved ranking also resulted in users being more engaged with the content.

7 CONCLUSION

In this work we have outlined the characteristics of ranking for push notifications that make it challenging and distinct from other ranking problems. We introduced a novel *expected regret* (ER) ranking loss designed for push notifications, which weights a pairwise loss to minimize the expected regret in the ranking. We compared ER against prior approaches in both a simulation and a real-world production setting. In both cases, we found significant improvements in performance over prior methods.

REFERENCES

- Aman Agarwal, Kenta Takatsu, Ivan Zaitsev, and Thorsten Joachims. 2019. A general framework for counterfactual learning-to-rank. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 5–14.
- Christopher Burges, Robert Ragno, and Quoc Le. 2006. Learning to rank with nonsmooth cost functions. *Advances in neural information processing systems* 19 (2006), 193–200.
- Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Huelender. 2005. Learning to rank using gradient descent. In *Proceedings of the 22nd international conference on Machine learning*. 89–96.
- 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*. 129–136.
- Olivier Chapelle, Donald Metzler, Ya Zhang, and Pierre Grinspan. 2009. Expected reciprocal rank for graded relevance. In *Proceedings of the 18th ACM conference on Information and knowledge management*. 621–630.
- Danqi Chen, Weizhu Chen, Haixun Wang, Zheng Chen, and Qiang Yang. 2012. Beyond ten blue links: enabling user click modeling in federated web search. In *Proceedings of the fifth ACM international conference on Web search and data mining*. 463–472.
- Wei Chen, Tie-Yan Liu, Yanyan Lan, Zhi-Ming Ma, and Hang Li. 2009. Ranking measures and loss functions in learning to rank. In *Advances in Neural Information Processing Systems*. 315–323.
- François Chollet et al. 2015. Keras. <https://keras.io>.
- Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. An experimental comparison of click position-bias models. In *Proceedings of the 2008 international conference on web search and data mining*. 87–94.
- Lucy Handley. 2019. Nearly three quarters of the world will use just their smartphones to access the internet by 2025. <https://www.cnbc.com/2019/01/24/smartphones-72percent-of-people-will-use-only-mobi>
- Eugene Ie, Chih-wei Hsu, Martin Mladenov, Vihan Jain, Sanmit Narvekar, Jing Wang, Rui Wu, and Craig Boutilier. 2019a. Recsim: A configurable simulation platform for recommender systems. *arXiv preprint arXiv:1909.04847* (2019).
- Eugene Ie, Vihan Jain, Jing Wang, Sanmit Navrekar, Ritesh Agarwal, Rui Wu, Heng-Tze Cheng, Morgane Lustman, Vince Gatto, Paul Covington, et al. 2019b. Reinforcement learning for slate-based recommender systems: A tractable decomposition and practical methodology. *arXiv preprint arXiv:1905.12767* (2019).
- Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)* 20, 4 (2002), 422–446.
- Thorsten Joachims. 2002. Optimizing search engines using clickthrough data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*. 133–142.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. 2020. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643* (2020).
- Cheng Li, Yue Lu, Qiaozhu Mei, Dong Wang, and Sandeep Pandey. 2015. Click-through prediction for advertising in twitter timeline. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1959–1968.

-
- David C Liu, Stephanie Rogers, Raymond Shiau, Dmitry Kislyuk, Kevin C Ma, Zhigang Zhong, Jenny Liu, and Yushi Jing. 2017. Related pins at pinterest: The evolution of a real-world recommender system. In *Proceedings of the 26th international conference on world wide web companion*. 583–592.
- Kuan Liu and Prem Natarajan. 2017. Wmrb: Learning to rank in a scalable batch training approach. *arXiv preprint arXiv:1711.04015* (2017).
- Tie-Yan Liu. 2011. Learning to rank for information retrieval. (2011).
- Martin Mladenov, Chih-Wei Hsu, Vihan Jain, Eugene Ie, Christopher Colby, Nicolas Mayoraz, Hubert Pham, Dustin Tran, Ivan Vendrov, and Craig Boutilier. 2021. RecSim NG: Toward Principled Uncertainty Modeling for Recommender Ecosystems. *arXiv preprint arXiv:2103.08057* (2021).
- Harrie Oosterhuis, Rolf Jagerman, and Maarten de Rijke. 2020. Unbiased learning to rank: counterfactual and online approaches. In *Companion Proceedings of the Web Conference 2020*. 299–300.
- Miles Osborne and Mark Dredze. 2014. Facebook, Twitter and Google Plus for breaking news: Is there a winner?. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 8.
- Rama Kumar Pasumarthi, Sebastian Bruch, Xuanhui Wang, Cheng Li, Michael Bendersky, Marc Najork, Jan Pfeifer, Nadav Golbandi, Rohan Anil, and Stephan Wolf. 2019. Tf-ranking: Scalable tensorflow library for learning-to-rank. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2970–2978.
- Sasa Petrovic, Miles Osborne, Richard McCreadie, Craig Macdonald, Iadh Ounis, and Luke Shrimpton. 2013. Can twitter replace newswire for breaking news?. In *Proceedings of the international AAAI conference on web and social media*, Vol. 7.
- David Rohde, Stephen Bonner, Travis Dunlop, Flavian Vasile, and Alexandros Karatzoglou. 2018. Recogym: A reinforcement learning environment for the problem of product recommendation in online advertising. *arXiv preprint arXiv:1808.00720* (2018).
- Sebastian Ruder. 2016. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747* (2016).
- David Sculley. 2010. Combined regression and ranking. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. 979–988.
- Nicolas Usunier, David Buffoni, and Patrick Gallinari. 2009. Ranking with ordered weighted pairwise classification. In *Proceedings of the 26th annual international conference on machine learning*. 1057–1064.
- Xuanhui Wang, Michael Bendersky, Donald Metzler, and Marc Najork. 2016. Learning to rank with selection bias in personal search. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. 115–124.
- Jason Weston, Samy Bengio, and Nicolas Usunier. 2011. Wsabie: Scaling up to large vocabulary image annotation. (2011).
- Jason Weston, Hector Yee, and Ron J Weiss. 2013. Learning to rank recommendations with the k-order statistic loss. In *Proceedings of the 7th ACM conference on Recommender systems*. 245–248.
- Qiang Wu, Christopher JC Burges, Krysta M Svore, and Jianfeng Gao. 2010. Adapting boosting for information retrieval measures. *Information Retrieval* 13, 3 (2010), 254–270.
- Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. Listwise approach to learning to rank: theory and algorithm. In *Proceedings of the 25th international conference on Machine learning*. 1192–1199.
- Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. 2015. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853* (2015).

-
- Xiao Xu, Fang Dong, Yanghua Li, Shaojian He, and Xin Li. 2020. Contextual-Bandit Based Personalized Recommendation with Time-Varying User Interests. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 6518–6525.
- Kevin P Yancey and Burr Settles. 2020. A Sleeping, Recovering Bandit Algorithm for Optimizing Recurring Notifications. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3008–3016.
- 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*. 271–278.
- Mu Zhu. 2004. Recall, Precision and Average Precision.

A ALTERNATIVE APPROACHES TO PSEUDO-CANDIDATE SETS

As outlined in section 4.2, in order to apply pairwise or listwise losses we grouped candidates set to the user into pseudo-candidate sets. Here we discuss two alternative approaches to constructing pseudo-candidate sets we considered.

The fundamental challenge is that, due to requirement to send a single document from a candidate set, the log data consists of user candidate pairs $(u_1, x_1, y_1), (u_2, x_2, y_2), \dots$ which differ in both the user attributes u and document x and user responses y . However, when ranking, we are considering various documents for a single user, so the candidate set will have the same user attributes for every document $(u, x_1), (u, x_2), \dots$. We can’t get this log because we can only send the user a single notification at a time, so we do not have the labels y for the other documents.

One approach would be to load a batch of examples and treat them all as a candidate set. However, one of the most important features that predict if a user will open a notification is their “user type,” which classifies a user by how actively they engage with the platform. Grouping all users in a batch would encourage the model to learn to simply pick out very active users and rank documents sent to them highly which would not work well for online ranking, where candidate documents are all being sent to the same user.

Another approach we tried is to group all notifications sent to the same user over some period of time together. However, we found this performed very poorly (data not shown) because the model learns to “cheat” by making use of user features that change based on their responses (e.g. features indicating if the user has been active on the platform recently). Concretely, even though document x_1 and x_2 were sent the same user, some user attributes can change in the interval between notifications so the log data looks like $(u, x_1), (u', x_2)$ where u' is the same user with some features changed. The model can exploit those feature changes to predict which document was opened in a pairwise loss by just focusing on the user features. This results in good offline performance on the pseudo-candidate sets, but abysmal performance when used for ranking over real candidate sets.

B HYPERPARAMETERS

For both simulation and online experiment, we use $k = 0.001$ for the weight capping in equation 8 and $\alpha = 0.3$ for the weight of pointwise loss equation 9.

Simulation: For the simulation, we use a single set of neural network structure and hyperparameters for a fair comparison. The simple MLP has a size of $(64, 32)$ and sigmoid activation functions. We use the Adam optimizer (Kingma and Ba, 2014) with learning rate 0.001 and rest parameters are default by the Keras API (Chollet et al., 2015). The batch size for all methods is 512 and the early stopping patience is 5.

Online Experiment: On top of the extracted features, we used a three-layer MLP with size $(256, 128, 64)$ and the activation function is LeakyReLU (Xu et al., 2015). We used a SGD optimizer (Ruder, 2016) with momentum 0.99999. For each method we use cross-validation to choose an optimal learning rate (from $\{1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}\}$) and batch size (from $\{64, 128, 256, 512, 1024\}$).