

1 Recommending items with attributes

2

3

4

5

6 Anonymous Author(s)

7 ABSTRACT

8

9 Recommendation models suggest items to users. We study a variant

10 of this task where each item has a set of attributes, such as

11 tags on an image, user reactions to a post, or foods in a meal. We

12 propose RANKFROMSETS, a flexible and scalable model for recom-

13 mending items with attributes. rfs treats item attributes as side

14 information and learns embeddings to discriminate items a user

15 will consume from items a user is unlikely to consume. We de-

16 velop theory connecting the rfs loss to optimal recall and show

17 that the learnable class of models for rfs is a superset of several

18 previously-proposed models. To scale rfs to large datasets, we pro-

19 pose a stochastic optimization framework based on negative sam-

20 pling. Within this framework, we suggest two specific approaches

21 that trade off computation for accuracy. In experiments on two

22 real-world datasets, rfs outperforms competing approaches while

23 also learning embeddings that reveal interpretable structure in user

24 behavior.

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

1 INTRODUCTION

29 Classical recommender system datasets contain a matrix where
30 each row is a user and each column is an item. Each entry in the ma-
31 trix indicates whether or not the user consumed the item. Modern
32 applications often gather rich side information about items in the
33 form of a set of attributes or tags. Intuitively, item attributes pro-
34 vide valuable side information for recommender systems. When
35 the number of items is large or the user-item matrix is sparse, lever-
36 aging the attribute information is crucial to achieving good per-
37 formance.

38 In practice, modeling item attributes in a recommender system is
39 not straightforward. Standard multiple-matrix factorization meth-
40 ods may fail if the attribute vocabulary or item set is too large.
41 At the same time, oversimplifying the problem risks sacrificing
42 the ability to capture nonlinear statistical dependencies between
43 attributes. For instance, a user may enjoy items tagged with at-
44 tributes A and B, B and C, or A and C, but not all three. Finding
45 the right balance between scalability and flexibility is therefore a
46 primary goal.

47 Even when a model can be scaled, it is not always clear how the
48 training procedure connects to the recommender system evalua-
49 tion metric. It may be that a model with optimal loss leads to sub-
50 optimal recommender system performance. For example, a matrix
51 factorization method may minimize mean squared error but the
52 recommender system may be evaluated on precision. While it is
53 plausible that minimizing error will improve precision, the con-
54 nection between the two is implicitly assumed in many methods.
55 Ideally, a well-designed recommender system should have an opti-
56 mization loss that matches its evaluation metric.

59 This paper proposes RANKFROMSETS (rfs), a scalable model for rec-
60 ommending items with attributes. rfs casts the recommendation
61 problem as binary classification. Given an item and a user, rfs
62 treats attributes as features and classifies whether or not the item
63 is likely to be consumed by the user. rfs learns embeddings for
64 each user and attribute; each item is represented as the mean of
65 its attribute embeddings. To scale to large datasets, rfs is trained
66 using a stochastic optimization procedure that randomly samples
67 items that are unlikely to be consumed.

68 rfs enjoys two benefits from framing the recommendation prob-
69 lem as classification. First, the rfs classification objective function
70 is directly tied to recommender recall, and we show that a classi-
71 fier with zero worst-case error achieves maximum recall. Second,
72 the model is provably flexible enough to learn any class of rec-
73 ommendation model based on set-valued side information. Thus,
74 the rfs model generalizes other models while retaining scalabil-
75 ity.

76 We study rfs on two datasets. We consider a dataset of 65k users
77 clicking on 637k papers posted to the arXiv; the attributes of a pa-
78 per are the unique words in its abstract. rfs outperforms a state-of-
79 the-art content-based matrix factorization model [7] on this dataset,
80 and the learned document embeddings capture the structure of
81 human-labeled topics. Second, we consider a large dataset of 55k
82 users logging 16M meals using the LoseIt! diet tracking app. Again,
83 rfs outperforms other baseline models, and the learned meal em-
84 beddings reveal a diverse set of recommendations that capture the
85 notion of taste.

2 BACKGROUND

86 We highlight two themes in research on recommendation models.
87 First, we describe several recommendation models that incorpo-
88 rate side information and note why they do not scale to datasets
89 with many items. We also describe models that optimize proxies
90 of ranking metrics and discuss why these proxies are not directly
91 tied to optimal recommendation evaluation metrics.

2.1 Recommendation with side information

95 Side information is included in recommendation models in sev-
96 eral ways; we focus on deep learning and matrix factorization ap-
97 proaches. Item side information can be modeled with deep repre-
98 sentations, or can be included in content-based matrix factoriza-
99 tion models as an additional matrix.

100 *Deep representations of side information.* Deep learning-based rec-
101 ommendation models incorporate side information in multiple ways
102 [34]. For example, items that have words as attributes can be rep-
103 resented using recurrent neural networks [1]. Lian et al. [14] use
104

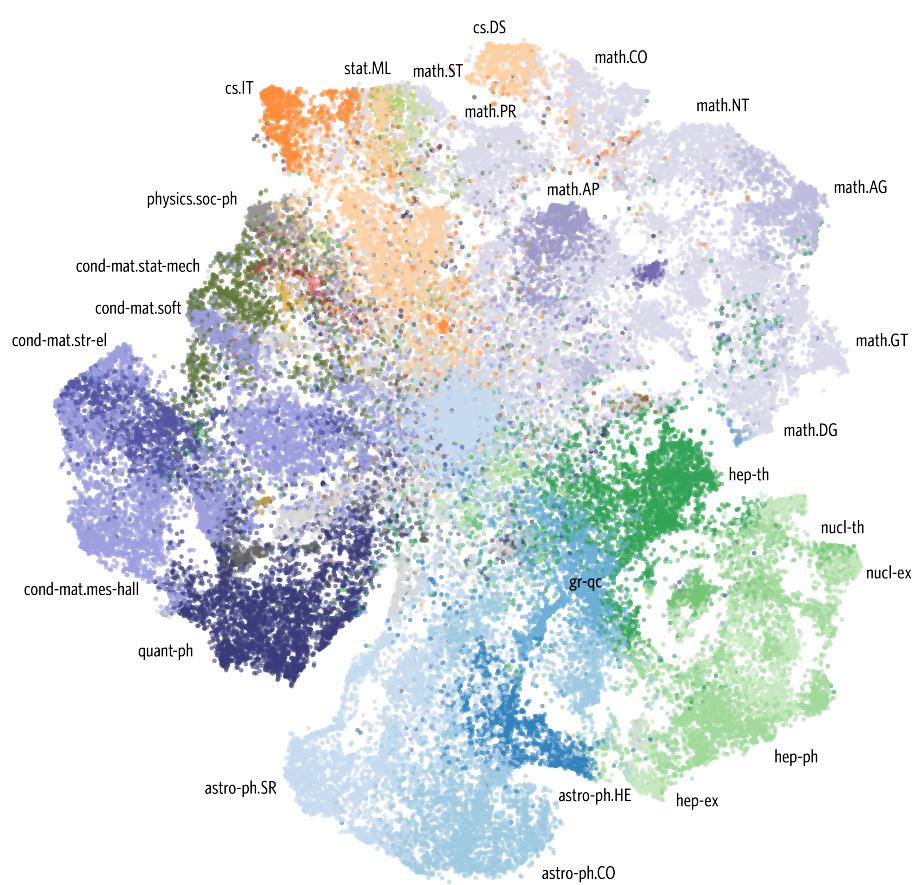


Figure 1: RANKFROMSETS trained on arXiv reading behavior clusters researchers by their most frequently-read arXiv category (best viewed in color). On this data, RANKFROMSETS is trained to recommend items using their attributes as described in Section 4.2; the set of attributes for a paper is the set of words in the abstract. A visualization of the user embeddings in the inner product regression function in Equation (4) yields an interpretable map of science. We use the t-SNE algorithm [18] for visualizing the high-dimensional embeddings in two dimensions. In this map of science, fields of study are related according to patterns in how people read papers in neighboring fields. Each marker represents a user embedding; the color assigned to a user is determined the user’s most-read arXiv category. The color assigned to a category is determined by the most-read categories across the arXiv, with similar colors assigned to similar fields according to the arXiv ontology. For an interactive version of this map, please visit the anonymous URL <https://app.box.com/v/rankfromsets>. The interactive version enables zoom and display of all 143 arXiv category labels to explore the relationships between different fields of science.

an attention mechanism to weight recommendations according to available item and user side information, and Dong et al. [6] use denoising autoencoders to model side information in a deep recommendation model, but requires fitting parameters for every item. An example of a more efficient approach is the method in Chen et al. [5], where embeddings are jointly learned for users, items, and item text for recommendation, but this method focuses on unsupervised pre-training of text representations. There are several examples of ‘tag-aware’ or ‘tag-based’ deep recommendation models [16, 36], such as Xu et al. [29], which focuses on data where users and items have different attributes and use autoencoders to learn user, item, and attribute representations. They use a cosine

similarity-based objective function which is not tied to a metric used to evaluate recommendation performance.

Matrix factorization with side information. Shi et al. [23] survey several matrix factorization methods that incorporate side information. Gopalan et al. [7] develop a Bayesian matrix factorization model for recommending items based on side information in the form of words in documents. Wang and Blei [28] develop a regression model that uses a topic model to incorporate side information into recommendations. There are also several ‘tag-based’ or ‘tag-aware’ content-based matrix factorization models [2, 17, 35]. Such

content-based matrix factorization methods maximize the conditional log-likelihood of the data (or a bound on the log-likelihood); optimizing these objective functions does not optimize an evaluation metric. Furthermore, all of these methods are not scalable to large numbers of items as they require learning unique parameters for every item. Specifically, such content-based matrix factorization methods require learning a matrix that has a row for every item. For items with attributes, it is often infeasible to store this matrix in memory or exploit efficient coordinate ascent optimization schemes that require processing this entire matrix.

2.2 Recommendation via ranking

Recommendation models can be trained on loss functions that approximate ranking-based evaluation metrics. We describe recommendation models trained on such loss functions and extensions that include side information.

Learning to rank. The literature on learning to rank includes models that optimize proxies of evaluation metrics, such as mean average precision, mean reciprocal rank, or discounted cumulative gain [15, 31]. For example, Bayesian personalized ranking models optimize a pairwise ranking objective function [21] that trains the model to rank items a user consumed higher than items a user did not consume. This objective is a heuristic motivated by an analogy to the receiver operating characteristic; a model trained on this objective does not provably maximize this metric. Song et al. [24] extend Bayesian personalized ranking using deep neural networks, but do not model side information.

Learning to rank with side information. Models that optimize proxies of ranking metrics that use side information include Shi et al. [22], where a smoothed approximation of mean average precision is used as a loss function. Yuan et al. [32] use a proxy of a ranking loss to fit a polynomial that models predictions of item consumption using item and side information features. Ying et al. [30] uses denoising autoencoders to represent item information in a model trained with a pairwise ranking loss. Cao et al. [4] use a ranking loss to jointly learn embeddings of items and attributes; they focus on the case where users interact directly with both attributes and items with said attributes. All of these models require learning unique parameters for every item, and do not scale to large numbers of items. An example of a scalable method that uses the Bayesian personalized ranking criterion is in Okura et al. [20], but this approach requires data with timestamps and negative item feedback.

3 THE RANKFROMSETS MODEL

RANKFROMSETS is a recommendation model that recommends items with attributes to users. Let $u \in \{1, \dots, N\}$ be a user, $m \in \{1, \dots, M\}$ be an item, and $y_{um} \in \{0, 1\}$ be a binary indicator where 1 indicates user u consumed item m . For each item m , there is an associated set of attributes $x_m \in \{0, 1\}^{|V|}$ from a vocabulary of V attributes.

We assume that a recommendation model is given a budget of K recommendations to be made for each user. In response, the recommender system produces a list of K distinct recommendations $\mathbf{r}_u = (r_{u1}, \dots, r_{uK})$ for each user. The goal of the recommendation task in this paper is to maximize the expected recall,

$$\text{Recall}@K = \mathbb{E}_u \left[\frac{\sum_{r \in \mathbf{r}_u} y_{ur}}{\sum_m y_{um}} \right], \quad (1)$$

where the expectation is over all users in the empirical data distribution \mathcal{D} . RANKFROMSETS (RFS) combines three techniques to maximize Equation (1). First, we cast recommendation as a classification task. Second, we learn user- and attribute-level embeddings. Statistical strength is shared between items with similar attributes by representing items as the mean of their attribute embeddings. Third, we scale RFS to large datasets by using a stochastic optimization-based negative sampling training procedure to fit the model.

3.1 Recommendation through classification

RANKFROMSETS (RFS) casts the recommendation problem as a classification task. Given a user-item pair (u, m) , RFS learns to predict the probability that item m will be consumed by user u ,

$$p(y_{um} = 1 | u, m) = \sigma(f(u, x_m)), \quad (2)$$

where x_m is the set of attributes associated with item m and σ is the sigmoid function. The recommendations made by RFS are the maximum likelihood set,

$$\mathbf{r}_u(K) = \operatorname{argmax}_{\mathbf{r} \in \mathbb{N}^K} \sum_{m \in \mathbf{r}} f(u, x_m). \quad (3)$$

To motivate treating recommendation as classification, we make the following observation.

PROPOSITION 1. *Let $u \in \mathcal{U}$ be a user, $x \in \mathcal{X}$ be an item, and $y(u, x) \in \{0, 1\}$ be an indicator of whether user u logged item x . Let \mathcal{E} be the worst-case error for binary classifier $\hat{y}(u, x)$ on any (u, x) pair drawn from the data \mathcal{D} ,*

$$\mathcal{E} = \max_{(u, x) \in \mathcal{D}} 1[\hat{y}(u, x) \neq y(u, x)].$$

A binary classifier with zero worst-case error ($\mathcal{E} = 0$) maximizes recommendation recall.

PROOF. A model with zero worst-case error is a perfect classifier: it assigns greater probability to positively-labeled datapoints than to negatively-labeled datapoints. In other words, it ranks positive examples above negative examples. Recall in Equation (1) is measured by the fraction of positively-labeled items in a ranking returned by the model. In a classifier that achieves zero worst-case error, positively-labeled datapoints must be ranked higher than other datapoints, maximizing recall. \square

Proposition 1 is simple, but conceptually important. Under the assumption that a perfect classifier exists, a consistent method for learning a classifier will be a consistent method for learning a recommendation system that targets expected recall. In practice, as

with any regression method, a perfect classifier is unachievable. Proposition 1 is a guiding principle rather than a finite-sample guarantee of maximal performance. As we show in Section 4, the classification approach of RFS performs well in practice.

3.2 Embedding users and items

For recommending items with attributes, Proposition 1 says that building a classifier such as RANKFROMSETS is optimal if we measure recommendation performance with recall. To parameterize the RFS classifier, we need to choose a regression function $f(u, x_m)$. A straightforward parameterization is an inner product,

$$f(u, x_m) = \theta_u^\top \left(\frac{1}{|x_m^+|} \sum_{j \in x_m^+} \beta_j + g(x_m) \right) + h(x_m). \quad (4)$$

Each element in the inner product regression function in Equation (4) has an intuitive interpretation:

- The user embedding $\theta_u \in \mathbb{R}^d$ captures the latent preferences for user u . This captures the individual-level tastes of a user and is analogous to the user preference vector in classical collaborative filtering or the row embedding in matrix factorization.
- The attribute embedding $\beta_j \in \mathbb{R}^d$ is the latent quality conveyed through item m having attribute j . The set x_m^+ contains only attributes with $x_{mj} = 1$. Attributes that are not associated with item m are ignored.
- The item embedding function $g(x_m) \in \mathbb{R}^d$ captures qualities not conveyed through the set of item attributes. This term in the regression function enables collaborative filtering by capturing unobserved patterns in item consumption such as popularity. We describe how to construct this function below.
- The item intercept function $h(x_m) \in \mathbb{R}$ makes an item more or less likely simply due to availability.

Scalable item embedding and item intercept functions. The parameterization of the item embedding function $g(x_m)$ depends on the size of the data. If the number of items is small, g can function as a lookup for unique intercepts for every item. However, if the number of items is so large that unique item intercepts lead to overfitting, we need a scalable parameterizations of item embeddings g using additional information about every item. For example, if the data consists of foods in meals, we can define a meal intercept as the mean of food intercepts, yielding a scalable item intercept function. The item intercept function $h(x_m)$ that maps item attributes to scalars is constructed in the same way. We study both of these choices in Section 4.

The inner product regression function in Equation (4) has several benefits. It requires computing a sum over only the attributes each item is associated with. This enables RFS to scale to large attribute vocabularies where traditional matrix factorization methods are intractable. Second, the embed-and-sum approach to set modeling is provably flexible. We describe deep variants of RFS and detail how RFS can approximate other recommendation models in the next section.

3.3 Deep extensions

The RFS inner product regression function in Equation (4) is a log-bilinear model. But there are several other choices of regression function, and we draw on the deep learning toolkit for classification to build two other example architectures. With finite data and finite compute, one architecture will outperform another, or prove insufficiently flexible to capture patterns in user consumption. The optimal choice of architecture depends on computational tradeoffs. We conclude this section by showing that in the regime of infinite data and compute, all the RFS architectures we propose can approximate other recommendation models that operate on set-valued input such as matrix factorization.

Parameterizing RFS using a neural network. As an alternative to the log-bilinear model in Equation (4), we can use a deep neural network as a regression function:

$$f(u, x_m) = \phi \left(\theta_u, \frac{1}{|x_m^+|} \sum_{j \in x_m^+} \beta_j, g(x_m) \right) + h(x_m), \quad (5)$$

where the deep network ϕ has weights and biases and takes as inputs the user embedding, sum of attribute embeddings, and item intercept. Such a neural network can represent functions that may or may not include the inner product in Equation (4) (*ex ante*, it is unclear whether a finite-depth, finite-width neural network can represent the inner product).

Parameterizing RFS with a residual network. Another possibility of regression function for RFS is a combination of Equations (4) and (5), using an idea borrowed from deep residual networks for image classification [9]. In this architecture a neural network ϕ with the same inputs as in Equation (5) is used to learn the residual in the inner product model:

$$f(u, x_m) = \theta_u^\top \left(\frac{1}{|x_m^+|} \sum_{j \in x_m^+} \beta_j + g(x_m) \right) + \phi + h(x_m). \quad (6)$$

The choice of regression function in RFS depends on the data. On finite data, with finite compute, one parameterization of RFS will outperform another. To demonstrate this, we simulated synthetic data from the same generative process RFS employs with a ground-truth regression function (a square kernel) where the residual and deep parameterizations outperformed the inner product architecture. These results are included in Appendix D, and motivate exploring other architectures than the three examples we give.

3.4 Generalization property

If we take a step back from the setting of finite data and compute, a bigger picture emerges, which reveals the choice of regression function in RFS does not matter. Any RFS architecture is sufficiently flexible to approximate other recommendation models that operate on set-valued input. Before applying the theoretical results derived in Zaheer et al. [33] to RFS, we first describe the class of recommendation models that operate on set-valued input.

465 *RFS and content-based matrix factorization are examples of permutation-*
 466 *invariant models.* The regression function f in RFS operates on set-
 467 valued input: the unordered collection of item attributes x_m^+ . A set
 468 is, by definition, permutation-invariant: it remains the same if we
 469 permute its elements. Functions that operate on set-valued inputs
 470 must also be permutation-invariant. RFS is permutation-invariant
 471 as the set of attributes associated with an item enter into Equa-
 472 tions (4) to (6) via summation. Other examples of permutation-
 473 invariant models are matrix factorization, recommendation mod-
 474 els based on word embeddings, and permutation-marginalized re-
 475 current neural networks. We show that these models are permutation-
 476 invariant in Section 4.1, and evaluate them in Section 4.
 477

478 *RFS can approximate other permutation-invariant models such as*
 479 *matrix factorization.* We apply a theorem from Zaheer et al. [33]
 480 to show that RFS can approximate other recommendation mod-
 481 els.

482 **PROPOSITION 2.** *Assume the vocabulary of attributes (set elements)*
 483 *is countable, $|V| < |\mathbb{N}_0|$. RFS can approximate any permutation-*
 484 *invariant recommendation model.*

485 The proof follows directly from Theorem 2 in Zaheer et al. [33]
 486 and we will not restate it here. (The only change to the proof is the
 487 mapping from set elements to one-hot vectors, $c: V \rightarrow \{0, 1\}^{|V|}$ to
 488 yield a unique representation of every object in the powerset.)

489 Proposition 2 means that any of the parameterizations in Equations
 490 (4) to (6) are flexible enough to approximate other prin-
 491 cipled recommendation models that leverage item attributes, such as
 492 content-based matrix factorization [7]. This proposition also sup-
 493 ports exploring other parameterizations of RFS that may have bet-
 494 ter computational or statistical properties.

495 3.5 Stochastic optimization

500 Denote the parameters of the RANKFROMSETS model by θ , and let
 501 \mathcal{D} be the empirical data distribution. The model parameters are
 502 the user embeddings, attribute embeddings, parameters of the item
 503 embedding g and item intercept h functions, and any weights and
 504 biases of the neural networks in Equations (5) and (6). The maxi-
 505 mum likelihood objective function is

$$506 \mathcal{L}(\theta, \mathcal{D}, \lambda_u) = \mathbb{E}_u [\mathbb{E}_{(x_m, y_{um}) \sim \mathcal{D} \mid u, y_{um}=1} [\log p(y_{um} = 1 \mid x_m; \theta)] \\ (7)$$

$$510 - \lambda_u \mathbb{E}_{(x_k, y_{uk}) \sim \mathcal{D} \mid u, y_{uk}=0} [\log p(y_{uk} = 0 \mid x_k; \theta)] \\ (8)$$

$$(9)$$

514 The parameter λ_u is a reweighting parameter for negatively-labeled
 515 datapoints, and it is derived as follows. This classification objective
 516 function derives from an assumption about the data-generating
 517 process. We assume that every user-item interaction is unique and
 518 every item is uniquely determined by its attributes. A user u does
 519 or does not consume an item m ; there is no noise in this process.
 520 This means there is a single label for every datapoint (x_m, y_{um}) .
 521 This assumption means that a classifier trained on Equation (7) is

522 consistent: in the infinite data limit, any reweighting of λ_u will pre-
 523 serve the ordering of items, as there is only a single label.
 524

525 *Stochastic optimization and per-user negative sample balance.* We
 526 use stochastic optimization to maximize Equation (7). In practice,
 527 the number of negatively-labeled datapoints is large, and we find
 528 that balancing the positive and negative labels per-user helps opti-
 529 mization. We describe two ways to fix λ_u below, using two nega-
 530 tive sampling schemes that are dependent on the choice of evalua-
 531 tion metric.

532 *User balance via corpus sampling.* Negative samples can be drawn
 533 uniformly over the entire corpus of items. If the item set is large,
 534 this can be an expensive procedure. However, in related applica-
 535 tions [27], corpus sampling was shown to substantially outperform
 536 other sampling strategies. This negative sampling scheme and re-
 537 sulting approximation of Equation (7) is similar to the objective
 538 functions in other recommender systems [10, 24].

539 On large datasets, it is infeasible to calculate recall for evaluation,
 540 as Equation (1) requires ranking every item for every user (e.g. in
 541 Section 4.3 we study a dataset with over 10M items). We define
 542 a scalable evaluation metric based on recall, and describe how it
 543 leads to a natural choice of negative sampling distribution.

544 *Sampled recall.* Start with held-out datapoints with positive labels,
 545 $(x_m, y_{um} = 1)$. Sampled recall is calculated as follows. For every
 546 held-out datapoint, $K-1$ datapoints with negative labels $(x_k, y_{uk} = 0)$ are
 547 sampled from the rest of the held-out data, which together
 548 yield a set of K datapoints. A recommendation model is used to
 549 rank the K datapoints r_{u1}, \dots, r_{uK} . Sampled recall is the fraction
 550 of the held-out datapoints that the model ranked in the top k :

$$551 \text{SampledRecall}@k = \mathbb{E}_{u,m} \left[\frac{\sum_{r \in \{r_{u1}, \dots, r_{uK}\}} y_{ur}}{K} \right], \quad (10)$$

552 where the expectation is over users and items in the held-out set
 553 of datapoints. This evaluation metric is scalable: instead of using
 554 a model to rank every item, it requires ranking K items. Sampled
 555 recall is 1 if $k = K$, as the held-out datapoints with $y_{um} = 1$ are
 556 in the list of K datapoints to be ranked.

557 *User balance via batch sampling.* When sampled recall is used as
 558 an evaluation metric, batch sampling is a natural way to draw neg-
 559 ative samples. Sampled recall is calculated on items drawn from
 560 other user’s data. We define batch sampling as generating negative
 561 samples by permuting mini-batch items. Besides corresponding to
 562 the sampled recall metric, this technique is memory-efficient, as it
 563 requires that only the current mini-batch be in memory.

564 In addition to scalability, both of the negative sampling procedures
 565 above have the advantage of implicitly balancing the classifier. As
 566 Veitch et al. [27] note, using stochastic gradient descent with nega-
 567 tive sampling is equivalent to a Monte Carlo approximation of the
 568 reweighted, or balanced, classification loss.

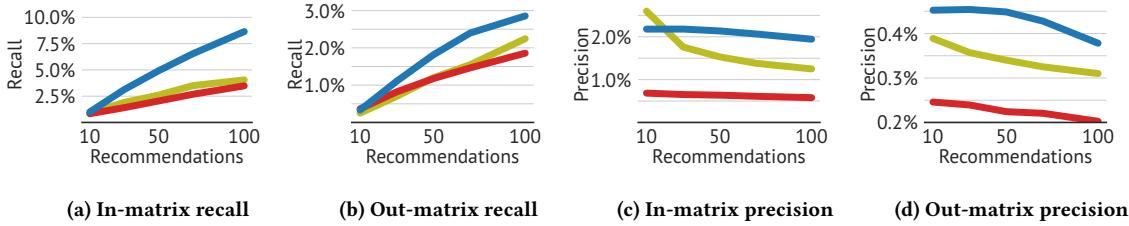


Figure 2: RANKFROMSETS with the inner product regression function in Equation (4) outperforms collaborative topic Poisson factorization (CTPF, described in Gopalan et al. [7]) and a word embedding model on recommending arXiv papers to scientists. In this dataset the items are documents and the attributes are the unique words in the abstracts. Recommendation performance is evaluated using precision and recall to match the evaluation in Gopalan et al. [7], and these metrics are reported on training (in-matrix) documents and cold-start (out-matrix) documents with no clicks in the training set.

4 EXPERIMENTS

We study the performance of RANKFROMSETS on two datasets. The first data consists of researcher reading behavior from the arXiv; the task is to recommend documents to scientists. The second is crowdsourced food consumption data from a diet tracking app, and the task is meal recommendation. On both benchmarks, RFS outperforms several baseline methods. An example implementation is in Appendix A and anonymized code for running experiments is at TODO (data cannot be shared, as releasing diet or arXiv data is associated with several privacy concerns).

4.1 Baselines

We compare RFS to a suite of recommendation models that rank items using their sets of attributes.

Word embedding model. We fit a word embedding model [19] on the training data using the implementation of fastText from Bojanowski et al. [3]. The dimensionality of the embeddings is set to match models being compared to. We use a context window in the model that is slightly different than in the original implementation. For an item with attributes x_m , the context for attribute $j \in x_m^+$ is the set of other attributes of the same item $j' \in x_m^+ : j' \neq j$. After learning the attribute embeddings β_j for $j \in V$, item embeddings are computed as the average of their attribute embeddings. Users are represented as the average of the embeddings of the items they consume. Recommendation is performed by cosine similarity of an item embedding and a user embedding.

Collaborative topic Poisson factorization. This is a probabilistic matrix factorization model of user consumption data. The recommendation model follows a latent variable scheme,

1. Document model:

- (a) Draw topics $\beta_{vk} \sim \text{Gamma}(a, b)$
- (b) Draw document topic intensities $\theta_{dk} \sim \text{Gamma}(c, d)$
- (c) Draw word count $w_{dv} \sim \text{Poisson}(\theta_d^T \beta_v)$.

2. Recommendation model:

- (a) Draw user preferences $\eta_{uk} \sim \text{Gamma}(e, f)$
- (b) Draw document topic offsets $\epsilon_{dk} \sim \text{Gamma}(g, h)$
- (c) Draw $r_{ud} \sim \text{Poisson}(\eta_u^T (\theta_d + \epsilon_d))$.

Gopalan et al. [7] develop a variational inference algorithm for posterior inference in this model; we use their implementation.

Collaborative topic Poisson factorization is a permutation-invariant recommendation model. To show that collaborative topic Poisson factorization [7] is permutation-invariant, consider the Poisson likelihood function over words w_{dv} . Conditional on the latent item representation θ_d and latent word representation β_v , every word in the document w_{dv} is independent; the joint probability of words in a document factorizes:

$$p(w_d | \theta_d, \beta_v) = \prod_{w_{dv} \in w_d} p(w_{dv} | \theta_d, \beta_v). \quad (11)$$

In this model, predictions are made using expectations under the posterior. The posterior is proportional to the log joint of the model, and the attributes of items (words in documents) enter into the model only via the above product. The product of the probability of words in a document is invariant to a reordering of the words in the document. Therefore, collaborative topic Poisson factorization is permutation-invariant.

Permutation-marginalized recurrent neural network. The Bernoulli distribution used in Equation (4) can be parameterized using a recurrent neural network. We treat attributes as a sequence and marginalize over all permutations of orderings,

$$p(y_{um} = 1 | x_m) = \frac{1}{|\pi(x_m)|} \sum_{\pi \in \pi(x_m)} \sigma \left(\phi(\theta_u, \{\beta_{\pi(1)}, \dots, \beta_{\pi(J)}\}) \right). \quad (12)$$

Here β_j are attribute embeddings, $\pi(x_m)$ denotes the set of all permutations of the set x_m , and ϕ is output of an LSTM recurrent neural network [11] projected to one dimension using a linear layer. This model is permutation-invariant as there is an explicit sum over all permutations.

697 4.2 Recommending research papers

698 We benchmark RFS on data of scientists reading research papers on
 699 the arXiv, where the goal is to recommend papers to scientists. The
 700 arXiv usage data represents one year of usage (2012) and consists
 701 of 65K users, 637K preprints, and 7.6M clicks.

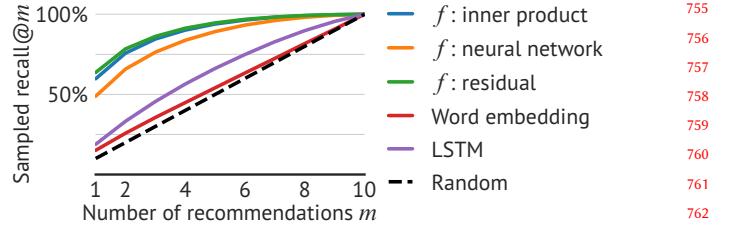
702
 703 *Evaluation.* We follow Gopalan et al. [7] and use the same test and
 704 validation splits, and the same set of held-out 10K users. To match
 705 Gopalan et al. [7], we compute precision in addition to recall. The
 706 held-out validation and test splits each consist of 20% of the ratings
 707 and 1% of the documents. In-matrix documents refer to documents
 708 that have clicks in the training data, while out-matrix or cold-start
 709 documents have no previous clicks.

710
 711 *Hyperparameters for RFS.* We test the stochastic gradient descent
 712 algorithm with and without momentum [25]. We use a linear learning
 713 rate decay that decays to zero in the maximum number of it-
 714 erations, 200k. We perform a grid search over learning rates of
 715 1, 5, 10, 15, 25 and momenta of 0.5, 0.9, 0.95, 0.99. The minibatch
 716 size is set to $\{2^{16}\}$. We use a single negative sample per datapoint,
 717 sampled uniformly over the entire dataset (corpus sampling is de-
 718 fined in Section 3.5). To match the hyperparameters in Gopalan
 719 et al. [7], we set the dimensionality of embeddings and intercepts
 720 to 100. Evaluation is performed every 20k iterations.

721
 722 *Hyperparameters for the permutation-marginalized recurrent neural
 723 network.* The embedding and hidden state sizes are fixed to 100.
 724 Evaluation is performed every 20k iterations. We grid search over
 725 learning rates of $10^{-1}, 10^{-2}, 10^{-3}$ with the Adam optimizer [12].
 726 If validation performance does not improve, we reload the best
 727 parameters and optimizer states, and decay the learning rate by
 728 0.9. We subsample single permutations during training and test-
 729 ing.

730
 731 *Results.* Figure 2 shows that RFS with the inner product parame-
 732 terization outperforms collaborative topic Poisson factorization in
 733 terms of in-matrix recall by over 90%. It also improves over collab-
 734 orative topic Poisson factorization (CTPF) in terms of out-matrix
 735 recall, out-matrix precision, and in-matrix precision (for the lat-
 736 ter, only when the number of recommendations is greater than
 737 30). The word embedding model performs about as well as CTPF
 738 in terms of recall, but performs worse in terms of precision. The
 739 permutation-marginalized recurrent neural network model per-
 740 formance was an order of magnitude worse than the other methods;
 741 we omit its performance.

742
 743 Qualitatively, RFS reveals patterns in usage of the arXiv. Figure 1
 744 is a dimensionality-reduced plot of the user embeddings that re-
 745 veals connections between fields of study. Scientists who focus on
 746 high energy physics, hep, neighbor specialists in differential geo-
 747 metry, math.DG; these areas share techniques. Machine learning re-
 748 searchers, stat.ML readers, neighbor statisticians or math.ST read-
 749 ers, highlighting the close connection between these fields. Plots
 750 for document embeddings show similar patterns. This illustrates
 751 how RFS captures rich patterns of interaction between users and
 752 items, while benefitting from information in the item attributes.



755
 756 **Figure 3: RANKFROMSETS outperforms models based on**
 757 **word embeddings permutation-marginalized recurrent neu-**
 758 **ral networks (denoted by LSTM) on meal recommendation.**
 759 **The recommendation models are trained on data from a food**
 760 **tracking app as described in Section 4.3 and are evaluated us-**
 761 **ing the sampled recall metric, Equation (10). The inner prod-**
 762 **uct, neural network, and residual regression functions for**
 763 **RANKFROMSETS are in Equations (4) to (6).**

764 4.3 Recommending meals

765 We evaluate RFS on data collected from the LoseIt! food tracking
 766 app. This app enables users to track their food intake to eat healthy.
 767 We use a year’s worth of data from 55k active users. This corre-
 768 sponds to 16M meals, where each meal is comprised of a subset of
 769 3M foods.

770
 771 *Preprocessing.* Each meal is a subset of 3M foods. We filter the vo-
 772 cabulary by keeping words that occur at least 20 times in the food
 773 names, resulting in 9963 words. A meal is represented as the union
 774 of the sets of words occurring in the food names.

775
 776 *Evaluation.* We hold out 1% of the items (meals) for validation and
 777 hold out 1% of the items for evaluating test performance. We use
 778 the sampled recall metric in Equation (10) with

779
 780 *Hyperparameters for RFS.* The embedding size is set to 128. For the
 781 neural network and residual models in Equations (5) and (6) the
 782 number of hidden layers is two, and the number of hidden units is
 783 set to 256. The item embeddings $g(x_m)$ (and item intercepts $h(x_m)$)
 784 are computed as the mean of learned food embeddings (intercepts).
 785 We use the RMSProp optimizer in Graves [8] and grid search over
 786 learning rates in $\{10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$. We use a batch size of
 787 64 and a single negative sample for every datapoint in a minibatch
 788 (batch sampling as defined in Section 3.5). Evaluation is performed
 789 every 50k iterations.

790
 791 *Hyperparameters for the permutation-marginalized recurrent neu-*
792 ral network. We grid search over batch sizes of $\{32, 64, 128\}$. For
 793 every item in a minibatch, we sample a single permutation of
 794 attributes to approximate the sum in Equation (12). We use a sin-
 795 gle negative sample per datapoint, and set the embedding and hid-
 796 den state sizes to 128. We use the Adam optimizer [12] and grid
 797 search over learning rates of $\{10^{-2}, 10^{-3}, 10^{-4}\}$. Evaluation is per-
 798 formed every 1k iterations. If the validation performance does not
 799 improve, we reload the best parameters and optimizer states, and
 800 decay the learning rate by 0.9.

Query meal	Nearest meal by cosine similarity	
Two scoops of Raisin Bran cereal, organic Moroccan green tea, almond milk, light honey, tap water, large banana, large strawberries	Vita Bee bread, salted butter, fresh medium tomatoes, large fried whole egg, small banana	871 872 873 874 875 876 877 878 879
Iceberg lettuce, cantaloupe cubes, diced honeydew melon, cherry tomatoes, olives, dry-cooked unsalted hulled sunflower seed kernels, chopped hard-boiled egg, cucumbers, dried cranberries, fat-free ranch dressing	Green leaf lettuce, chopped sweet red bell peppers, crumbled feta cheese, large hard-boiled egg, chopped cucumber, oil-roasted salted sunflower seeds, sliced radishes, sliced strawberries, pitted Calamata olives, fat-free balsamic vinegar	875 876 877 878 879
Boston roast pork, mackerel, artichoke hearts, spinach, pimiento-stuffed Manzanilla olives, carrots, mushrooms, pepperoni ranch dressing	Broiled top round steak, tomatoes, cucumber, baby yellow squash, zucchini, black olives, extra virgin olive oil	880 881 882 883 884 885 886
Meatloaf with tomato sauce, chopped sweet red bell peppers, extra virgin olive oil, cooked asparagus spears, sweet potatoes, orange, cantaloupe cubes	Chicken breast, breadcrumbs, fresh tomatoes, shredded green leaf lettuce, extra virgin olive oil, spinach, chopped yellow onion, sweet large yellow bell peppers, whole mushrooms, chili peppers, vinaigrette	882 883 884 885 886 887 888 889 890
Ciabatta bun, cooked skinless chicken breast, fresh baby spinach, shredded iceberg lettuce, shredded mozzarella cheese, ketchup, frozen yogurt bar	Small whole wheat submarine roll, broiled round roast beef, roasted light turkey meat without skin, fresh medium tomatoes, honey smoked ham, shredded iceberg lettuce, sliced mozzarella cheese	887 888 889 890

Table 1: **RANKFROMSETS**, trained on food consumption data, provides diverse recommendations. We fit the model to data from a food tracking app as described in Section 4.3; items are meals and attributes are the ingredients in the meal. We represent meals as the mean of their attribute embeddings, and use cosine similarity to compute the nearest neighbors of meals. This reveals that **RANKFROMSETS** uncovers latent patterns of consumption in the attribute embeddings that can be leveraged to improve recommendations. For example, the second-last query meal is a mix of meat, vegetables, and fruit, and the nearest neighbor meal is a different meat and a side of salad. The last query meal is a sandwich, and its nearest neighbor is also a sandwich, but with different ingredients.

Results. Figure 3 plots the sampled recall. This shows that RFS outperforms the permutation-marginalized recurrent neural network and the word embedding model. We tried to fit the collaborative topic Poisson factorization model in [7] to this dataset, but the code released by the authors was unable to scale to this size of data.

Qualitatively, RFS learns an interpretable representation of items, as shown by nearest neighbors of meals in Table 1. In this table, we display breakfast, lunch, and dinner meals, alongside their nearest neighbors. We find that the nearest neighbors are also breakfast, lunch, and dinner meals respectively, showing that the attribute embeddings learned by the model can be used to explore qualitative patterns in the learned latent space.

5 DISCUSSION

The task of recommending items with attributes is difficult for several reasons. It is unclear how to incorporate set-valued side information into models that scale to large numbers of items and attributes. In addition, existing recommendation models that leverage item attributes are not directly tied to evaluation metrics. We developed **RANKFROMSETS** (RFS), a scalable recommendation model for items with attributes. Theoretically, we showed that optimizing its loss function optimizes recall, and that it can approximate other well-founded recommendation models such as content-based matrix factorization. Empirically, RFS outperformed competitors and

scaled large datasets, and was very simple to implement as shown in Appendix A.

Simple theory and implementation. It is surprising that RFS outperformed the collaborative topic Poisson factorization model in Gopalan et al. [7], as RFS is a much simpler model. RFS is simpler than CTPF conceptually; it does not require approximate posterior inference, but mere binary classification. RFS is also simpler in terms of implementation: a scalable example is a few dozen lines of python in Appendix A, compared to several thousand lines of C++ released by the authors of CTPF.

Architectures. There is a wealth of deep learning models that can be used to parameterize RFS. For example, attention-based transformer networks can operate on set-valued input [13] and are an interesting choice of architecture for RFS.

Side information. The RFS can easily accommodate various metadata available about users, items, or attributes. Attribute-level or user-level side information is particularly interesting, and when available, should lead to improved recommendation performance.

Regularization. Further performance gains from regularizing RFS are also interesting avenues for future work. It would be interesting to test L2 regularization for sparsity and dropout for parsimony.

929 *Pretrained embeddings.* We note that the word embedding baseline
 930 performed almost as well as collaborative topic Poisson factoriza-
 931 tion, which is surprising given that the word embedding model is
 932 much simpler. Further performance gains in RFS should be possi-
 933 ble by initializing user and item attribute embeddings to pretrained
 934 word embeddings as in Chen et al. [5].

935
 936 *Theory (approximation).* Proposition 2 means that RFS can approx-
 937 imate the output of any order-invariant model such as matrix fac-
 938 torization. However, there is no guarantee that fit with maximum
 939 likelihood estimation using the objective in Equation (7), RFS will
 940 approximate a specific model. We leave the development of such
 941 approximation techniques to future work.
 942

943
 944 *Theory (ranking metrics).* How well does binary classification per-
 945 form for other ranking-based recommendation metrics, such as
 946 non-discounted cumulative gain? Analyzing this question is more
 947 difficult, and we leave this to future work. We conjecture that a dif-
 948 ferent loss function should allow a similar proof to Proposition 1.
 949

950 REFERENCES

- 951
 952 [1] Trapit Bansal, David Belanger, and Andrew McCallum. 2016. Ask the GRU:
 953 Multi-task Learning for Deep Text Recommendations. *ACM Recommender Sys-
 954 tems*.
- 955 [2] Toine Bogers. 2018. *Tag-Based Recommendation*. Springer International Publish-
 956 ing, 441–479.
- 957 [3] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. En-
 958 riching Word Vectors with Subword Information. *Association for Computational
 959 Linguistics*.
- 960 [4] Da Cao, Liqiang Nie, Xiangnan He, Xiaochi Wei, Shunzhi Zhu, and Tat-Seng
 961 Chua. 2017. Embedding Factorization Models for Jointly Recommending Items
 962 and User Generated Lists. In *ACM SIGIR (SIGIR ’17)*. ACM, New York, NY, USA,
 963 10.
- 964 [5] Ting Chen, Liangjie Hong, Yue Shi, and Yizhou Sun. 2017. Joint Text Embedding
 965 for Personalized Content-based Recommendation. *arXiv:1706.01084*.
- 966 [6] Xin Dong, Lei Yu, Zhonghuo Wu, Yuxia Sun, Lingfeng Yuan, and Fangxi Zhang.
 967 2017. A Hybrid Collaborative Filtering Model with Deep Structure for Recom-
 968 mender Systems. *AAAI Conference on Artificial Intelligence*.
- 969 [7] Prem Gopalan, Laurent Charlin, and David M. Blei. 2014. Content-based Recom-
 970 mendations with Poisson Factorization. *Advances in Neural Information Process-
 971 ing Systems*.
- 972 [8] Alex Graves. 2013. Generating Sequences With Recurrent Neural Networks.
 973 *arXiv:1308.0850*.
- 974 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual
 975 Learning for Image Recognition. *IEEE Conference on Computer Vision and Pattern
 976 Recognition*.
- 977 [10] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng
 978 Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th Interna-
 979 tional Conference on World Wide Web*. International World Wide Web Confer-
 980 ences Steering Committee, 173–182.
- 981 [11] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neu-
 982 ral Computation* 9, 8, 1–32.
- 983 [12] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Opti-
 984 mization. *International Conference on Learning Representations*.
- 985 [13] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam R. Kosiorek, Seungjin Choi, and
 986 Yee Whye Teh. 2018. Set Transformer. *arXiv:1810.00825*.
- 987 [14] Jianxun Lian, Fuzheng Zhang, Xing Xie, and Guangzhong Sun. 2018. Towards
 988 Better Representation Learning for Personalized News Recommendation: a Multi-
 989 Channel Deep Fusion Approach. *International Joint Conference on Artifi-
 990 cial Intelligence*.
- 991 [15] Junjie Liang, Jinlong Hu, Shoubin Dong, and Vasant G. Honavar. 2018. Top-
 992 N-Rank: A Scalable List-wise Ranking Method for Recommender Systems.
 993 *arXiv:1812.04109*.
- 994 [16] Nan Liang, Hai-Tao Zheng, Jin-Yuan Chen, Arun Sangaiah, and Cong-Zhi Zhao.
 995 2018. TRSDL: Tag-Aware Recommender System Based on Deep Learning.
 996 Intelligent Computing Systems. *Applied Sciences*.
- 997 [17] Benedikt Loepp, Tim Donkers, Timm Kleemann, and Jürgen Ziegler. 2019. Inter-
 998 active recommending with Tag-Enhanced Matrix Factorization (TagMF). *Inter-
 999 national Journal of Human-Computer Studies*. Advances in Computer-Human
 999 Interaction for Recommender Systems.
- 999 [18] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-
 999 SNE. *Journal of Machine Learning Research* 9, 2579–2605.
- 999 [19] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed
 999 Representations of Words and Phrases and their Compositionality. *Neural Infor-
 999 mation Processing Systems*.
- 999 [20] Shumpei Okura, Yukihiro Tagami, Shingo Ono, and Akira Tajima. 2017. Embed-
 999 ding-based News Recommendation for Millions of Users. In *Proceedings
 999 of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and
 999 Data Mining (KDD ’17)*. ACM, New York, NY, USA, 1933–1942.
- 999 [21] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-
 999 Thieme. 2009. BPR: Bayesian Personalized Ranking from Implicit Feedback. *Uncer-
 999 tainty in Artificial Intelligence*.
- 999 [22] Yue Shi, Alexandros Karatzoglou, Linas Baltrunas, Martha Larson, Alan Hanjalic,
 999 and Nuria Oliver. 2012. TFMAP: Optimizing MAP for Top-N Context-Aware
 999 Recommendation. *ACM Special Interest Group on Information Retrieval*.
- 999 [23] Yue Shi, Martha Larson, and Alan Hanjalic. 2014. Collaborative Filtering Beyond
 999 the User-Item Matrix: A Survey of the State of the Art and Future Challenges.
 999 *Comput. Surveys* 47, 1, 1–45.
- 999 [24] Bo Song, Xin Yang, Yi Cao, and Congfu Xu. 2018. Neural Collaborative Ranking.
 999 *ACM Conference on Information and Knowledge Management*, 10.
- 999 [25] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. 2013. On
 999 the Importance of Initialization and Momentum in Deep Learning. *International
 999 Conference on Machine Learning*.
- 999 [26] Adrianus Willem van der Vaart. 1998. *Asymptotic Statistics*. Cambridge Univer-
 999 sity Press.
- 999 [27] Victor Veitch, Morgane Austern, Wenda Zhou, David M. Blei, and Peter Orbanz.
 999 2019. Empirical Risk Minimization and Stochastic Gradient Descent for Rela-
 999 tional Data. *International Conference on Artificial Intelligence and Statistics*.
- 999 [28] Chong Wang and David M. Blei. 2011. Collaborative Topic Modeling for Recom-
 999 mending Scientific Articles. *ACM Knowledge Discovery and Data Mining*.
- 999 [29] Zhenghua Xu, Thomas Lukasiewicz, Cheng Chen, Yishu Miao, and Xiangwu
 999 Meng. 2017. Tag-Aware Personalized Recommendation Using a Hybrid Deep
 999 Model. *International Joint Conference on Artificial Intelligence*, 3196–3202.
- 999 [30] Haochoo Ying, Liang Chen, Yuwen Xiong, and Jian Wu. 2016. Collaborative
 999 Deep Ranking: A Hybrid Pair-Wise Recommendation Algorithm with Implicit
 999 Feedback. *Pacific-Asia Conference on Knowledge Discovery and Data Mining*.
- 999 [31] Lu Yu, Chuxu Zhang, Shichao Pei, Guolei Sun, and Xiangliang Zhang. 2018. WalkRanker: A Unified Pairwise Ranking Model With Multiple Relations for
 999 Item Recommendation. *AAAI Conference on Artificial Intelligence*.
- 999 [32] Fajie Yuan, Guibing Guo, Joemon M. Jose, Long Chen, Haitao Yu, and Weinan
 999 Zhang. 2016. Optimizing Factorization Machines for Top-N Context-Aware Rec-
 999 ommendations. *Web Information Systems Engineering*.
- 999 [33] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabás Póczos, Ruslan
 999 Salakhutdinov, and Alexander J. Smola. 2017. Deep Sets. *Neural Information
 999 Processing Systems*.
- 999 [34] Shuai Zhang, Lina Yao, and Aixin Sun. 2017. Deep Learning Based Recommender
 999 System: A Survey and New Perspectives. *arXiv:1707.07435*.
- 999 [35] Yi Zhen, Wu-Jun Li, and Dit-Yan Yeung. 2009. TagiCoFi: Tag Informed Collabo-
 999 rative Filtering. *ACM Recommender Systems*, 8.
- 999 [36] Yi Zuo, Jiulin Zeng, Maoguo Gong, and Licheng Jiao. 2016. Tag-aware recom-
 999 mender systems based on deep neural networks. *Neurocomputing*. Big Learning
 999 in Social Media Analytics.

A REPRODUCIBILITY SUPPLEMENT: CODE

999 We give an example implementation of RANKFROMSETS with the
 999 inner product regression function in Equation (4) in python with
 999 the PyTorch package (version 1.0.0). This shows how easy it is to
 999 implement RFS; we used the same implementation to achieve state-
 999 of-the-art results in Section 4.2.

	# recommendations	10	30	50	70	100
	Recall (percent)	0.02	0.06	0.11	0.14	0.21

Table 2: The permutation-marginalized LSTM model performs poorly on the arXiv dataset of researcher reading behavior in terms of in-matrix recall. We report the held-out recall averaged over 100 users due to computational constraints of the model. It is much lower than that of other methods, but better than random.

B REPRODUCIBILITY SUPPLEMENT: LSTM

C REPRODUCIBILITY SUPPLEMENT: PERMUTATION-INVARIANT MODELS

Matrix factorization permutation-invariant. We show that collaborative topic Poisson factorization [7] is permutation-invariant. The likelihood of the attributes (words) associated with items (documents) in this model factorizes, and multiplication is permutation-invariant. Conditional on the latent item representation θ_d and latent word representation β_v , every word in the document w_{dv} is independent. The joint probability of the words in a document factorizes:

$$p(w_d | \theta_d, \beta_v) = \prod_{w_{dv} \in w_d} p(w_{dv} | \theta_d, \beta_v). \quad (13)$$

In this model, predictions are made using expectations under the posterior. The posterior is proportional to the log joint of the model, and the attributes of items (words in documents) enter into the model only via the above product. The product of the probability of words in a document is invariant to a reordering of the words in the document. Therefore, collaborative topic Poisson factorization is permutation-invariant.

D REPRODUCIBILITY SUPPLEMENT: SIMULATION STUDY

With finite data and a finite number of parameters, the optimal parameterization of RFS is dependent on the data-generating distribution. Recall that observations of user-item interactions are generated by a Bernoulli distribution with logit function f . We describe a choice of logit function f that leads to the residual and deep architectures in Equations (5) and (6) outperforming the inner product architecture in Equation (4) in terms of predictive performance. Note that the number of parameters across architectures must be equal for a fair comparison.

We simulate data from the following generative process:

1. **For every user u :**

- (a) Draw user embedding $\theta_u \sim \text{Normal}(0, I)$.

2. **For every item attribute j :**

- (a) Draw attribute embedding $\beta_j \sim \text{Normal}(0, I)$.

3. **For every item m :**

- (a) Draw item topics $\theta_m \sim \text{Dirichlet}(\alpha)$

- (b) Draw number of item attributes $M \sim \text{Poisson}(\lambda)$

- (c) Draw nonzero item attributes $x_m^+ \sim \text{Multinomial}(M, \theta_m)$.

4. **For every user, item observation:**

- (a) $y_{um} \sim \text{Bernoulli}\left(y_{um}; \sigma(f(\theta_u, x_m^+))\right)$

The logit function f is the square kernel:

$$f(\theta_u, x_m^+) = \left(\theta_u^\top \sum_{j \in x_m^+} \beta_j \right)^2.$$

Before sampling from the Bernoulli, we standardize the logits output by f across users and subtract 7 to achieve sparse user-item observations.

Simulation setup. We set the Dirichlet parameter to be $\alpha = 0.01$ and the Poisson rate to be $\lambda = 20$. We generate data for 1k users, 5k item attributes, 30k items, and hold out 100 users for each of the validation and test sets. We fix the momentum to 0.9 [25] and grid search over stochastic gradient descent learning rates of 10, 1, 0.1, 0.01 and over two learning rate decay schedules. The first linear learning rate decay goes to zero over 100k iterations, while the second divides the learning rate by 10 if the validation in-matrix recall does not improve (evaluation is performed every 500 iterations). We run the grid search on one instance of data generated from this model, then for the best performing hyperparameters for each model trained on this instance, we regenerate data 30 times and average results over these synthetic datasets.

Regression function	Inner product	Deep	Residual
Recall	0.29 ± 0.15	0.32 ± 0.14	0.33 ± 0.18

Table 3: A simulation study demonstrating that the choice of parameterization of RANKFROMSETS is data-dependent. We report the in-matrix recall averaged over 100 users and 30 replications of the simulation where data is regenerated. The residual model in Equation (6) outperforms the deep model, Equation (5), and the inner product model in Equation (4).

Simulation results. The results in Table 3 demonstrate that the residual model outperforms both the deep and inner product architectures for data generated by the above generative process.

1161 *Generalization.* The above example shows that the choice of archi-
1162 ture in RFS is data-dependent. To ensure that the model does
1163 not overfit as new users or items are included in the training data,
1164 we need to compare the number of parameters to the number of
1165 datapoints. A model with parameters the size of the training data
1166 can overfit by memorizing the training data. For generalization to
1167 be possible, overfitting can be avoided if the number of parame-
1168 ters grows slower than the size of the data. The technical backing
1169 for this comes from asymptotic statistics and the concept of sieved
1170 likelihoods. Specifically, the maximum likelihood estimation pro-
1171 cedure with the objective function in Equation (7) can be replaced
1172 by maximization of a sieved likelihood function. The ‘sieve’ refers
1173 to filtering information as the number of parameters (in this case,
1174 user and item representations) grows with the number of obser-
1175 vations. The sieved likelihood function enables the analysis of as-
1176 ymptotic behavior as the number of users grows $U \rightarrow \infty$ and the
1177 number of items grows $I \rightarrow \infty$. An example of a technique to grow
1178 the number of parameters in a way that supports generalization is
1179 given in Chapter 25 of Vaart [26].

1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218

1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276