

# Self-supervised Learning for Large-scale Item Recommendations

Tiansheng Yao, Xinyang Yi, Derek Zhiyuan Cheng, Felix Yu, Ting Chen, Aditya Menon

Lichan Hong, Ed H. Chi, Steve Tjoa, Jieqi (Jay) Kang, Evan Ettinger

{tyao,xinyang,zcheng,felixyu,iamtingchen,adityakmenon}@google.com

{lichan,edchi,stevetjoa,jaykang,eettinger}@google.com

Google Inc.

## ABSTRACT

Large scale recommender models find most relevant items from huge catalogs, and they play a critical role in modern search and recommendation systems. To model the input space with large-vocab categorical features, a typical recommender model learns a joint embedding space through neural networks for both queries and items from user feedback data. However, with millions to billions of items, the power-law user feedback makes labels very sparse for a large amount of long-tail items.

Inspired by the recent success in self-supervised representation learning research in both computer vision and natural language understanding, we propose a multi-task self-supervised learning (SSL) framework for large-scale item recommendations. The framework is designed to tackle the label sparsity problem by learning more robust item representations. Furthermore, we propose two self-supervised tasks applicable to models with categorical features within the proposed framework: (i) Feature Masking (FM) and (ii) Feature Dropout (FD).

We evaluate our framework using two large-scale datasets with 500M and 1B training examples respectively. Our results demonstrate that the proposed framework outperforms traditional supervised learning only models and state-of-the-art regularization techniques in the context of item recommendations. The SSL framework shows larger improvement with less supervision compared to the counterparts. We also apply the proposed techniques to a web-scale commercial app-to-app recommendation system, and significantly improve top-tier business metrics via A/B experiments on live traffic. Our online results also verify our hypothesis that our framework indeed improves model performance on slices that lack supervision.

## 1 INTRODUCTION

In the last few years, neural network based models have emerged to the main stage of modern recommendation systems throughout the industry (see, e.g., [1, 8, 11, 20, 35, 44, 47]), and academia ([9, 36]). Compared to conventional approaches like matrix factorization [2, 22, 25, 26], gradient boosted decision trees [5, 23, 33], and logistic regression based recommenders [21], these deep models handle categorical features more effectively, enable more complex data representations, and introduce more non-linearity to better fit the complex data for recommenders.

A particular recommendation task we focus on in this paper is to identify the most relevant items given a query from a huge item catalog. This general problem of *large-scale item recommendations* has been widely adopted in various applications. Depending on the type of the query, a recommendation task could be: (i) personalized recommendation, when the query is a user; (ii) item to item

recommendation, when the query is also an item; and (iii) search, when the query is a piece of free text. To model the interactions between a query and an item, a well-known approach leverages *embedding-based neural networks*. The recommendation task is typically formulated as an extreme classification problem [11] where each item is represented as a dense vector in the output space. More recently, a line of research [27, 44] adopts two-tower DNNs (Figure 1) where a neural network encodes a set of item features to an embedding. These factorized model architectures make it possible to convert the problem of finding top-k items to Maximum-Inner-Product-Search (MIPS) [10], solvable in sublinear complexity. This also makes it suitable for serving recommendations in real time.

Embedding-based deep models typically have large amount of parameters because they are built with high-dimensional embeddings that represent high cardinality sparse features such as topics or item IDs. In many existing literature, the loss functions for training these models are formulated as a supervised learning problem. The supervision comes from the collected label (e.g., click) between a pair of query and item. Accordingly, similar to other applications, training deep models in recommendations requires large amount of training data. One might argue that collecting a large amount of data is not an issue given the ubiquitous internet services and increasing amount of user data. However, when it comes to modeling a huge catalogue of items in the order of millions (e.g., songs and apps [32]) to even billions (e.g., videos on YouTube [11]), there are still lots of challenges, including:

- **Highly-skewed data distribution:** The interaction between queries and items are often highly skewed in a power-law distribution [34]. So a small percentage of the items gets most of the interactions. This will always leave the training data for long-tail items very sparse.
- **Lack of explicit user feedback:** Although modern recommendation systems in the industry collect hundreds of billions of implicit user actions like clicks and thumb-ups, these systems still lack high quality user explicit feedback like item ratings, feedback for user happiness, and relevance scores.

Self-supervised learning (SSL) offers a different angle to learn deep representation via unlabeled data. The basic idea is to enhance training data with various data augmentations, and supervised tasks to predict or reconstitute the original examples as auxiliary tasks. Self-supervised learning has been widely used in the areas of Computer Vision (CV) [17, 29, 37] and Natural Language Understanding (NLU) [13, 28]. An example work [29] in CV proposed to rotate images at random, and train a model to predict how each augmented input image was rotated. In NLU, masked language task was introduced in the BERT model, to help improve pre-training of language models. Similarly, other pre-training tasks like predicting

surrounding sentences and linked sentences in Wikipedia articles have also been used in improving dual-encoder type models in NLU [4]. Compared to conventional supervised learning, self-supervised learning provides complementary objectives eliminating the prerequisite of collecting labels manually. In addition, SSL enables autonomous discovery of good semantic representations by exploiting the internal relationship of input features.

Despite the wide adoption in computer vision and natural language understanding, the application of self-supervised learning in the field of recommendation systems is less well studied. The closest line of research studies a set of regularization techniques [19, 27, 46], which are designed to force learned representations (i.e., output layer (embeddings) of a multi-layer perceptron), of different examples to be farther away from each other, and spread out in the entire latent embedding space. Although sharing similar spirit with SSL, these techniques do not explicitly construct SSL tasks. In contrast to models in CV or NLU applications, recommendation model takes extremely sparse input where high cardinality categorical features are one-hot (or multi-hot) encoded, such as the item IDs or item categories [35]. These features are typically represented as learnable embedding vectors in deep models. As most models in computer vision and NLU deal with dense input, existing methods for creating SSL tasks are not directly applicable to the sparse models in recommendation systems. More recently, a line of research studies self-supervised learning improving sequential user modeling in recommendations [31, 42, 48]. Different from these works, this paper focuses on item representation learning, and shows how SSL can help improve generalization in the context of long-tail item distribution. Moreover, in contrast to using SSL on a certain sequential user feature, we design new SSL tasks and demonstrate their effectiveness for learning with a set of heterogeneous categorical features, which we believe is a more general setup for other types of recommendation models such as multitask ranking models (e.g., [47]).

In this paper, we propose to leverage self-supervised learning based auxiliary tasks to improve item representations, especially with long-tail distributions and sparse data. Different from CV or NLU applications, input space of recommendation model is highly *sparse* and represented by a set of categorical features (e.g. item ids) with large cardinality. For such sparse models, we propose a new SSL framework, where the key idea is to: (i) augment data by masking input information; (ii) encode each pair of augmented examples by a two-tower DNN; and (iii) apply a contrastive loss to learn representations of augmented data. The goal of contrastive learning is to let augmented data from the same example be discriminated against others. Note that the two-tower DNN for contrastive learning and the one for encoding query and item can share a certain amount of model parameters. We will discuss the details in Section 3.

Our contribution is four-fold:

- **SSL framework:** We present a model architecture agnostic self-supervised learning framework for sparse neural models in recommendations. The auxiliary self-supervised loss and the primary supervised loss are jointly optimized via a multitask learning framework. We focus on using this framework for efficiently scoring a large corpus of items, which is also

known as item retrieval in two-stage recommenders [11]. We believe it would also shed light on designing SSL for other types of models such as ranking models [8].

- **SSL tasks:** We propose two generalizable approaches to construct auxiliary tasks, (i) Feature Masking and (ii) Feature Dropout, within the proposed framework. The two tasks are tailored for heterogeneous categorical features that are common in recommender models.
- **Offline experiments:** On one public dataset and one industry scale dataset for recommendation systems, we demonstrate that introducing SSL as an auxiliary task can significantly improve model performance, especially when labels are scarce. Comparing to the state-of-art non-SSL regularization techniques [19, 27, 46], we demonstrate that SSL consistently performs better, and improves model performance even when non-SSL regularization does not bring any additional gains.
- **Live experiment in a web-scale recommender:** We apply the proposed SSL technique in a fairly strong two-tower app-to-app recommendation model used in a large-scale real-world system. Live A/B testing shows significantly improved top-tier metrics for adopting our SSL framework. We especially see bigger improvements for slices without much supervision.

## 2 RELATED WORK

*Self-supervised Learning and Pre-training.* Various unsupervised / self-supervised learning tasks have been studied in the computer vision community [24]. Popular SSL tasks include: (i) predicting image rotations [17]; (ii) predicting relative patch locations [37]; (iii) predicting next video frames [39]; and (iv) leveraging contrastive loss [14] etc. In natural language understanding, SSL like pre-training tasks such as next sentence/word prediction and masked-LM have been widely used [13].

In NLU, for dual-encoder models, Chang et al. shows that pre-training tasks better aligned with the final task are more helpful than generic tasks such as next sentence prediction and masked-LM. The pre-training tasks are designed to leverage large-scale public NLU content, such as Wikipedia. In this paper, we also use the dual-encoder model architecture. Different from the above, the proposed self-supervision tasks do not require the use of a separate data source.

In computer vision, the closest line of research is SimCLR [6] which also utilizes self-supervised learning and contrastive learning for visual representation learning. Different with SimCLR and other works [3, 7] in vision, here we propose augmentations that are more suitable and tailed for categorical features for recommendations, instead of relying on image-specific augmentations such as image cropping, rotation and color distortion. In addition, the proposed framework does not require multi-stage training schedules (such as pre-training then fine-tuning).

*Spread-out and Instance-based Discriminative Learning.* Zhang et al. [46] and Wu et al. [41] use spread-out regularization for improving generalization of deep models. Specifically, in [46], a regularization promoting separation between random instances is shown to improve training stability and generalization. In [41], one

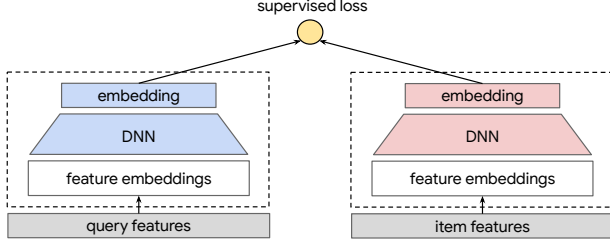


Figure 1: Model architecture: Two-tower DNN with query and item representations.

objective is to train a classifier treating each instance as its own class, therefore promoting large instance separations in the embedding space. Both the above approaches are studied for computer vision applications.

*Neural Modeling for Large-scale Item Recommendations.* Deep learning has led to many successes in building industry-scale recommender systems such as video suggestion [11], news recommendation [38], and visual discovery in social networks [30, 45]. To handle the challenge of scoring a large number of items, many recommender systems follow a two-stage design that first relies on a *candidate generation* (retrieval) model to find thousands of relevant items given a query, and then uses a finer-granularity *ranking model* to identify the best items to show to the user [8, 47]. On the retrieval side, factorized architecture with query and item embeddings are widely used due to its scalability for scoring a large catalog of items. The item scores are typically computed by dot-product between query and item embeddings so that finding top-k items can be converted to MIPS problem [10] with sublinear time complexity. One popular factorized structure is softmax-based multi-class classification model. The work in [11] treats the retrieval task as an extreme multi-class classification trained with multi-layer perceptron (MLP) model using sampled softmax as loss function. Such models leverages item ID as the only item feature, and thus might suffer from cold-start issue. More recently, a line of research [27, 44] considers applying two-tower DNNs on retrieval problems, which is also known as dual-encoder [18, 43], where item embeddings are constructed by a MLP from ID and other categorical metadata features. The self-supervised approach proposed is applicable to both ranking and retrieval models. In this paper we focus on using SSL for retrieval models, particularly, on improving item representations in two-tower DNNs.

*Self-supervised Learning in Sequential Recommendations.* In recommender systems, a line of research has been recently studied for utilizing self-supervised learning for sequential recommendation. Self-supervised learning tasks are designed to capture information among user history [48] and learn more robust disentangled user representation [31] in user sequential recommendation. Moreover, Xin et al. shows combining SSL with reinforcement learning is effective to capture long-term user interest in sequential recommendation. Different from the above, our proposed SSL framework is focusing on improving item representation with long-tail distributions. The proposed SSL tasks do not require modeling sequential

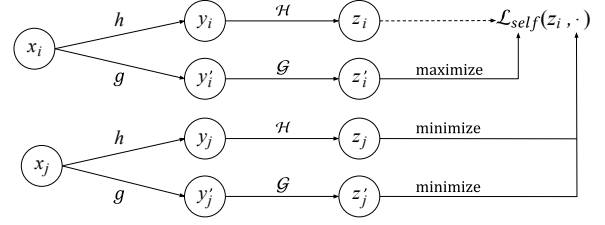


Figure 2: Self-supervised learning framework illustration. Two data augmentations  $h$  and  $g$  are applied to the original input; Encoders  $\mathcal{H}$  and  $\mathcal{G}$  are applied to the augmented examples  $y_i$  and  $y'_i$  to generate embeddings  $z_i$  and  $z'_i$ . The SSL loss  $\mathcal{L}_{self}$  w.r.t.  $z_i$  is optimized towards maximizing the similarity with  $z'_i$  while minimizing the similarity between  $z_j$  and  $z'_j$ .

information and are generally applicable to deep models with heterogeneous categorical features.

### 3 SELF-SUPERVISED LEARNING FRAMEWORK FOR LARGE-SCALE ITEM RECOMMENDATIONS

In this section, we present our framework of self-supervised learning for deep neural net models for recommenders using large-vocab categorical features. Particularly, a general self-supervised learning framework is introduced in Section 3.1. In Section 3.2, we present two methods to construct SSL tasks and elaborates on their connections to spread-out regularization. Finally, in Section 3.3, we describe how to use SSL to improve factorized models (i.e., two-tower DNNs as shown in Figure 1), via a multi-task learning framework. It is worth noting that the SSL framework we propose is model architecture agnostic. In this paper, we focus on two-tower models, however the same learning should apply to other model architectures like matrix factorization and multilayer perceptron.

#### 3.1 Framework

Inspired by the SimCLR framework [6] for visual representation learning, we adopt similar contrastive learning algorithms for learning representations of categorical features. The basic idea is two folds: first, we apply different data augmentation for the same training example to learn representations; and then use contrastive loss function to encourage the representations learned for the same training example to be similar. Contrastive loss was also applied in training two-tower DNNs (see e.g., [27, 44]), although the goal there was to make positive item agree with its corresponding queries.

We consider a batch of  $N$  training examples  $x_1, \dots, x_N$ , where  $x_i \in \mathcal{X}$  represents a set of features for example  $i$ . In the context of recommenders, an example indicates a query, an item or a query-item pair. Suppose there are a pair of transform function  $h, g : \mathcal{X} \rightarrow \mathcal{X}$  that augment  $x_i$  to be  $y_i, y'_i$  respectively,

$$y_i \leftarrow h(x_i), y'_i \leftarrow g(x_i). \quad (1)$$

Given the same input of example  $i$ , we want to learn different representations  $y_i, y'_i$  after augmentation to make sure the model

still recognizes that both  $y_i$  and  $y_i$  represent the same input  $i$ . In other words, the contrastive loss learns to minimize the difference between  $y_i, y_i'$ . In the mean time, for different example  $i$  and  $j$ , the contrastive loss maximizes the difference between the representations learned  $y_i, y_j'$  after data different augmentations. Let  $z_i, z_i'$  denote the embeddings of  $y_i, y_i'$  after encoded by two neural networks  $\mathcal{H}, \mathcal{G} : \mathcal{X} \rightarrow \mathbb{R}^d$ , that is

$$z_i \leftarrow \mathcal{H}(y_i), z_i' \leftarrow \mathcal{G}(y_i'). \quad (2)$$

We treat  $(z_i, z_i')$  as positive pairs, and  $(z_i, z_j')$  as negative pairs for  $i \neq j$ . Let  $s(z_i, z_j') = \langle z_i, z_j' \rangle / (\|z_i\| \cdot \|z_j'\|)$ . To encourage the above properties, we define a loss function for each example in a way similar to the widely used softmax cross entropy:

$$\mathcal{L}_{self}(\mathcal{H}, \mathcal{G}) := -\frac{1}{N} \sum_{i \in [N]} \log \frac{\exp(s(z_i, z_i')/\tau)}{\sum_{j \in [N]} \exp(s(z_i, z_j')/\tau)}. \quad (3)$$

where  $\tau$  is a tunable hyper-parameter for the softmax temperature. The above loss function learns a robust embedding space such that similar items are close to each other after data augmentation, and random examples are pushed farther away. The overall framework is illustrated in Figure 2.

*Encoder Architecture.* For input examples with categorical features,  $\mathcal{H}, \mathcal{G}$  are typically constructed with an input layer and a multi-layer perceptron (MLP) built on top of it. The input layer is often a concatenation of normalized dense features and multiple sparse feature embeddings, where the sparse feature embeddings are learnt representations stored in embedding tables (In contrast, the input layers for computer vision and language models directly work on raw inputs). In order to make SSL facilitate the supervised learning task, we share the embedding table of sparse features for both neural networks  $\mathcal{H}, \mathcal{G}$ . Depending on the technique for data augmentation  $(h, g)$ , the MLPs of  $\mathcal{H}$  and  $\mathcal{G}$  could also be fully or partially shared.

*Connection with Spread-out Regularization.* In the special case where  $(h, g)$  are identical map and  $\mathcal{H}, \mathcal{G}$  are the same neural network, loss function in equation (3) is then reduced to

$$-N^{-1} \sum_i \log \frac{\exp(1/\tau)}{\exp(1/\tau) + \sum_{j \neq i} \exp(s(z_i, z_j)/\tau)}$$

which encourages learned representations of different examples to have small cosine similarity. The loss is similar to the spread-out regularization in [46], except that the original proposal uses square loss, i.e.,  $N^{-1} \sum_i \sum_{j \neq i} \langle z_i, z_j \rangle^2$ , instead of softmax. Spread-out regularization has been proven to improve generalization of large-scale retrieval models. In Section 4, we show that by introducing specific data augmentations, using SSL-based regularization can further improve model performance compared to spread-out regularization.

### 3.2 Self-supervised Learning Tasks

The most critical question to answer is how to design the transformation functions  $h$  and  $g$ . A good transformation and data augmentation should make minimal amount of assumptions on the data such that it can be applicable to a large variety of tasks. Compared to previous works for SSL in visual representation learning that rely

on image cropping and color distortion for augmentation, here we propose augmentation tasks that are tailored to categorical features that are widely applicable in recommendation systems.

Inspired by masked-LM used in pre-training for NLU models (e.g. BERT [13]), we present two approaches to construct the auxiliary SSL task. Note that in the rest of this paper, we focus on applying the self-supervised learning tasks on the item side, although the tasks can also be applied on the queries. These tasks can be easily applied to all scenarios where there are multiple features per item.

- **Feature Masking Task:**  $(h, g)$  are designed to mask a subset of input features, so that they would only use partial inputs to learn representations. For example, for movie recommendations with movie features  $\{\text{movie\_id}, \text{genre}, \text{category}\}$ , one set of feature could be  $\{\text{genre}, \text{category}\}$ , and another set could be  $\{\text{movie\_id}\}$ . By masking through feature splitting, this approach aims to learn the internal relationship between two sets of features. The learned feature embeddings or intermediate representations in neural network can then be shared with down-stream tasks for retrieval and ranking. Note that it's possible to split the whole feature set into  $k$  groups, and optimize  $k(k-1)/2$  SSL loss functions like Equation (3). With feature masking, as input layer sizes of network  $\mathcal{H}, \mathcal{G}$  are no longer same, we use different MLP parameters, but share the input feature embeddings between  $\mathcal{H}$  and  $\mathcal{G}$ .
- **Feature Dropout Task:** As demonstrated in [40], randomly dropping out features during training time can improve cold-start recommendation by learning not to overly rely on certain input features. Following this idea, we propose to apply random feature dropout as feature augmentation, i.e.,  $(h, g)$  randomly retain each categorical feature with certain probability. This is essentially stochastic feature selection. By applying two independent dropout processes, the SSL task aims to use a random set of input features to predict another random set. This method can avoid manual design with prior knowledge used in the first approach proposed above. When a categorical feature is dropped out, it is treated as missing value, and represented by a default out-of-vocab embedding. In this case, we use the same network for  $\mathcal{H}, \mathcal{G}$ . This network can be also shared with the main task in multi-task training as introduced below.

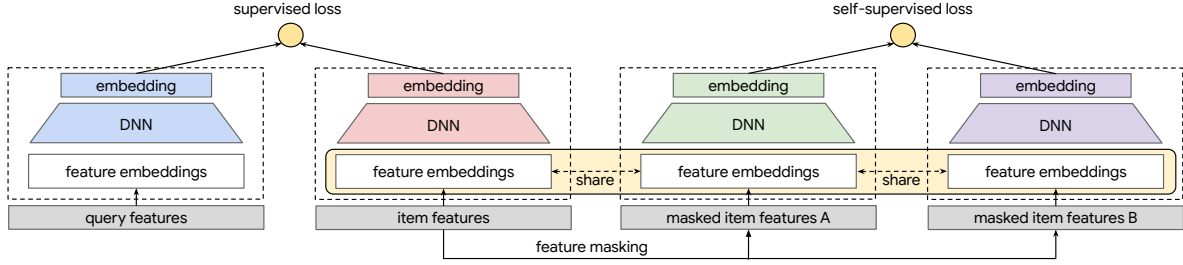
### 3.3 Multi-task Training

To enable SSL learned representations to help improve the learning for the main supervised task such as regression or classification, we leverage a multi-task training strategy where the main supervised task and the auxiliary SSL task are jointly optimized. Intuitively, a linear weighted sum approach is applied:

$$\mathcal{L} = \mathcal{L}_{main} + \alpha \cdot \mathcal{L}_{self}, \quad (4)$$

where  $\mathcal{L}_{main}$  is the loss function from main task, and  $\alpha$  controls to what extent SSL regularization is applied. An alternative training strategy could be a two-phase training plan where the model firstly pre-trains on  $\mathcal{L}_{self}$ , and then fine tunes by the main task  $\mathcal{L}_{main}$ . We leave the experimentation of this method to future work.





**Figure 3: Model architecture: Two-tower Model with feature masking (FM) SSL task.** In the FM task, we mask the item features by two separate sets  $\Theta_A$  and  $\Theta_B$ . The masked item features are used to learn the item embeddings in the FM task. The feature embeddings are shared with the supervised task, but DNN modules in different colors do not share parameters.

As mentioned in Section 2, we focus on improving recommender models built with two-tower DNNs. In practice, we found two-tower models quite effective for recommendation / retrieval tasks. Compared to matrix factorization models, it works very well to (1) learn large-scale sparse feature representations; and thus (2) improves model generalization to mitigate cold-start issues. Compared to a multilayer perceptron (MLP) model, model inference is extremely efficient for two-tower models due to the simple computation of dot product and highly efficient systems for nearest neighbor search.

We consider a batch of positive query-item pairs  $\{(q_i, c_i)\}_{i=1}^N$ , and the goal is to retrieve the top K most relevant items from a large corpus. We use two neural networks to map query and item to embedding vectors  $\{(q_i, c_i)\}_{i=1}^N$  respectively. Similar to the contrastive loss in (3), the loss for main task is

$$\mathcal{L}_{main} := -\frac{1}{N} \sum_{i \in [N]} \log \frac{\exp(s(q_i, c_i)/\tau)}{\sum_{j \in [N]} \exp(s(q_i, c_j)/\tau)}. \quad (5)$$

It is worth noting that the SSL task can be applied on either query or item side. Figure 3 and Figure 4 demonstrates the model architecture for both feature masking and feature dropout tasks to facilitate item representation learning.

## 4 OFFLINE EXPERIMENTS

In this section, we provide empirical results to demonstrate the effectiveness of our proposed self-supervised framework both in academic open dataset and in actual large-scale recommendation products. The experiments are designed to answer the following research questions.

- **RQ1:** We have seen the successes of SSL in both computer vision and natural language understanding. Does the proposed SSL Framework improve deep models in recommendations?
- **RQ2:** SSL is designed to improve primary supervised task by learning representations from unlabeled examples. What is the impact of the amount of training data on the improvement from SSL?
- **RQ3:** Under the multi-task learning framework, how does the self-supervised loss interact with the supervised loss? Specifically, how does weight  $\alpha$  on the SSL loss in Equation (3) affect model quality?

- **RQ4:** Which SSL approach works better? How does SSL compare with non-SSL regularization such as spread-out and dropout?
- **RQ5:** For feature dropout, how does dropout rate affect the model quality?

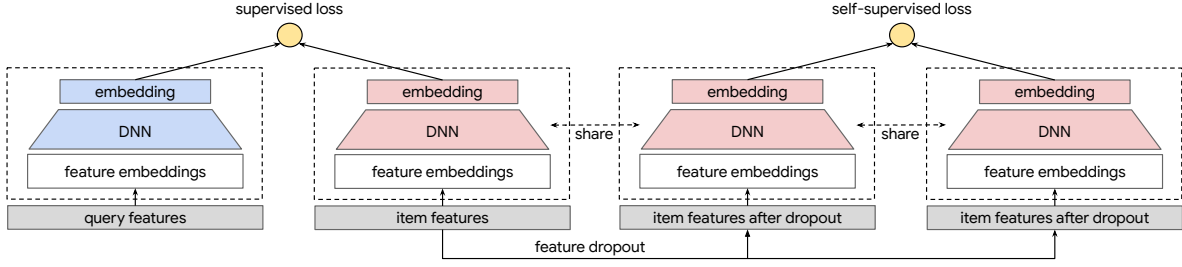
The above questions are addressed in order from Section 4.3 - 4.8.

### 4.1 Datasets

We conduct experiments on two large-scale datasets that both come with rich set of item metadata features. We formulate their primary supervised task as an item-to-item recommendation problem to study the effects of SSL on training recommender (in this case, retrieval) models. Table 1 shows some basic stats for the two datasets. We use 500M and 1B training examples from the two datasets respectively to train and evaluate our models.

**Wikipedia [16].** The first dataset focuses on the problem of link prediction between Wikipedia pages. It consists of pairs of pages  $(x, y) \in \chi \times \chi$ , where  $x$  indicates a *source page*, and  $y$  is a *destination page* linked from  $x$ . The goal is to predict the set of pages that are likely to be linked to a given source page from the whole corpus of web pages. Each page is represented by a feature vector  $x = (x_{id}, x_{ngrams}, x_{cats})$ , where all the features are categorical. Here,  $x_{id}$  denotes the one-hot encoding of the page URL,  $x_{ngrams}$  denotes a bag-of-words representation of the set of n-grams of the page's title, and  $x_{cats}$  denotes a bag-of-words representation of the categories that the page belongs to. We partitioned the dataset into training and evaluation using a (90%, 10%) split, following the same treatment in [27] and [44].

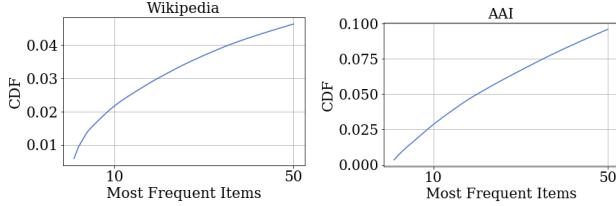
**App-to-App Install (AAI).** The AAI dataset was collected on the app landing pages from a commercial mobile app store. On a particular app's (seed app) landing page, the app installs (candidate apps) from the section of recommended apps were collected. Each training example represents a pair of seed-candidate pairs denoted as  $(x_{seed}, x_{candidate})$  and their metadata features. The goal is to recommend highly similar apps given a seed app. This is also formulated as an item-to-item recommendation problem via a multi-class classification loss. Note that we only collect positive examples, i.e.,  $x_{candidate}$  is an installed app from the landing page of  $x_{seed}$ . All the impressed recommended apps with no installs are all ignored since we consider them more like weak positives instead of negatives



**Figure 4: Model architecture: Two-tower model with feature dropout (FD) SSL task. In the FD task, we randomly dropout item feature values to learn item embeddings. The whole item tower (in red) is shared with the supervised task.**

Dataset	# queries	# items	# examples
Wikipedia	5.3M	5.3M	490M
AAI	2.4M	2.4M	1B

**Table 1: Corpus sizes of the Wikipedia and the AAI datasets.**



**Figure 5: CDF of most frequent items in the Wikipedia and AAI datasets. The distribution is extremely dominated by popular items. For example, the top 50 items out of the 2.4M items already constitute 10% of data in the AAI dataset.**

for building retrieval models. Each item (app) is represented by a feature vector  $\mathbf{x}$  with the following features:

- *id*: Application id as a one-hot categorical feature.
- *developer\_name*: Name of the app developer as a one-hot categorical feature.
- *categories*: Semantic categories of the app as a multi-hot categorical feature.
- *title\_unigram*: Uni-grams of the app title as a multi-hot categorical feature.
- *popularity*: App popularity as a dense feature.
- *quality*: App quality as a dense feature.

Figure 5 shows the CDF of most frequent items for the two datasets, indicating a highly skewed data distribution. For example, the top 50 items in the AAI dataset collectively appeared roughly 10% in the training data. If we consider a naive baseline (i.e., *TopPopular* recommender [12]) that recommends the most frequent top-K items for every query, the CDF of the  $K$ -th frequent item essentially represents the *Recall@K* metric of such baseline. This suggests a naive *TopPopular* recommender achieves *Recall@50*  $\approx 0.1$  for AAI and *Recall@50*  $\approx 0.05$  for Wikipedia. We present that all the proposed methods outperform this baseline by a large margin in Section 4.

## 4.2 Experiment Setup

**Backbone Network.** For the main task that predicts relevant items given the query, we use the two-tower DNN to encode query and items features (see Figure 1) as the *backbone network*. The item-to-item recommendation problem is formalized as a multi-class classification problem, using the batch softmax loss presented in Equation (5) as the loss function. Note that the proposed framework is architecture agnostic and applicable to networks with sparse categorical features. We choose the two-tower DNN architecture due to its effectiveness in industrial scale recommendation systems [11, 44].

**Hyper-parameters.** For the backbone two-tower DNN, we search the set of hyper-parameters such as the learning rate, softmax temperature ( $\tau$ ) and model architecture that gives the highest *Recall@50* on the validation set. Note that the training batch size in batch softmax is critical for model quality as it determines the number of negatives used for each positive item. Throughout this section, we use batch sizes 1024 and 4096 for Wikipedia and AAI respectively. We find that by further increasing batch size, the quality improvement becomes marginal. We then fine-tuned the number of hidden layers, hidden layer sizes and softmax temperature  $\tau$  for the baseline models. Specifically, for Wikipedia dataset, we end up with a softmax temperature  $\tau = 0.07$ , and *hidden\_layers* with sizes [1024, 128]. For AAI, we use  $\tau = 0.06$  and *hidden\_layers* [1024, 256]. Note that the dimension of last hidden layer is also the dimension of final query and item embeddings.

**SSL Networks and Tasks.** We fix the backbone network architecture in training with SSL tasks. For each of the SSL tasks, we have:

- (1) *FeatureMasking*: The SSL item encoders have the same final embedding dimensions with the primary task as plotted in Figure 3. For Wikipedia, we have *hidden\_layers* = [128]; for AAI, *hidden\_size* = [256]. We use a fixed set of features as mask  $\Theta_A$  and the rest features as another set of mask  $\Theta_B = \Theta'_A$ . For Wikipedia, we have  $\Theta_A = \{url\}$ ; for AAI, we have  $\Theta_A = \{id, developer\_name\}$ .
- (2) *FeatureDropout*: The SSL item encoders are identical to the item tower (red boxes in Figure 4). We randomly drop out values at the same dropout rate for all features.

Full Dataset				
Method	Recall@10	Recall@50	MAP@10	MAP@50
Baseline	0.0537	0.1930	0.0171	0.0229
DO [40]	0.0535	0.1912	0.0172	0.0229
SO [46]	0.0549	0.1956	0.0176	0.0235
FD	<b>0.0605</b>	0.2056	<b>0.0198</b>	<b>0.0259</b>
FM	0.0539	0.1934	0.0172	0.0230
FD+FM	0.0596	<b>0.2061</b>	0.0193	0.0255

Table 2: Experiment results on the Wikipedia dataset.

*Training and Evaluation.* For model training, we use Adagrad [15] optimizer with learning rate 0.01. To evaluate the recommendation performance given a seed item, we compute and find the top  $K$  items with the highest cosine similarity from the whole corpus and evaluate the quality based on the  $K$  retrieved items. Note this is a relatively challenging task, given the sparsity of the dataset and large number of items in the corpus. We adopt popular standard metrics  $Recall@K$  and mean average precision ( $MAP@K$ ) to evaluate recommendation performance [20]. For each configuration of experiment results, we ran the experiment 5 times and report the average.

### 4.3 Effectiveness of SSL

To answer **RQ1**, we first evaluate the impact on model quality for a combination of SSL tasks within the proposed framework. For the baseline method, besides using the task without SSL, we also consider training the main task trained with 2 popular regularization techniques: dropout (DO), and spread-out regularization (SO). Note that *DO* is a popular regularization technique in industrial-size recommendation systems to enable more more generalization. We include it as a baseline for the Feature Dropout (FD) task to isolate the potential improvement from using dropout only. These SSL techniques are also complementary to each other, and can also be added on top of the baseline methods. In our experiments, we also evaluate the performance of combining the FD and FM tasks. It's worth noting that, although SSL techniques are applicable to improve both query and item representation learning, in our experiments, we only add SSL tasks on the item towers. Specifically, we compare:

- *Baseline*: Backbone two-tower DNN model.
- *DO* [40]: Backbone model with random feature dropout training on the item tower (no SSL task) for regularization.
- *SO* [46]: Backbone model with spread-out regularization on the item tower.
- *FM*: Backbone model with *Feature Masking* (FM) as the SSL task.
- *FD*: Backbone model with *Feature Dropout* (FD) as the SSL task.
- *FD+FM*: Backbone model with FD and FM used together as two SSL tasks.

We conduct a large scale study to evaluate the effectiveness of the proposed framework. For each of the above methods, we experiment with different choices of SSL weights  $\alpha \in \{0.1, 0.3, 1.0, 3.0\}$ , and report the best results. In Section 4.6, we show the effect of different

Full Dataset				
Method	Recall@10	Recall@50	MAP@10	MAP@50
Baseline	0.2891	0.5163	0.1298	0.1407
DO [40]	0.2516	0.4620	0.1121	0.1222
SO [46]	0.2879	0.5119	0.1301	0.1408
FD	0.2972	0.5233	0.1346	0.1455
FM	0.3031	0.5338	0.1372	0.1483
FD+FM	<b>0.3125</b>	<b>0.5464</b>	<b>0.1419</b>	<b>0.1531</b>

Table 3: Experiment results on the AAI dataset.

weights. For methods with random dropout rate including *FD* and *DO*, we fix the dropout rate to be 0.3. In Section 4.8, we show the effect of different dropout rates on *FD* and *DO*.

We report the results for Wikipedia and AAI in Table 2 and Table 3 respectively. For each method, we report the metric averaged over five runs. The best result is highlighted in bold. We observe that with **full datasets**:

- Both FD and FM, when used as single SSL task, improve outperforms baseline and DO. This helps answer **RQ1** that the proposed SSL framework and tasks indeed improves model performance for recommenders.
- FD consistently outperforms spread-out regularization (SO) in both datasets. FM outperforms SO on the AAI. Our hypothesis is that the way we construct these auxiliary SSL-based tasks adds additional supervised information to the model training. And could be why we are seeing that they outperform spread-out regularization (SO) approach that does not use additional information.
- By only applying dropout (DO) on the backbone network, we find that model performance on both datasets is not improved. It demonstrates that the gain is from the SSL task with dropout as data augmentation, instead of from only applying feature dropout as regularization.
- We also see encouraging results when combining both tasks together (FD+FM). On the AAI dataset, FM+FD gives the best performance, showing good complementary between the two SSL tasks. To partially address **RQ4**, we observe that FD consistently outperforms FM, and we have seen cases that hybrid approach that combines both FD and FM could potentially outperform FD only.

### 4.4 Data Sparsity

We study the impact of data sparsity on SSL tasks to address **RQ2**. For AAI, we repeat the experiments in section 4.3 after down-sampling the original dataset to 10% and 1%. For Wikipedia dataset, we only down-sample it to 10% as the absolute values of retrieval metrics become too small to draw any conclusions when we just use 1% data.

Table 4 and Table 5 show the model performance on the aforementioned sparser datasets for Wikipedia and AAI respectively. SSL tasks do see larger improvement when sparsity increases. In particular, the best SSL approach (*FD*) on the full Wikipedia dataset improves Recall@10 by 22.7% compared to the baseline, while the relative improvement is 57.4% on the 10% dataset. Similar trend is

10% Dataset				
Method	Recall@10	Recall@50	MAP@10	MAP@50
Baseline	0.0237	0.0924	0.0077	0.0105
DO	0.0272	0.1046	0.0089	0.0120
SO	0.0254	0.0978	0.0083	0.0112
FD	<b>0.0373</b>	<b>0.1294</b>	<b>0.0125</b>	<b>0.0164</b>
FM	0.0261	0.1023	0.0084	0.0115
FD+FM	0.0345	0.1243	0.0115	0.0153

Table 4: Experiment results on 10% of the Wikipedia dataset.

10% Dataset				
Method	Recall@10	Recall@50	MAP@10	MAP@50
Baseline	0.2827	0.4993	0.1274	0.1379
DO	0.2642	0.4810	0.1173	0.1276
SO	0.2879	0.5061	0.1308	0.1413
FD	0.2987	0.5202	0.1364	0.1471
FM	0.2995	0.5235	0.1369	0.1476
FD+FM	<b>0.3023</b>	<b>0.5281</b>	<b>0.1379</b>	<b>0.1487</b>

1% Dataset				
Method	Recall@10	Recall@50	MAP@10	MAP@50
Baseline	0.2600	0.4357	0.1201	0.1287
DO	0.2616	0.4483	0.1194	0.1285
SO	0.2836	0.4838	0.1305	0.1402
FD	0.2832	0.4858	0.1308	0.1406
FM	0.2733	0.4756	0.1254	0.1352
FD+FM	<b>0.3071</b>	<b>0.5128</b>	<b>0.1430</b>	<b>0.1531</b>

Table 5: Experiment results on 10% and 1% of AAI Dataset.

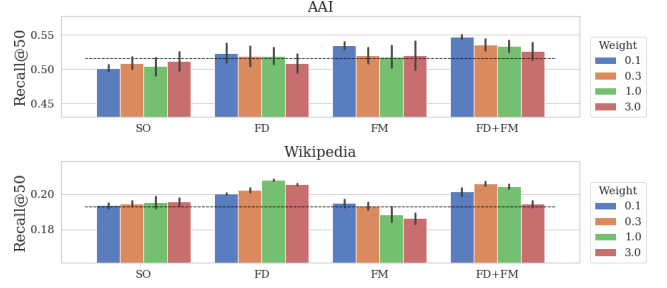
observed for the AAI dataset as shown in Table 3, especially on the 1% dataset although not as obvious as the gains on the Wikipedia dataset. One hypothesis is that the AAI dataset is even more skewed towards popular items as illustrated in Figure 5. In addition, it also includes lots of duplicated pairs of (*query*, *item*), where the pairs in Wikipedia are all unique. When we down-sample the AAI dataset to 10%, many of the (*query*, *item*) pairs are still preserved due to duplication, so the information is still well kept. When the data is down-sampled to 1%, the remaining data is truly sparse of the original supervised information, and thus leads to more obvious improvements from SSL.

It’s worth nothing that, FD consistently outperforms DO and the gap is larger as data becomes sparser. This demonstrates that having dropout for data augmentation in SSL tasks is more effective than directly applying dropout in supervised task.

As a summary, these findings answer research questions raised in RQ2 that the proposed SSL framework improves model performance more with even less supervision.

#### 4.5 Item Slice Analysis

In this section, we further break down the overall performance into head and tail item slices so that we can better understand where does the gains for SSL come from. Our hypothesis is that SSL

Figure 6: Impact of SSL regularization weight  $\alpha$ . The vertical dash line indicates the baseline model’s metric.

generally helps improve the performance for slices without much supervision (e.g., tail items).

Specifically, for AAI, we partition the test data examples into two sub-sets, where the *head* set consists of the top 10% most frequent items, and the rest of the test examples are treated as *tail*. For Wikipedia, we follow the data partitions in [27], where test examples containing items not included in the training set are treated as *tail*, and the rest test examples are treated as *head*. We evaluate on the tail and head test sets separately and report the results in Table 6 and Table 7.

We observe that the proposed SSL methods not only improve the overall performance, but also improve the performance for both head and tail item recommendations. Echo to the data sparsity analysis in Section 4.4, we indeed observe the significantly larger improvements on the tail slice than the head slice. For Wikipedia, the best performing SSL method (FD) improves Recall@10 by 20.97% relatively for tail items, compared to 5.41% on the head items. For AAI, the gap is even larger – the best SSL method improves over 141.42% of the Recall@10 on tail items. These results show that the proposed SSL techniques are effective in alleviating the label sparsity in the long-tail item distribution, and significantly improves the representation learned for tail items.

Note for AAI dataset, since labels were collected from actual user engagement, the ground-truth may suffer from lack of implicit feedback and bias from existing system. We also report live experiment results in Section 5 to further demonstrate the effectiveness of the proposed SSL approaches in a web-scale system.

#### 4.6 Effect of SSL weight $\alpha$

To address RQ3, we study the effect of SSL regularization strength  $\alpha$  presented in equation (3). We perform a line search of  $\alpha$  in {0.1, 0.3, 1.0, 3.0} to study the SSL tasks’ sensitivity to  $\alpha$ .

Figure 4.6 summarizes the Recall@50 evaluated on the Wikipedia and AAI dataset w.r.t. various SSL tasks and  $\alpha$ . We observe that with increasing  $\alpha$ , the model performance is worse than the baseline after certain threshold. This is expected, since large SSL weight  $\alpha$  leads to the multitask loss  $\mathcal{L}$  dominated by  $\alpha \cdot \mathcal{L}_{self}$  in Equation 4. In future work, we plan to reduce the tuning of  $\alpha$  by either providing a prior or exploring the alternative training schedule where query and item representations are first pre-trained with SSL tasks, and then fine-tuned with primary supervised task.



Method	Tail		Head		Overall	
	Recall@10 (Rel %)	Recall@50 (Rel %)	Recall@10 (Rel %)	Recall@50 (Rel %)	Recall@10 (Rel %)	Recall@50 (Rel %)
Baseline	0.0472	0.1621	0.0610	0.2273	0.0537	0.1930
FD	<b>0.0571 (+20.97)</b>	0.1812 (+11.78)	<b>0.0643 (+5.41)</b>	<b>0.2326 (+2.33)</b>	<b>0.0605 (+12.76)</b>	<b>0.2056 (+6.53)</b>
FM	0.0473 (+0.21)	0.1622 (+0.06)	0.0613 (+0.49)	0.228 (+0.31)	0.0539 (+0.37)	0.1934 (+0.21)
FD+FM	0.056 (+18.64)	<b>0.1828 (+12.77)</b>	0.0635 (+4.10)	0.232 (+2.07)	0.0596 (+10.99)	0.2061 (+6.79)

Table 6: Results of Wikipedia on tail, head and overall item slices. Relative improvements over the baseline are in parenthesis.

Method	Tail		Head		Overall	
	Recall@10 (Rel %)	Recall@50 (Rel %)	Recall@10 (Rel %)	Recall@50 (Rel %)	Recall@10 (Rel %)	Recall@50 (Rel %)
Baseline	0.0169	0.0376	0.2976	0.5082	0.2891	0.5163
FD	<b>0.0408 (+141.42)</b>	<b>0.2066 (+449.47)</b>	0.2983 (+0.24)	0.5247 (+3.25)	0.2972 (+2.80)	0.5233 (+1.36)
FM	0.0321 (+89.94)	0.1866 (+396.28)	0.3043 (+2.25)	0.5353 (+5.33)	0.3031 (+4.84)	0.5338 (+3.39)
FD + FM	0.0380 (+124.85)	0.2045 (+443.88)	<b>0.3149 (+5.81)</b>	<b>0.5082 (+8.15)</b>	<b>0.3125 (+8.09)</b>	<b>0.5464 (+5.83)</b>

Table 7: Results of AAI on tail, head and overall item slices. Relative improvements over the baseline are in parenthesis.

#### 4.7 Comparison of Different SSL Tasks

To address RQ4, we compare the performances of different SSL tasks and their combinations. As mentioned in Section 4.3, Table 2 and 3 show that *FD* consistently outperforms *SO* and *DO* in all experiments.

By comparing *FD* with *SO*, where both methods introduce a contrastive learning auxiliary task, we see introducing the feature variants leads to more robust item representation, which is especially true when dataset is sparse.

By comparing *FD* with *DO*, where both methods introduce dropout to increase feature variants, we see the technique is more effective when applied on top of the proposed SSL framework, instead of the primary supervised task.

By comparing a single SSL task (*FD* or *FM*) with the combined approach *FD + FM*, we see *FD* performs the best on Wikipedia while *FD + FM* performs the best on AAI. Our hypothesis is this is due to the *complimentary* of SSL tasks – the *FM* task focuses on learning feature interactions while the *FD* task focuses on feature robustness to variants. The AAI dataset contains a larger number of high cardinality sparse features (such as the *categories* and *title\_unigram*) compared to Wikipedia, thus the complementary effects of combining *FM* and *FD* is more obvious. This also suggests future work in designing more orthogonal or complementary SSL tasks in recommender systems besides *FD* and *FM*, which facilitates more effective feature learning.

#### 4.8 Dropout Rate

In this subsection, we discuss the effect of the dropout rate on the Feature Dropout (*FD*) task, to tackle research question RQ5. Recall that in Table 2 - 5, we report results with fixed dropout rate 0.3. Now we show the results with varied dropout rates. In particular, we experiment with SSL weight  $\alpha = 1.0$  and vary the dropout rate from 0.1 to 0.9. Results are reported in Figure 7.

We observe that, as we increase the dropout rate, the model performance of *DO* continues to deteriorate. For most choices of  $\alpha$  (except  $\alpha = 0.1$ ), *DO* is worse than the baseline (the dashed

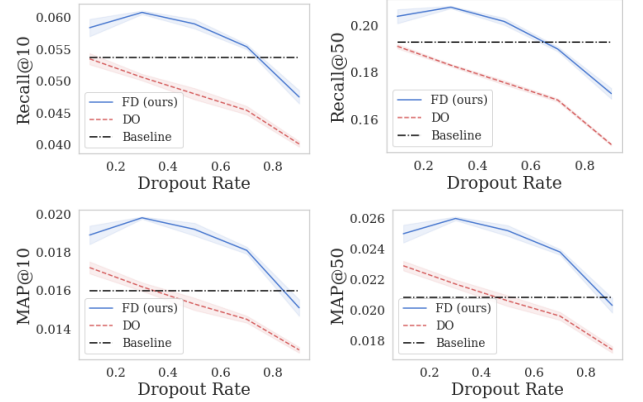
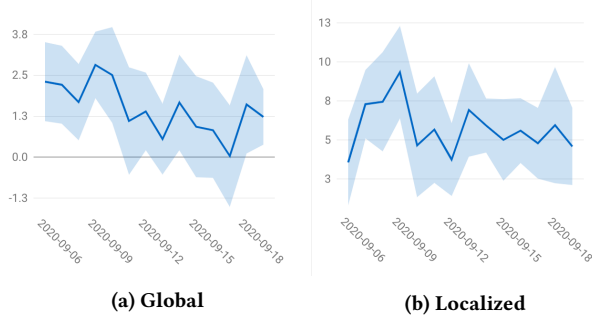


Figure 7: Relationship between dropout rate and evaluation metrics on the Wikipedia dataset.

horizontal line in Figure 7. For the SSL task with *FD*, the model performance peaks when dropout rate is around 0.3 and then deteriorates when we further improve dropout rate. The model starts to under-perform the baseline when we have a dropout rate over 0.7. This observation aligns with our expectation, since when we have dropout rate equal to 0, the *FD* task becomes *SO*; when we have dropout rate equal to 1.0, *FD* task is trained with no features, and thus hurts main task. Moreover, this suggests the effectiveness of combining the dropout on self-supervised task is more effective than directly applying it on the primary supervised task.

### 5 ONLINE EXPERIMENT: SSL-BASED MODELS IN LIVE TRAFFIC

In this section, we describe how we apply our proposed SSL framework to a web-scale commercial app recommender system. Specifically, given an app as the query, the system identifies similar apps given the query. One of the models surfacing this recommendation



**Figure 8: Top Business Metric Improvement Percentage (y-axis) over Days (x-axis) in Online Experiments: (a) improvement globally; (b) improvement on a localized market.**

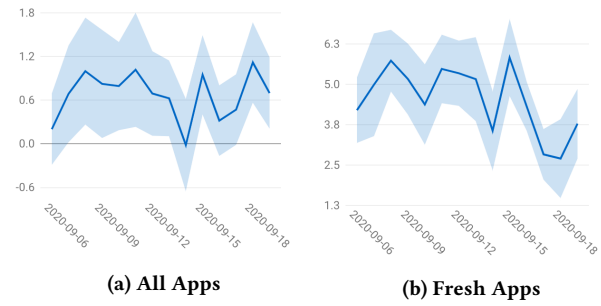
is trained on the AAI dataset as described in Section 4.1, with the same backbone network structure as the two-tower DNN structure in Figure 1 (with modifications). For a natural extension of the offline experiments conducted on Section 4.3 for the AAI experiments, we conducted an A/B experiment for investigating the synergy of deploying the best SSL-based model online. While we already presented improved offline metrics on this dataset, in many real-world systems, offline studies might not align with live impact due to 1) lack of implicit feedback, since the offline evaluation data is collected via user engagement history based on the production system; 2) failing to capture product’s multiple objectives optimization goal, where it’s very likely that recommending more engaging apps hurts other business goals. Therefore, this experiment is critical in demonstrating the effectiveness of the proposed framework in real-world settings.

In our live A/B testing, we add the best performing SSL task with the same set of hyper-parameters on top of the existing well-tuned two-tower DNN model used in production. In a time frame of 14 days, the model improved the overall business metrics significantly, with +0.67% increase in key user engagement (Figure 9a) and +1.5% increase in top business metric (Figure 8a). To echo the study on data sparsity conducted in Section 4.4 and Section 4.5, we see significant improvements on two slices: 1) cold-starting for fresh apps: the model improves +4.5% on user engagement for fresh apps (Figure 9b); and (2) international countries that have sparser training data compared to major markets: we see significant +5.47% top business metric gains (Figure 8b right). Again, both of these results verify our hypothesis that our SSL framework indeed significantly improves model performance for slices without much supervision. Given the results, the SSL empowered model was successfully launched in the current production system.

## 6 CONCLUSION

In this paper, we proposed a model architecture agnostic self-supervised learning (SSL) framework for large-scale neural recommender models. Within the SSL framework, we also proposed two SSL tasks: (i) Feature Masking (FM), and (ii) Feature Dropout (FD), and verified their effectiveness through both offline and online experiments.

For future works, we plan to investigate how different training schemes impact the model quality. One direction is to first pre-train



**Figure 9: Top User Engagement Improvement Percentage (y-axis) over Days (x-axis) in Online Experiments: (a) improvement on all apps; (b) improvement on fresh apps.**

on SSL task to learn query and item representations and fine-tune on primary supervised tasks. Alternatively, it would be interesting to look at SSL with multi-task learning with additional data source with significantly different data distribution. Furthermore, for the two proposed tasks Feature Masking and Feature Dropout, we plan to explore non-parametric based approaches in order to reduce the human intervention in designing feature selection and masking.

## REFERENCES

- [1] Xavier Amatriain and Justin Basilico. 2016. Past, Present, and Future of Recommender Systems: An Industry Perspective. In *Proceedings of the 10th ACM Conference on Recommender Systems* (Boston, Massachusetts, USA) (RecSys ’16). Association for Computing Machinery, New York, NY, USA, 211–214.
- [2] Alex Beutel, Ed H. Chi, Zhiyuan Cheng, Hubert Pham, and John Anderson. 2017. Beyond Globally Optimal: Focused Learning for Improved Recommendations. In *Proceedings of the 26th International Conference on World Wide Web* (Perth, Australia) (WWW ’17). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 203–212.
- [3] L. Beyer, X. Zhai, A. Oliver, and A. Kolesnikov. 2019. S4L: Self-Supervised Semi-Supervised Learning. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 1476–1485.
- [4] Wei-Cheng Chang, Felix X. Yu, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. 2020. Pre-training Tasks for Embedding-based Large-scale Retrieval. In *International Conference on Learning Representations*.
- [5] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Francisco, California, USA) (KDD ’16). Association for Computing Machinery, New York, NY, USA, 785–794.
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. A Simple Framework for Contrastive Learning of Visual Representations. arXiv:2002.05709 <https://arxiv.org/abs/2002.05709>
- [7] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. 2020. Big Self-Supervised Models are Strong Semi-Supervised Learners. arXiv preprint arXiv:2006.10029 (2020).
- [8] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. 2016. Wide & Deep Learning for Recommender Systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems* (Boston, MA, USA) (DLRS 2016). Association for Computing Machinery, New York, NY, USA, 7–10.
- [9] Evangelia Christakopoulou and George Karypis. 2018. Local Latent Space Models for Top-N Recommendation. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (London, United Kingdom) (KDD ’18). Association for Computing Machinery, New York, NY, USA, 1235–1243.
- [10] Edith Cohen and David D. Lewis. 1997. Approximating Matrix Multiplication for Pattern Recognition Tasks. In *Proceedings of the Eighth Annual ACM-SIAM Symposium on Discrete Algorithms* (New Orleans, Louisiana, USA) (SODA ’97). Society for Industrial and Applied Mathematics, USA, 682–691.
- [11] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems* (Boston, Massachusetts, USA) (RecSys ’16). Association for

- Computing Machinery, New York, NY, USA, 191–198.
- [12] Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. 2019. Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems* (Copenhagen, Denmark) (*RecSys '19*). Association for Computing Machinery, New York, NY, USA, 101–109.
  - [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). ACL, 4171–4186.
  - [14] Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. 2014. Discriminative Unsupervised Feature Learning with Convolutional Neural Networks. In *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger (Eds.). Curran Associates, Inc., 766–774.
  - [15] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *J. Mach. Learn. Res.* 12, null (July 2011), 2121–2159.
  - [16] Wikimedia Foundation. [n.d.]. Wikimedia. <https://dumps.wikimedia.org/>
  - [17] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. 2018. Unsupervised Representation Learning by Predicting Image Rotations. In *International Conference on Learning Representations*.
  - [18] Daniel Gillick, Alessandro Presta, and Gaurav Singh Tomar. 2018. End-to-End Retrieval in Continuous Space. *CoRR abs/1811.08008* (2018). [arXiv:1811.08008](http://arxiv.org/abs/1811.08008) <http://arxiv.org/abs/1811.08008>
  - [19] Chuan Guo, Ali Mousavi, Xiang Wu, Daniel N Holtmann-Rice, Satyen Kale, Sashank Reddi, and Sanjiv Kumar. 2019. Breaking the Glass Ceiling for Embedding-Based Classifiers for Large Output Spaces. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 4943–4953.
  - [20] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural Collaborative Filtering. In *Proceedings of the 26th International Conference on World Wide Web* (Perth, Australia) (*WWW '17*). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 173–182.
  - [21] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, and Joaquin Quiñero Candela. 2014. Practical Lessons from Predicting Clicks on Ads at Facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising* (New York, NY, USA). Association for Computing Machinery, New York, NY, USA, 1–9.
  - [22] Jonathan L. Herlocker, Joseph A. Konstan, Loren G. Terveen, and John T. Riedl. 2004. Evaluating Collaborative Filtering Recommender Systems. *ACM Trans. Inf. Syst.* 22, 1 (Jan. 2004), 5–53.
  - [23] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 3146–3154.
  - [24] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. 2019. Revisiting Self-Supervised Visual Representation Learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 1920–1929.
  - [25] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix Factorization Techniques for Recommender Systems. *Computer* 42, 8 (Aug. 2009), 30–37.
  - [26] Yehuda Koren and Robert M. Bell. 2015. *Advances in Collaborative Filtering*. Springer, 77–118.
  - [27] Walid Krichene, Nicolas Mayoraz, Steffen Rendle, Li Zhang, Xinyang Yi, Lichan Hong, Ed Chi, and John Anderson. 2019. Efficient Training on Very Large Corpora via Gramian Estimation. In *International Conference on Learning Representations*.
  - [28] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *International Conference on Learning Representations*.
  - [29] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. 2016. Learning Representations for Automatic Colorization. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV (Lecture Notes in Computer Science, Vol. 9908)*, Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (Eds.). Springer, 577–593.
  - [30] 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* (Perth, Australia) (*WWW '17 Companion*). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 583–592.
  - [31] Jianxin Ma, Chang Zhou, Hongxia Yang, Peng Cui, Xin Wang, and Wenwu Zhu. 2020. Disentangled Self-Supervision in Sequential Recommenders. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) (*KDD '20*). Association for Computing Machinery, New York, NY, USA, 483–491. <https://doi.org/10.1145/3394486.3403091>
  - [32] Klaas Bostels Mark Levy. 2010. Music Recommendation and the Long Tail. In *1st Workshop On Music Recommendation And Discovery (WOMRAD)*, *ACM RecSys*, 2010.
  - [33] Rishabh Mehrotra, Mounia Lalmas, Doug Kenney, Thomas Lim-Meng, and Golli Hashemian. 2019. Jointly Leveraging Intent and Interaction Signals to Predict User Satisfaction with Slate Recommendations. In *The World Wide Web Conference* (San Francisco, CA, USA) (*WWW '19*). Association for Computing Machinery, New York, NY, USA, 1256–1267.
  - [34] Staša Milojević. 2010. Power Law Distributions in Information Science: Making the Case for Logarithmic Binning. *J. Am. Soc. Inf. Sci. Technol.* 61, 12 (Dec. 2010), 2417–2425.
  - [35] Maxim Naumov, Dheevatsa Mudigere, Hao-Jun Michael Shi, Jianyu Huang, Narayanan Sundaraman, Jongsoo Park, Xiaodong Wang, Udit Gupta, Carole-Jean Wu, Alisson G. Azzolini, Dmytro Dzhulgakov, Andrey Malleval, Ilia Cherniavskii, Yinghai Lu, Raghuraman Krishnamoorthi, Ansha Yu, Volodymyr Kondratenko, Stephanie Pereira, Xianjie Chen, Wenlin Chen, Vijay Rao, Bill Jia, Liang Xiong, and Misha Smelyanskiy. 2019. Deep Learning Recommendation Model for Personalization and Recommendation Systems. *CoRR abs/1906.00091* (2019).
  - [36] Wei Niu, James Caverlee, and Haokai Lu. 2018. Neural Personalized Ranking for Image Recommendation. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining* (Marina Del Rey, CA, USA) (*WSDM '18*). Association for Computing Machinery, New York, NY, USA, 423–431.
  - [37] Mehdi Noroozi and Paolo Favaro. 2016. Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VI (Lecture Notes in Computer Science, Vol. 9910)*, Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (Eds.). Springer, 69–84.
  - [38] Shumpei Okura, Yukihiko Tagami, Shingo Ono, and Akira Tajima. 2017. Embedding-Based News Recommendation for Millions of Users. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Halifax, NS, Canada) (*KDD '17*). Association for Computing Machinery, New York, NY, USA, 1933–1942.
  - [39] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. 2015. Unsupervised Learning of Video Representations using LSTMs. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015 (JMLR Workshop and Conference Proceedings, Vol. 37)*, Francis R. Bach and David M. Blei (Eds.). JMLR.org, 843–852.
  - [40] Maksims Volkovs, Guangwei Yu, and Tomi Poutanen. 2017. DropoutNet: Addressing Cold Start in Recommender Systems. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 4957–4966.
  - [41] Zhirong Wu, Yuanjun Xiong, Stella Yu, and Dahua Lin. 2018. Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination. *CoRR abs/1805.01978* (2018). [arXiv:1805.01978](http://arxiv.org/abs/1805.01978) <http://arxiv.org/abs/1805.01978>
  - [42] Xin Xin, Alexandros Karatzoglou, I. Arapakis, and J. Jose. 2020. Self-Supervised Reinforcement Learning for Recommender Systems. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (2020).
  - [43] Yinfei Yang, Steve Yuan, Daniel Cer, Sheng-yi Kong, Noah Constant, Petr Pilar, Heming Ge, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Learning Semantic Textual Similarity from Conversations. In *Proceedings of The Third Workshop on Representation Learning for NLP*. ACL, 164–174.
  - [44] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, and Ed Chi. 2019. Sampling-Bias-Corrected Neural Modeling for Large Corpus Item Recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems* (Copenhagen, Denmark) (*RecSys '19*). Association for Computing Machinery, New York, NY, USA, 269–277.
  - [45] Andrew Zhai, Dmitry Kislyuk, Yushi Jing, Michael Feng, Eric Tzeng, Jeff Donahue, Yue Li Du, and Trevor Darrell. 2017. Visual Discovery at Pinterest. In *Proceedings of the 26th International Conference on World Wide Web Companion* (Perth, Australia) (*WWW '17 Companion*). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 515–524.
  - [46] Xu Zhang, Felix X. Yu, Sanjiv Kumar, and Shih-Fu Chang. 2017. Learning Spread-Out Local Feature Descriptors. In *IEEE International Conference on Computer Vision, Venice, Italy, October 22-29, 2017*. IEEE Computer Society, 4605–4613.
  - [47] Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumthekar, Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. 2019. Recommending What Video to Watch next: A Multitask Ranking System. In *Proceedings of the 13th ACM Conference on Recommender Systems* (Copenhagen, Denmark) (*RecSys '19*). Association for Computing Machinery, New York, NY, USA, 43–51.
  - [48] Kun Zhou, Haibo Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhong yu Wang, and Jirong Wen. 2020. S3-Rec: Self-Supervised Learning for Sequential Recommendation with Mutual Information Maximization. *Proceedings of the 29th ACM International Conference on Information & Knowledge Management* (2020).