



# Learning an Adaptive Meta Model-Generator for Incrementally Updating Recommender Systems

Danni Peng  
danni001@ntu.edu.sg  
Nanyang Technological University  
Singapore

Jie Zhang  
zhangj@ntu.edu.sg  
Nanyang Technological University  
Singapore

Sinno Jialin Pan  
sinnopan@ntu.edu.sg  
Nanyang Technological University  
Singapore

Anxiang Zeng  
zeng0118@ntu.edu.sg  
Nanyang Technological University  
Singapore

## ABSTRACT

Recommender Systems (RSs) in real-world applications often deal with billions of user interactions daily. To capture the most recent trends effectively, it is common to update the model incrementally using only the newly arrived data. However, this may impede the model's ability to retain long-term information due to the potential overfitting and forgetting issues. To address this problem, we propose a novel Adaptive Sequential Model Generation (ASMG) framework, which generates a better serving model from a sequence of historical models via a meta generator. For the design of the meta generator, we propose to employ Gated Recurrent Units (GRUs) to leverage its ability to capture the long-term dependencies. We further introduce some novel strategies to apply together with the GRU meta generator, which not only improve its computational efficiency but also enable more accurate sequential modeling. By instantiating the model-agnostic framework on a general deep learning-based RS model, we demonstrate that our method achieves state-of-the-art performance on three public datasets and one industrial dataset.

## CCS CONCEPTS

• Information systems → Recommender systems; • Computing methodologies → Lifelong machine learning.

## KEYWORDS

Incremental Training, Continual Learning, Meta Learning

### ACM Reference Format:

Danni Peng, Sinno Jialin Pan, Jie Zhang, and Anxiang Zeng. 2021. Learning an Adaptive Meta Model-Generator for Incrementally Updating Recommender Systems. In *Fifteenth ACM Conference on Recommender Systems (RecSys '21)*, September 27–October 1, 2021, Amsterdam, Netherlands. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3460231.3474239>

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

*RecSys '21, September 27–October 1, 2021, Amsterdam, Netherlands*

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8458-2/21/09...\$15.00

<https://doi.org/10.1145/3460231.3474239>

## 1 INTRODUCTION

Recommender Systems (RSs) serve to identify the products or services that could be of potential interest to users by leveraging personal interaction records and various contextual information. With the prevalence of deep learning in recent years, sophisticated deep models are also adopted in RSs owing to their outstanding performance. As the RS models become increasingly complex and cumbersome, how to update them efficiently and effectively in a data-streaming environment (which is often the case in real-world applications) has emerged as a challenge [25, 29].

To capture the most recent interest drift of user and change in item popularity, it is important to periodically update the model as new data arrives. Various methods have been proposed to update conventional RS algorithms in the streaming setting [4, 14, 17, 22, 23]. From the perspective of data utilization, the model can be updated using only the newly arrived data (termed **incremental update**) or using the most recent  $w$  periods of data (termed **batch update**), where  $w$  is the size of the sliding window [12, 13, 25]. In practice, incremental update is often preferred over batch update. Apart from the more efficient training, incremental update can also result in better performance most of the time, as training exclusively with the most recent data can help capture the fast-changing trends and short-term user interests more accurately, which are highly valuable for the near future prediction [1, 6, 9].

However, there exists the problem of forgetting for incremental update, which limits the performance of models updated in this manner [5, 12, 24, 30]. Since the model is only trained to fit the data from the current period, it is difficult for it to retain and consolidate the long-term memory, as the past knowledge learned is constantly overwritten by the new knowledge representing a distinct distribution. In fact, the forgetting issue is not exclusive to RS incremental update. It generally exists when updating neural networks with data from a different distribution [10]. Continual learning is a field of study that specifically tackles this problem, which has developed effective methods for applications in Computer Vision (CV) and Natural Language Processing (NLP). However, when it comes to RSs incremental update, directly adopting these methods may not lead to desirable outcomes. This is because the two problems have different ultimate goals: continual learning aims to perform well for the current task without sacrificing performance on previous tasks, while incremental update of RSs only cares about the performance on future tasks [29]. Hence, the main focus of incremental update

lies in how to effectively transfer past knowledge especially useful for future predictions.

To overcome the specific forgetting issue in RSs incremental update, two lines of research have been developed inspired by some of the methods in continual learning [13], namely **sample-based approach** and **model-based approach**. Despite the effectiveness of both approaches, there exist some major limitations. For the sample-based approach, although the individual samples are carefully selected to be the ‘informative’ ones, they can hardly represent the big picture of the overall distribution. To tackle this, the model-based approach which considers transfer of knowledge between past and present models was recently proposed. However, a major limitation of the existing model-based methods is that, none of them explores the potential of the long-term sequential patterns exist in model evolution, which can be very useful information for generating a better model for future serving.

Inspired by this, we propose a novel **Adaptive Sequential Model Generation (ASMG)** framework for incrementally updating RSs, which encodes a sequence of historical models to generate a better up-to-date serving model via a meta generator. The meta generator is adaptively trained to optimize for the next period data to extract past knowledge that is specially useful for future serving. For the design of the meta generator, we propose to employ Gated Recurrent Units (GRUs) to leverage its ability to capture long-term dependencies. We also introduce some novel strategies to train the GRU meta generator, which not only improve the training efficiency, but also enable better sequential modeling ability. More specifically, we propose to train the GRU meta generator concurrently at multiple steps on the truncated sequence by continuing on a previously learned hidden state. To demonstrate its effectiveness, we instantiate the model-agnostic updating framework on a general deep learning-based Embedding&MLP model, and conduct extensive experiments compared with state-of-the-art updating methods<sup>1</sup>. Our main contributions are summarized as follows:

- We propose an ASMG framework which generates a better model by encoding a sequence of historical models.
- We introduce a GRU-based meta generator design that is capable of capturing the sequential dependencies. We further develop some training strategies to improve its computational efficiency and sequential modeling ability.
- We demonstrate the effectiveness of our method by conducting extensive experiments on three widely-used public datasets and one industrial production dataset from Lazada.

## 2 RELATED WORK

### 2.1 Continual Learning

Methods developed for continual learning generally fall under three categories: experience replay, regularization-based methods, and model fusion. Experience replay prevents forgetting by reusing past samples together with the new data to update the model [18]. Generative methods are later on developed to alleviate the burden of storing real data [21]. Regularization-based methods retain past knowledge by constraining the parameters update based on some

measure of importance [10, 28]. Some works also employ distillation techniques to regularize the update direction [16, 26]. Model fusion supports continual learning of tasks by gradually incorporating sub-networks. Both expandable [19] and fixed-size network methods [11] are developed. Though these methods have shown promising results in CV and NLP, how to adapt them to incrementally update RS models is non-trivial, as they lack mechanism to explicitly optimize for the future performance. To overcome the specific forgetting problem in RSs incremental update, two lines of research have been developed in recent years, which are closely related to some of the methods in continual learning [13].

### 2.2 Sample-based Approach

This line of works relies on reusing historical samples to preserve past memory, analogous to experience replay in continual learning. Most often, a reservoir is maintained to keep a random sample of history [5]. After that, heuristics are designed to select samples from the reservoir for model updating, prioritising recency [1, 30] or the extent of being forgotten [8, 15, 24]. Sample-based approach can be seen as an intermediate between incremental update and training on full historical data, aiming to find a balance between short-term interest and long-term memory. However, an apparent limitation of this approach is that the individual samples, though carefully selected with some measure of informativeness, are insufficient to represent the big picture of overall distribution. In contrast, the historical models learned from the historical data can be seen as a better summarization of the past knowledge. Hence, model-based approach which focuses on transferring knowledge between past and present models was later on developed.

### 2.3 Model-based Approach

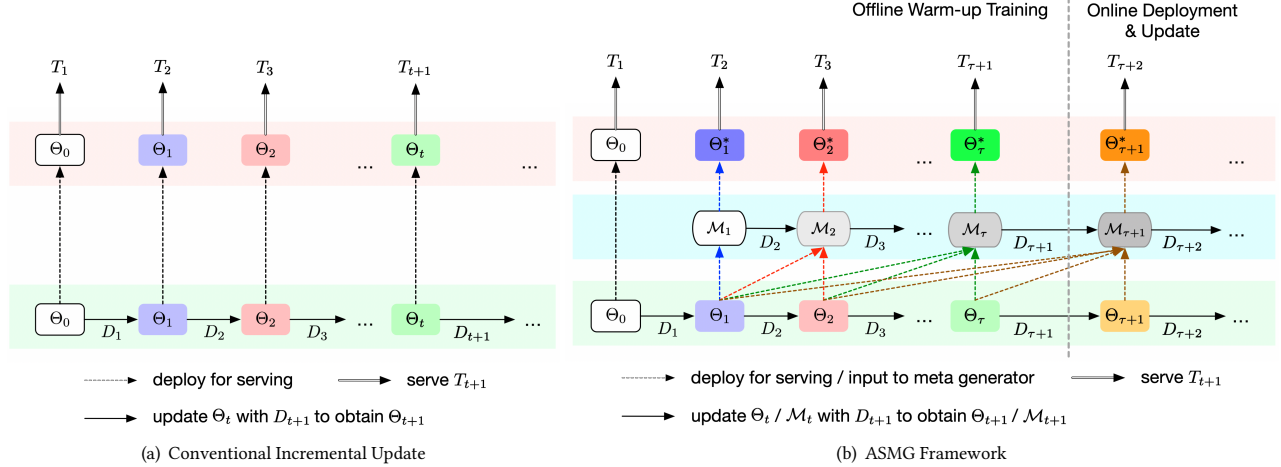
This approach aims to overcome forgetting by transferring knowledge between past and present models. Similar to the regularization-based methods in continual learning, [25] incorporates a knowledge distillation loss to regularize the model update, using previous incremental model as the teacher model. [27] and [13] also employ the distillation techniques and specifically target at GNN-based and session-based RSs respectively. However, simply constraining the change in networks does not guarantee that the knowledge useful for future prediction is extracted. A recently proposed method termed *Sequential Meta-Learning (SML)* [29] achieves this by introducing a mechanism to explicitly optimize for the next period data. It devises a transfer module to combine the past and present models, which is adaptively trained in a sequential manner to optimize for future serving. However, the major limitation of both SML and the distillation methods is that they only consider transfer between models of two consecutive incremental periods, ignoring the long-term sequential patterns that may be highly valuable for generating a better model for future serving.

## 3 METHODOLOGY

### 3.1 Problem Description

Conventional incremental update of RS models involves restoring the last period model as initialization and updating it with the newly arrived data to obtain the most recent model. Formally, let  $D_t$  denote the dataset collected from period  $T_t$  for  $t \in \{1, 2, 3, \dots\}$ ,

<sup>1</sup>Source code can be found at <https://github.com/danni9594/ASMG>.



**Figure 1: Comparison between conventional incremental update and after applying ASMG framework.**

whereby a period can be a certain length of time (e.g., a day, a week) or until a certain amount of data is collected. The model  $\Theta_t$  updated from  $\Theta_{t-1}$  is obtained by minimizing loss on  $D_t$ :

$$\Theta_t = \arg \min_{\Theta} \mathcal{L}(\Theta | D_t, \Theta_{t-1}), \quad (1)$$

where  $\mathcal{L}(\Theta | D_t, \Theta_{t-1})$  is the loss on  $D_t$  with  $\Theta_{t-1}$  as initialization (for recommendation tasks, the loss can be pointwise log loss or pairwise BPR loss). The initial model  $\Theta_0$  can be obtained by random initialization or pre-training with sufficient amount of historical data. Without any post-processing, the updated model  $\Theta_t$  is directly deployed to serve for period  $T_{t+1}$ , from which  $D_{t+1}$  can be obtained for training the next period model  $\Theta_{t+1}$ . The overall procedure of conventional incremental update is illustrated in Figure 1(a).

Although this updating method is widely adopted due to its efficiency and effectiveness, the model trained in this manner can easily overfit to the current data and forget past patterns learned. To address this issue, we propose an Adaptive Sequential Model Generation (ASMG) framework to apply on a sequence of incrementally updated models.

### 3.2 ASMG Framework

For the ASMG framework, instead of direct deployment of the updated model, we introduce a meta generator to take in a sequence of historical models (including the newly updated one), and generate a better model to be deployed for serving the next period. The input sequence of models are obtained from regular incremental update, which means that the model at each period is only trained on the data collected from that period. This allows us to obtain a sequence of models that are highly representative of the respective period. Mining the temporal trends exist in this sequence can have great potential in generating a better model for the next period's serving. Formally, let  $\mathcal{M}_t$  denote the meta generator that takes in a sequence of models until  $\Theta_t$ , the output model  $\Theta_t^*$  used for serving  $T_{t+1}$  is obtained by:

$$\Theta_t^* = \mathcal{M}_t(\Theta_{1:t}), \quad (2)$$

where  $\Theta_{1:t}$  is a sequence of models of length  $t$  from  $\Theta_1$  to  $\Theta_t$ . Note that  $\Theta_0$  is excluded from the input sequence, as it is obtained from

random initialization or pre-training instead of regular incremental update.

The desired properties of an ideal meta generator include 1) good capability of modeling the temporal dependencies exist in the sequence, and 2) stable performance in generating a good model for the next period's serving. The first property is related to the specific network design of the meta generator, which will be discussed in details in the next section. The second property, on the other hand, can be achieved via proper optimization process of the meta generator.

To generate model that is particularly good at serving the next period, we need to ensure that past knowledge that is specially useful for the next period's serving is extracted. To achieve this, we propose to update the meta generator by adaptively optimizing the output model towards the next period data. More specifically, let  $\omega_t$  denote the parameters of  $\mathcal{M}_t$ , i.e.,  $\mathcal{M}_t = \mathcal{M}_{\omega_t}$ , the parameters  $\omega_{t+1}$  are updated from  $\omega_t$  by optimizing  $\Theta_t^*$  towards  $D_{t+1}$ :

$$\begin{aligned} \omega_{t+1} &= \arg \min_{\omega} \mathcal{L}(\Theta_t^* | D_{t+1}, \omega_t) \\ &= \arg \min_{\omega} \mathcal{L}(\mathcal{M}_{\omega}(\Theta_{1:t}) | D_{t+1}, \omega_t), \end{aligned} \quad (3)$$

where  $\mathcal{L}(\Theta_t^* | D_{t+1}, \omega_t)$  is the loss on  $D_{t+1}$  with  $\omega_t$  as initialization. The loss function here is the same as that in (1).

To ensure that the meta generator is sufficiently trained before it can be deployed for online model generation, we perform  $\tau$  periods of meta generator warm-up training prior to the online deployment. During the online deployment phase, the meta generator will still need to be adaptively updated in order to keep up with the most recent sequential trends. The overall procedure of the ASMG framework is depicted in Figure 1(b) and described in Algorithm 1.

### 3.3 GRU Meta Generator

There are many possible designs for the meta generator. In this work, we propose to employ Gated Recurrent Unit (GRUs) [3] to leverage its ability to capture the sequential patterns. For simplicity, we apply the GRU meta generator coordinate-wise on the base model  $\Theta$ . That is to say, for each individual  $\theta \in \Theta$ , we generate the output parameter  $\theta^*$  by applying the GRU meta generator on a sequence of its historical values, independently of parameters

**Algorithm 1:** ASMG Framework

---

**Input:** Sequence of historical datasets  $\{D_t\}_{t=2}^{\tau+1}$  and incrementally updated models  $\{\Theta_t\}_{t=1}^{\tau+1}$   
**Output:** Trained meta generator  $\mathcal{M}_t$  for  $t \geq 2$ , and incrementally updated model  $\Theta_t$  for  $t \geq \tau + 2$   
**Offline Warm-up Training:**  
 Randomly initialize  $\omega_1$   
**for**  $t = 1$  **to**  $\tau$  **do**  
    $\omega_{t+1} \leftarrow \arg \min_{\omega} \mathcal{L}(\mathcal{M}_{\omega}(\Theta_{1:t})|D_{t+1}, \omega_t)$  // update the meta generator  
**end**  
**Online Deployment & Update:**  
**while**  $t \geq \tau + 1$  **do**  
   Generate output model  $\Theta_t^* = \mathcal{M}_t(\Theta_{1:t})$ , deploy  $\Theta_t^*$  to serve for  $T_{t+1}$  and obtain  $D_{t+1}$   
    $\Theta_{t+1} \leftarrow \arg \min_{\Theta} \mathcal{L}(\Theta|D_{t+1}, \Theta_t)$  // update the base model  
    $\omega_{t+1} \leftarrow \arg \min_{\omega} \mathcal{L}(\mathcal{M}_{\omega}(\Theta_{1:t})|D_{t+1}, \omega_t)$  // update the meta generator  
**end**

---

at all the other coordinates. The final serving model is therefore  $\Theta^* = \{\theta^*s\}$ . More specifically, at each step  $i \in \{1, \dots, t\}$ , the hidden state  $h_i \in \mathbb{R}^d$  is computed from the last step hidden state  $h_{i-1} \in \mathbb{R}^d$  and  $\theta \in \mathbb{R}$  of the current step input model  $\Theta_i \in \mathbb{R}^n$ :

$$\begin{aligned}
 r_i &= \sigma(W_r \cdot [h_{i-1}, \theta]), \\
 z_i &= \sigma(W_z \cdot [h_{i-1}, \theta]), \\
 \tilde{h}_i &= \tanh(W_{\tilde{h}} \cdot [r_i \odot h_{i-1}, \theta]), \\
 h_i &= (1 - z_i) \odot h_{i-1} + z_i \odot \tilde{h}_i,
 \end{aligned} \tag{4}$$

where  $r_i, z_i \in \mathbb{R}^d$  are gates that control how much of past and present information to be retained,  $W_r, W_z, W_{\tilde{h}} \in \mathbb{R}^{d \times (d+1)}$  are learnable weights,  $\sigma(\cdot)$  is the sigmoid function,  $[\cdot, \cdot]$  denotes concatenation, and  $\odot$  denotes element-wise multiplication. The initial hidden state  $h_0$  is zero-initialized.

The output parameter  $\theta^* \in \mathbb{R}$  of the final serving model  $\Theta_i^* \in \mathbb{R}^n$  at step  $i$  is obtained from the respective hidden state  $h_i$  via a linear transformation:

$$\theta^* = W \cdot h_i + b, \tag{5}$$

where  $W \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  are also learnable parameters. All the learnable parameters are shared across  $t$  steps.

If we assign one unique GRU meta generator (i.e., unshared set of learnable parameters) at each coordinate in the base model, the space complexity of the GRU meta generator will be  $\mathcal{O}(nd^2)$  where  $d$  is the GRU hidden size and  $n$  is the base model parameter size. Note that the size of the GRU meta generator is unrelated to the sequence length.

### 3.4 Training Strategies

In this section, we introduce two training strategies specifically tailored to the GRU meta generator to further improve its training efficiency and sequential modeling ability.

**3.4.1 Training GRU Meta Generator on Truncated Sequence.** One main drawback of GRUs is that the computation time increases linearly with sequence length. Hence, applying GRU meta generator on the full sequence of historical models is not practical, as the sequence length will continuously increase as time goes by,

i.e., the computational complexity is  $\mathcal{O}(tnd^2)$ , linear in the current time step  $t$ . An alternative is simply to truncate the sequence and start training from an intermediate step with zero-initialized input hidden state. However, using a zero-initialized input hidden state inevitably discards earlier models and prevents the retention of longer-term information.

To tackle this problem, we propose to train the GRU meta generator on truncated sequence by continuing on a previously learned hidden state. By doing so, we not only improve training efficiency with shorter sequence, but also preserve long-term memory by reusing a hidden state learned from the past. For consistency, we fix the sequence length at  $k$  for some  $1 \leq k \leq t$  for all training periods. Specifically, when training  $\mathcal{M}_{t+1}$ , we take in a sequence of  $k$  most recent models  $\Theta_{t-(k-1):t}$  as inputs, and a learned hidden state  $h_{t-k}$  obtained from the training of the previous meta generator  $\mathcal{M}_t$  as the initial hidden state for the current training. Consistently applying this strategy at every training period endows a nice effect that consolidates all the previous models  $\Theta_{1:t-k}$  into hidden state  $h_{t-k}$ . Formally, this strategy modifies the way of generating output model in (2) to  $\Theta_t^* = \mathcal{M}_t(\Theta_{t-(k-1):t}, h_{t-k})$ . The overall optimization is as follows:

$$\begin{aligned}
 \omega_{t+1} &= \arg \min_{\omega} \mathcal{L}(\Theta_t^*|D_{t+1}, \omega_t) \\
 &= \arg \min_{\omega} \mathcal{L}(\mathcal{M}_{\omega}(\Theta_{t-(k-1):t}, h_{t-k})|D_{t+1}, \omega_t).
 \end{aligned} \tag{6}$$

Here,  $\mathcal{M}_{\omega}$  is the GRU meta generator that maps many to one. This strategy is illustrated in Figure 2 from (a) to (b).

#### 3.4.2 Training GRU Meta Generator at Multiple Steps Concurrently.

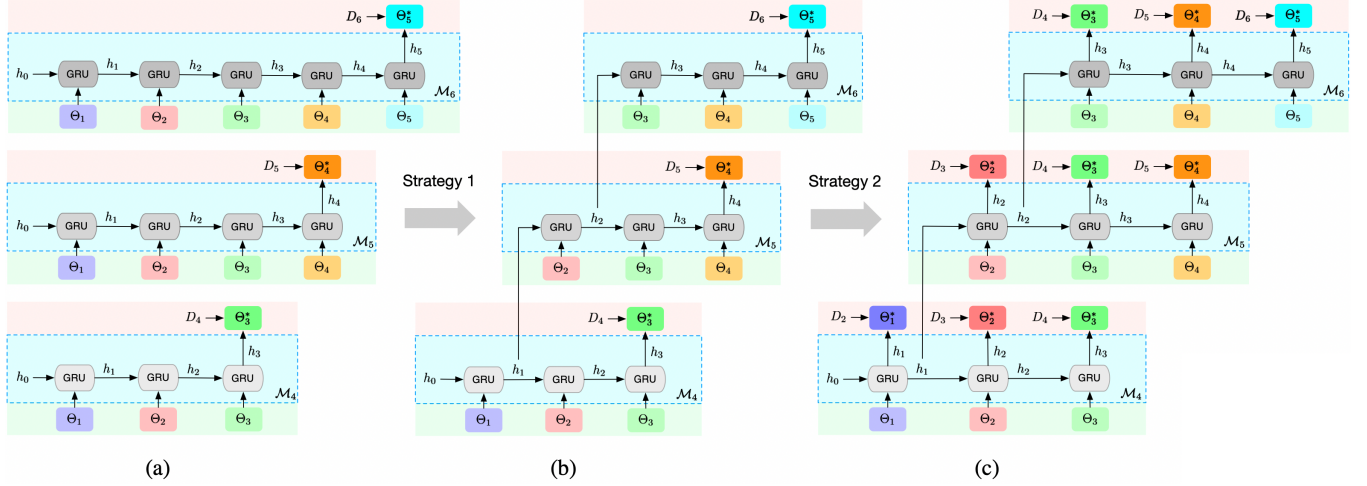
To enable more accurate sequential modeling, we further propose a second strategy to perform concurrent training at multiple steps. The intuition behind is to enforce that the hidden states generated at different steps are on the right track and able to produce model that serves well for the respective next period. More specifically, when training  $\mathcal{M}_{t+1}$ , instead of optimizing the last model  $\Theta_t^*$  towards  $D_{t+1}$  only, we optimize all  $\{\Theta_i^*\}_{i=t-(k-1)}^t$  towards all  $\{D_{i+1}\}_{i=t-(k-1)}^t$  concurrently. To take into account that the later data is less seen before and hence more informative, we assign greater weight  $\lambda$  to loss of the more recent data. Further incorporating this strategy results in the following optimization:

$$\begin{aligned}
 \omega_{t+1} &= \arg \min_{\omega} \sum_{i=t-(k-1)}^t \lambda_i \mathcal{L}(\Theta_i^*|D_{i+1}, \omega_t) \\
 &= \arg \min_{\omega} \sum_{i=t-(k-1)}^t \lambda_i \mathcal{L}(\mathcal{M}_{\omega}(\Theta_{t-(k-1):i}, h_{t-k})|D_{i+1}, \omega_t),
 \end{aligned} \tag{7}$$

where  $\lambda_i$  is the weight for loss  $\mathcal{L}_i = \mathcal{L}(\Theta_i^*|D_{i+1}, \omega_t)$ , computed from a linear decay function, i.e.,  $\lambda_{t-(k-j)} = \frac{j}{\sum_{j'=1}^k j'}$  for  $j \in \{1, 2, \dots, k\}$ .

This strategy is illustrated in Figure 2 from (b) to (c).

Note that for this strategy, we slightly violate the proposed ASMG framework by training the meta generator on data from multiple periods. However, we later on show that this strategy is able to further boost the performance and outperform all the existing model updating methods, including those that also involve multiple periods of data. Furthermore, this strategy should be distinguished from batch update, as the data from earlier periods



**Figure 2: Training of GRU meta generator under the ASMG framework for three consecutive periods (i.e., training of  $M_4$ ,  $M_5$  and  $M_6$ ).** Transition from (a) to (b) illustrates the first strategy, which trains the GRU meta generator on truncated sequence (here we use sequence length  $k = 3$ ) by continuing on a previously learned hidden state. Transition from (b) to (c) illustrates the second strategy, which performs concurrent training by optimizing data at multiple steps, with greater weights assigned to the more recent data. ASMG-GRUsingle and ASMG-GRUmulti are depicted by (b) and (c) respectively.

are used to train the meta generator instead of the base model. The base model  $\Theta_j$ , as the inputs to meta generator, are still trained by regular incremental update.

We term our proposed method that employs GRU meta generator under the ASMG framework as **ASMG-GRU**. We further differentiate the two strategies by naming the first one that trains on a single period of data as **ASMG-GRUsingle**, and the second one that trains on multiple periods of data as **ASMG-GRUmulti**.

### 3.5 Instantiation on Embedding&MLP Base Model

The proposed ASMG-GRU is model-agnostic. To test its effectiveness, we instantiate it on a general deep learning-based Embedding&MLP model, a network architecture that most of the Click-Through Rate (CTR) prediction models developed in recent years are based on [2, 7, 20, 31]. The model mainly comprises embedding layers that transform high-dimensional sparse features into low-dimensional dense vectors, and Multi-Layer Perceptron (MLP) layers that learn the interaction of the concatenated feature embeddings.

More specifically, given user  $i$  with  $P$  categorical features, the concatenated feature embedding  $u_i$  is obtained by:

$$u_i = [E_{u1}z_{i1}, \dots, E_{uP}z_{iP}], \quad (8)$$

where  $E_{up} \in \mathbb{R}^{d_e \times d_{up}}$  with  $d_e \ll d_{up}$  is the embedding matrix for categorical feature  $p \in \{1, \dots, P\}$  of user, and  $z_{ip} \in \mathbb{R}^{d_{up}}$  is the one-hot or multi-hot vector for categorical feature  $p$  of user  $i$ .

Similarly, for item  $j$  with  $Q$  categorical features, the concatenated feature embedding  $v_j$  is obtained by:

$$v_j = [E_{v1}z_{j1}, \dots, E_{vQ}z_{jQ}], \quad (9)$$

where  $E_{vq} \in \mathbb{R}^{d_e \times d_{vq}}$  with  $d_e \ll d_{vq}$  is the embedding matrix for categorical feature  $q \in \{1, \dots, Q\}$  of item, and  $z_{jq} \in \mathbb{R}^{d_{vq}}$  is the one-hot or multi-hot vector for categorical feature  $q$  of item  $j$ .

The predicted preference of user  $i$  on item  $j$  is then obtained by passing the concatenated feature embeddings through a two-layer MLP, followed by a sigmoid function to normalize the score:

$$\hat{y} = \sigma(\text{MLP}([u_i, v_j])). \quad (10)$$

The parameters of the Embedding&MLP base model  $\Theta$  therefore include both the embedding matrices  $\{E_u, E_v\}$ , and the weights and biases of  $\text{MLP}(\cdot)$ . For CTR prediction, we adopt log loss to update the base model  $\Theta$  as described in (1):

$$\mathcal{L}(\Theta|D_t) = -\frac{1}{|D_t|} \left( \sum_{(i,j) \in D_t^+} \log(\hat{y}) + \sum_{(i,j) \in D_t^-} \log(1 - \hat{y}) \right), \quad (11)$$

where  $D_t^-$  is obtained by randomly sampling 1 unobserved interaction for each observed interaction in  $D_t^+$ .

To apply the proposed ASMG-GRU, we assign one unique meta generator to each individual weight and bias parameter in the MLP layers. However, for the embedding layers of sparse ID features (i.e., user ID and item ID), the parameter size is large while the training samples for each ID are limited. Assigning one unique meta generator for each parameter in the embedding layers significantly increases the complexity and can easily lead to overfitting. Hence, we assign one unique meta generator at each dimension of the embedding vector and share it across all IDs. Take for instance, for the  $p$ -th user embedding matrix  $E_{up} \in \mathbb{R}^{d_e \times d_{up}}$ , we assign  $d_e$  distinct meta generators along the first dimension, and share them across  $d_{up}$  IDs along the second dimension. The space complexity of the GRU meta generators for  $E_{up}$  will therefore be significantly reduced from  $O(d_e d_{up} d^2)$  to  $O(d_e d^2)$ , as  $d_{up}$  is often large.

## 4 EXPERIMENTS

### 4.1 Experimental Settings

**4.1.1 Datasets.** Experiments are conducted on three widely-used public datasets and one private industrial production dataset collected from Lazada mobile app.

- **Tmall**<sup>2</sup> contains user-item interaction records from Tmall.com. We extract 31-day data from October 2014, splitting into 31 periods by treating each day as a period. We select the top 50,000 users with the most interactions and filter out items with less than 20 interactions. In the end, the pre-processed dataset contains 3,047,101 observed interactions with 49,986 users and 43,571 items from 634 categories.
- **Sobazaar**<sup>3</sup> is also an e-commerce dataset that contains user-item interaction records from Sobazaar mobile app from September 2014 to December 2014. We split the 4-month data into 31 periods such that each period has equal amount of data. This results in around 4 days per period. In the end, the pre-processed dataset contains 842,660 observed interactions with 1,7126 users and 2,4785 items.
- **MovieLens**<sup>4</sup> contains explicit user ratings for movies on a scale of 0 to 5. We use the largest 25M dataset and extract data from 2014 to 2018. We split the 5-year data into 31 periods based on equal amount of data, which gives around 2 months per period. For movie recommendations, it makes sense to have each period to span longer time (as compared to the e-commerce domain), as the user's taste in movies is more of a long-term interest and less time-sensitive. To make it suitable for our binary classification base model, we follow [24] and label ratings above 3 as positive, and the rest as negative. In the end, the pre-processed dataset contains 6,840,091 samples (including both positive and negative samples) with 43,181 users and 51,142 movies from 20 genres.
- **Lazada** is an industrial production dataset collected from Lazada mobile app, the largest e-commerce platform in South-east Asia. Samples are constructed from the impression logs of users, with clicked records labeled as positive and unclicked records labeled as negative. A 31-day data from December 2020 is extracted and split into 31 periods by treating each day as a period. We randomly select 10,000 users for our experiments. In the end, the pre-processed dataset contains 6,659,828 samples (including both positive and negative samples) with 10,000 users and 43,878 items from 3,860 categories.

For all datasets, we collect the most recent 30 positive feedbacks of each user to form the user sequence feature.

**4.1.2 Baselines.** We compare the following model updating methods applying on the same Embedding&MLP model.

- **Incremental Update (IU):** This method updates the model incrementally using only the new data  $D_t$ .
- **Batch Update (BU-w):** This method updates the model using the most recent  $w$  periods of data  $\{D_{t-(w-1)}, \dots, D_t\}$ .

We experiment the window size  $w$  for 3, 5, and 7 periods. Note that IU is equivalent to BU-1.

- **SPMF** [24]: This method is a sample-based approach. It maintains a reservoir of historical samples, from which the samples with lower predicted ranks are selected to mix with the new data  $D_t$  for current model updating. We tune the reservoir size as a fraction of the total dataset size from 0.1 to 0.5.
- **IncCTR** [25]: This method is a model-based approach that incorporates a knowledge distillation loss when training the current model with  $D_t$ . We implement the version that uses the previous incremental model as teacher model. We tune the weight assigned to the knowledge distillation loss from 0.1 to 0.5.
- **SML** [29]: This method is another model-based approach that has recently achieved the state of the art. It employs a CNN-based transfer module to leverage the previous model while training the current model on  $D_t$ . The transfer of knowledge is only between models of two consecutive periods. SML is originally proposed to instantiate on Matrix Factorization (MF) base model. For a fairer comparison, we also implement this version and term it **SML-MF**. We tune the number of CNN filters in  $\{2, 5, 8, 10\}$ , and the MLP hidden size in  $\{10, 20, 40, 80\}$ .
- **ASMG-Linear:** To compare with the proposed GRU meta generator design under ASMG framework, we implement another fixed-length meta generator that linearly combines the sequence of models, i.e.,  $\Theta_t^* = \sum_{i=t-(k-1)}^t \alpha_i \Theta_i$ , where  $\alpha_i$  is the weight assigned to model at period  $i$ . We tune the sequence length  $k$  from 1 to 6.
- **ASMG-GRU:** This is our proposed method to employ the GRU meta generator under the ASMG framework. Incorporating the proposed training strategies, we compare **ASMG-GRUsingle** and **ASMG-GRUmulti** against the above baselines. We tune the sequence length  $k$  from 1 to 6, and the GRU hidden size in  $\{4, 8, 12, 16\}$ .

For all the compared methods, we use Adam optimizer to update the trainable parameters (i.e., the base model, the meta generator under ASMG, and the transfer module under SML) and tune the learning rate in  $\{1e-2, 1e-3, 1e-4\}$ . We also tune the number of epochs for each update period in  $\{1, 5, 10, 20\}$ . The batch size is set to be 256 for the Sobazaar dataset and 1024 for other datasets.

**4.1.3 Evaluation Protocols.** All datasets are split into 31 periods. The first 10 periods are used to pre-train an initial model  $\Theta_0$ . Model updating will be conducted for the remaining 20 periods. As mentioned earlier in section 3.2, we need to reserve some periods for warm-up training of the meta generator. For input sequence of length  $k$ , the training of the meta generator can only start from period  $10+k$ , as we will need to first obtain a sequence of incrementally updated models of length  $k$  from period 11 to  $10+k$ . Hence, to allow for at least 5 periods of effective warm-up training for a maximum sequence length of  $k=6$  in our experiments, we split the 20 periods into train/validation/test sets by 10/3/7. Evaluation of the model updated at each test period is based on its prediction performance for the respective next period. The metrics used are **AUC** (Area Under ROC) and **LogLoss** (binary cross-entropy), which are

<sup>2</sup><https://tianchi.aliyun.com/dataset/dataDetail?dataId=42>

<sup>3</sup><https://github.com/hainguyen-telenor/Learning-to-rank-from-implicit-feedback>

<sup>4</sup><https://grouplens.org/datasets/movielens/>

good measurements of CTR prediction performance. We conduct 5 runs of model updating experiment for each compared method.

## 4.2 Performance Comparison

Table 1 presents the overall performance for the compared methods, from which we have the following observations:

- BU methods with  $w \geq 2$  are not comparable to IU in terms of serving for the next period, and the degradation increases for larger window size. This observation is consistent with [6, 25], which have identified that in non-stationary environment, batch update is inferior as compared to incremental update. This may be due to that the newly collected data can better represent the latest trends, which are highly related to the upcoming period. Expanding the window size affects the recency of the training data, and hence will be less effective in capturing the latest trends and fast-changing short-term user interests for the near future prediction.
- SPMF and IncCTR only improve over IU marginally. The former is a sample-based approach. The latter is a model-based approach but only utilizes the historical model indirectly via knowledge distillation. The fact that SML performs the best among all the baseline methods demonstrates that directly incorporating the previous model to derive the transfer patterns can be very effective for generating a better current model. Moreover, the mechanism of SML to train the transfer module by optimizing towards the next period data helps extract knowledge that is particularly useful for future prediction.
- SML-MF is not comparable to other methods due to a different choice of base model. This justifies our choice of a deep learning-based Embedding&MLP model which is more sophisticated in modeling complex feature interaction. Another benefit of this base model of our choice is that it is general in the sense that many state-of-the-art RS models are developed from the same architectural paradigm. Being able to achieve success on this model implies that similar effects may also be obtained for other deep learning models with similar architectures.
- Under the proposed ASMG framework, the improvement of ASMG-Linear over IU is minor, and it is sometimes outperformed by other baselines that do not make use of the historical models. This shows that choosing the right design for the meta generator is important, and a simple linear combination is not suitable/sufficient to model the sequential effects. Also, having to drop the models outside of the sliding window limits its ability to accumulate longer-term information.
- Both ASMG-GRUsingle and ASMG-GRUmulti outperform all the baselines by a significant margin of around 1% in AUC on all four datasets. Note that in the industry, 1% increase in offline AUC is considered very significant [7, 31], and it is likely to translate to greater improvement in online CTR [2]. The strong results demonstrate the benefits of 1) directly incorporating the historical models and learning by optimizing towards the future (vs SPMF and IncCTR), 2) involving a longer sequence for more informative sequential

mining (vs SML), and 3) employing the GRU-based meta generator which is a better design for sequential modeling (vs ASMG-Linear).

Figure 3 shows the disentangled performance at each test period. We can see that both ASMG-GRUsingle and ASMG-GRUmulti consistently outperforms all the other methods at each test period. We purposely divide the periods of different datasets to have different time spans, so that we can test the robustness of our method to different model updating frequencies. Furthermore, we observe that there are drastic fluctuations across periods, which indicates the dynamically changing sequential trends and validates the necessity to adaptively update the meta generator.

## 4.3 Ablation Study

In this section, we focus on ASMG-GRUmulti and test the efficacy of its various designs in regard to the training of GRU meta generator. We compare ASMG-GRUmulti with the following variants, each disables a different design:

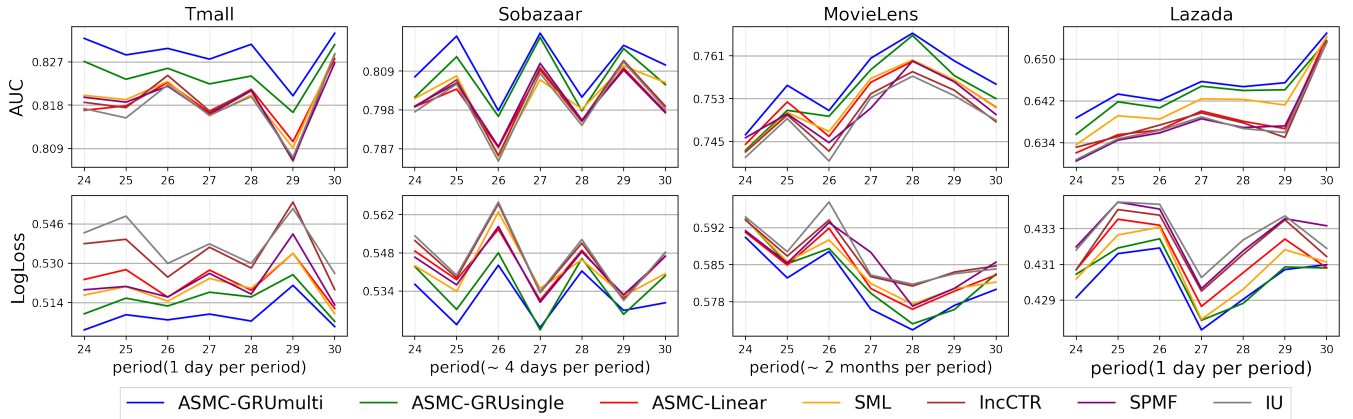
- **ASMG-GRUfull:** This variant trains GRU meta generator on the full sequence of historical models. Illustration of this variant is shown in Appendix A Figure 5(a).
- **ASMG-GRUzero:** This variant uses zero-initialized input hidden state for every period's training instead of continuing on a previously learned hidden state. Illustration of this variant is shown in Appendix A Figure 5(b).
- **ASMG-GRUunif:** This variant assigns uniform weights to the loss at different time steps. Illustration of this variant is the same as ASMG-GRUmulti, as the weights assigned are not explicitly shown in the diagram.
- **ASMG-GRUsingle:** This variant performs single-step training at the last output model towards the newly collected data only. Illustration of this variant is shown in Figure 2(b).

**4.3.1 Prediction Performance.** Table 2 presents the prediction performance of four variants. Firstly, by comparing ASMG-GRUmulti with ASMG-GRUfull, we see that they yield comparable performance on MovieLens and Lazada datasets. This demonstrates that ASMG-GRUmulti is able to maintain the performance while achieving better efficiency by training on truncated sequence. We also observe that for Tmall and Sobazaar, the performance of ASMG-GRUfull slightly drops as compared to ASMG-GRUmulti. This suggests that training on full-length sequence may not always be helpful, as the models at earlier stage may no longer be relevant. Secondly, we see that ASMG-GRUmulti clearly outperforms ASMG-GRUzero. This shows that reusing the previously learned hidden state can help to preserve long-term information. Thirdly, ASMG-GRUmulti also outperforms ASMG-GRUsingle, signifying the benefit of multi-step concurrent training. Finally, ASMG-GRUmulti and ASMG-GRUunif give very similar performance, implying that a simple average of losses may be good enough. The weight decay strategy can be applied on a case-by-case basis.

**4.3.2 Computational Efficiency.** To show that ASMG-GRUmulti serves to improve efficiency as compared to ASMG-GRUfull, we report the training time of GRU meta generator for both methods at the last test period (i.e., period 30) in Table 3. Experiments are conducted using NVIDIA GeForce GTX 1080 Ti with 11GB memory.

**Table 1: Average AUC & LogLoss over 7 test periods for four datasets. The results are reported in mean  $\pm$  std, computed from 5 runs of model updating experiment. imp% indicates the relative improvement over the IU baseline.**

Method	Tmall		Sobazaar		MovieLens		Lazada	
	AUC $\uparrow$	imp%	AUC $\uparrow$	imp%	AUC $\uparrow$	imp%	AUC $\uparrow$	imp%
IU	0.8180 $\pm$ 0.0007	-	0.7998 $\pm$ 0.0007	-	0.7494 $\pm$ 0.0002	-	0.6381 $\pm$ 0.0001	-
BU-3	0.8107 $\pm$ 0.0009	-0.89%	0.7913 $\pm$ 0.0009	-1.06%	0.7379 $\pm$ 0.0003	-1.53%	0.6332 $\pm$ 0.0002	-0.77%
BU-5	0.8002 $\pm$ 0.0009	-2.18%	0.7824 $\pm$ 0.0012	-2.18%	0.7280 $\pm$ 0.0005	-2.86%	0.6287 $\pm$ 0.0004	-1.47%
BU-7	0.7938 $\pm$ 0.0005	-2.96%	0.7781 $\pm$ 0.0007	-2.71%	0.7212 $\pm$ 0.0003	-3.76%	0.6203 $\pm$ 0.0002	-2.79%
SPMF	0.8187 $\pm$ 0.0006	0.09%	0.8007 $\pm$ 0.0004	0.11%	0.7511 $\pm$ 0.0002	0.23%	0.6381 $\pm$ 0.0002	0.00%
IncCTR	0.8190 $\pm$ 0.0007	0.12%	0.8009 $\pm$ 0.0006	0.14%	0.7502 $\pm$ 0.0003	0.11%	0.6388 $\pm$ 0.0003	0.11%
SML	0.8194 $\pm$ 0.0007	0.17%	0.8021 $\pm$ 0.0012	0.29%	0.7522 $\pm$ 0.0009	0.37%	0.6416 $\pm$ 0.0011	0.55%
SML-MF	0.7628 $\pm$ 0.0013	-6.75%	0.7782 $\pm$ 0.0017	-2.70%	0.7242 $\pm$ 0.0012	-3.36%	0.6100 $\pm$ 0.0016	-4.40%
ASMG-Linear	0.8190 $\pm$ 0.0006	0.12%	0.8002 $\pm$ 0.0008	0.05%	0.7524 $\pm$ 0.0002	0.40%	0.6390 $\pm$ 0.0001	0.14%
ASMG-GRUsingle	0.8241 $\pm$ 0.0010	0.75%	0.8055 $\pm$ 0.0017	0.71%	0.7539 $\pm$ 0.0009	0.60%	0.6439 $\pm$ 0.0005	0.91%
ASMG-GRUmulti	<b>0.8289 <math>\pm</math> 0.0010</b>	1.33%	<b>0.8108 <math>\pm</math> 0.0017</b>	1.38%	<b>0.7564 <math>\pm</math> 0.0009</b>	0.93%	<b>0.6452 <math>\pm</math> 0.0005</b>	1.11%
	LogLoss $\downarrow$	imp%	LogLoss $\downarrow$	imp%	LogLoss $\downarrow$	imp%	LogLoss $\downarrow$	imp%
IU	0.5382 $\pm$ 0.0011	-	0.5466 $\pm$ 0.0017	-	0.5871 $\pm$ 0.0002	-	0.4327 $\pm$ 0.0001	-
BU-3	0.5518 $\pm$ 0.0009	-2.53%	0.5536 $\pm$ 0.0013	-1.28%	0.5949 $\pm$ 0.0004	-1.33%	0.4342 $\pm$ 0.0001	-0.35%
BU-5	0.5615 $\pm$ 0.0015	-4.33%	0.5653 $\pm$ 0.0009	-3.42%	0.6056 $\pm$ 0.0007	-3.15%	0.4361 $\pm$ 0.0003	-0.79%
BU-7	0.5732 $\pm$ 0.0012	-6.50%	0.5783 $\pm$ 0.0017	-5.80%	0.6105 $\pm$ 0.0004	-3.99%	0.4379 $\pm$ 0.0002	-1.20%
SPMF	0.5220 $\pm$ 0.0007	3.01%	0.5425 $\pm$ 0.0005	0.75%	0.5857 $\pm$ 0.0001	0.24%	0.4327 $\pm$ 0.0001	0.00%
IncCTR	0.5344 $\pm$ 0.0010	0.71%	0.5458 $\pm$ 0.0007	0.15%	0.5865 $\pm$ 0.0003	0.10%	0.4321 $\pm$ 0.0001	0.14%
SML	0.5198 $\pm$ 0.0009	3.42%	0.5418 $\pm$ 0.0017	0.88%	0.5843 $\pm$ 0.0008	0.48%	0.4309 $\pm$ 0.0003	0.42%
SML-MF	0.5822 $\pm$ 0.0019	-8.18%	0.5713 $\pm$ 0.0021	-4.52%	0.6106 $\pm$ 0.0014	-4.00%	0.4390 $\pm$ 0.0005	-1.46%
ASMG-Linear	0.5226 $\pm$ 0.0012	2.90%	0.5430 $\pm$ 0.0013	0.66%	0.5840 $\pm$ 0.0002	0.53%	0.4314 $\pm$ 0.0000	0.30%
ASMG-GRUsingle	0.5154 $\pm$ 0.0018	4.24%	0.5370 $\pm$ 0.0017	1.76%	0.5828 $\pm$ 0.0011	0.73%	0.4304 $\pm$ 0.0003	0.53%
ASMG-GRUmulti	<b>0.5079 <math>\pm</math> 0.0018</b>	5.63%	<b>0.5309 <math>\pm</math> 0.0017</b>	2.87%	<b>0.5806 <math>\pm</math> 0.0011</b>	1.11%	<b>0.4299 <math>\pm</math> 0.0003</b>	0.65%

**Figure 3: Prediction performance at each test period for four datasets, averaged over 5 runs of model updating experiment. Note that we omit BU and SML-MF here, as they are not comparable to the rest.**

For ASMG-GRUmulti, we set sequence length  $k = 3$ . For ASMG-GRUfull, the sequence length increases to 20 at period 30. From the results, we see that the training time of ASMG-GRUfull is about 6.5~7 times longer than that of ASMG-GRUmulti, consistent with the fact that the training time is linear in sequence length. The results empirically show that fixing the sequence length can greatly improve the computational efficiency, especially at the later training periods.

#### 4.4 Sensitivity Analysis

We further conduct sensitivity analysis to investigate the effects of two important factors, input sequence length and GRU hidden size,

on the performance of ASMG-GRUmulti. Here we only present the analysis in AUC. For results in LogLoss, see Appendix B.

**4.4.1 Effect of Input Sequence Length.** To study the effect of sequence length, we fix the hidden size at 4 and vary the sequence length from 1 to 6. Figure 4(a) shows the performance of ASMG-GRUmulti *w.r.t* different sequence lengths. Firstly, we see that ASMG-GRUmulti is able to outperform all the baselines regardless of sequence length. Even for  $k = 1$  (i.e., applying ASMG-GRUmulti on the most recent model only), ASMG-GRUmulti is still slightly better than the strongest baseline SML, which is a model-based method that considers transfer between 2 periods (i.e., sequence length is 2). This can be attributed to the GRU design which better

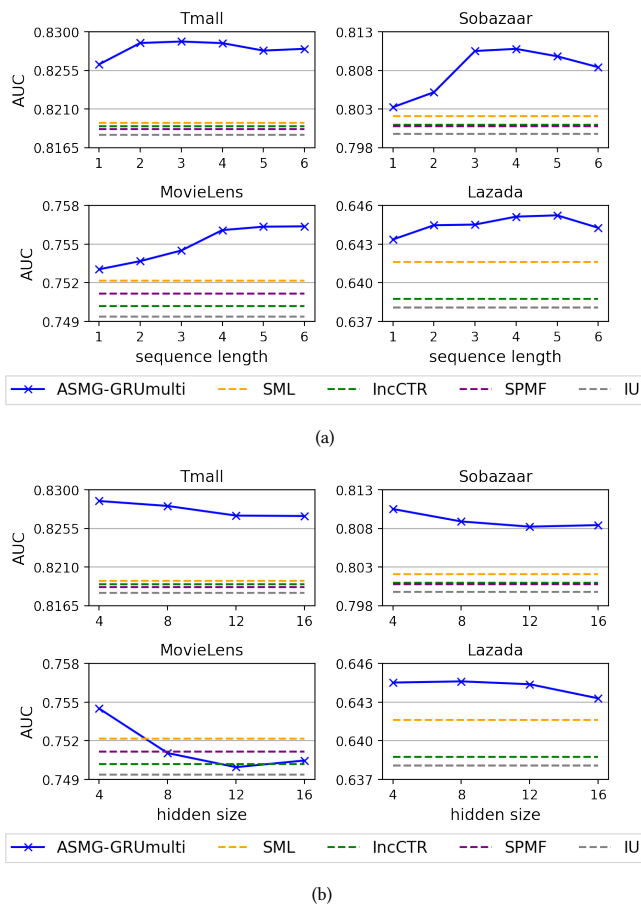


**Table 2: Prediction performance (average AUC & LogLoss over 7 test periods) of ASMG-GRUmulti and its four variants.**

Variant	Tmall		Sobazaar		MovieLens		Lazada	
	AUC $\uparrow$	LogLoss $\downarrow$	AUC $\uparrow$	LogLoss $\downarrow$	AUC $\uparrow$	LogLoss $\downarrow$	AUC $\uparrow$	LogLoss $\downarrow$
ASMG-GRUfull	0.8267	0.5108	0.8083	0.5323	<b>0.7565</b>	0.5811	0.6452	0.4299
ASMG-GRUzero	0.8224	0.5162	0.8079	0.5343	0.7550	0.5818	0.6440	0.4303
ASMG-GRUunif	0.8284	0.5102	0.8091	0.5324	0.7563	0.5811	0.6449	0.4300
ASMG-GRUsingle	0.8241	0.5154	0.8055	0.5370	0.7539	0.5828	0.6439	0.4304
ASMG-GRUmulti	<b>0.8289</b>	<b>0.5079</b>	<b>0.8108</b>	<b>0.5309</b>	0.7564	<b>0.5806</b>	<b>0.6452</b>	<b>0.4299</b>

**Table 3: Training time (in minutes) of GRU meta generator for the two methods at the last test period (i.e., period 30).**

	Tmall	Sobazaar	MovieLens	Lazada
ASMG-GRUfull	93.6	23.1	59.3	42.6
ASMG-GRUmulti	13.8	3.4	8.7	6.2

**Figure 4: Prediction performance (average AUC over 7 test periods) of ASMG-GRUmulti w.r.t (a) different input sequence lengths and (b) different GRU hidden sizes.**

models the sequential patterns, and also the strategy to continue from a previously learned hidden state, which helps retain longer memory. Secondly, we observe from the trends that as long as the sequence length is long enough for effective training of GRUs (e.g., the minimum length should be 2 for Tmall and Lazada, 3 for Sobazaar, and 4 for MovieLens), the performance tends to stabilize

at a satisfactory level. However, it is also noticed that further increasing the sequence length may degrade the performance (e.g., downward trend observed for Sobazaar after 4 and Lazada after 5). Longer sequence length is also not desirable due to the expensive computation. Therefore, the optimal range of sequence length should be from 3 to 5.

**4.4.2 Effect of GRU Hidden Size.** We also study the effect of GRU hidden size on the performance of ASMG-GRUmulti. By fixing the sequence length at 3, we vary the hidden size in  $\{4, 8, 12, 16\}$ . The results are shown in Figure 4(b). We can see that the performance generally drops for more than 4 hidden units. For MovieLens, the degradation is so severe that the performance becomes worse than some of the baselines. This implies that overfitting may occur if the meta generator is over-parameterized by increasing the GRU hidden size. Furthermore, complex GRU meta generator is also expensive to train. Hence, we suggest that setting the hidden size to 4 is a good starting point.

## 5 CONCLUSION

In this work, we study the model updating strategies for RSs, aiming to achieve recency with incremental update while still being able to prevent forgetting of past knowledge. To this end, we propose a novel ASMG framework, which produces a better serving model by encoding a sequence of incrementally updated models in the past via a meta generator. We introduce a meta generator design based on GRUs and further develop some strategies to improve its training efficiency and sequential modeling ability. By instantiating the proposed method on a general deep learning-based RS model, we demonstrate that it achieves the state of the art on four real-world datasets. In the future, we will investigate more effective meta generator designs considering different base model architectures.

## ACKNOWLEDGMENTS

This research is supported, in part, by Alibaba Group through Alibaba Innovative Research (AIR) Program and Alibaba-NTU Singapore Joint Research Institute (JRI), Nanyang Technological University, Singapore.

## REFERENCES

- [1] Chen Chen, Hongzhi Yin, Junjie Yao, and Bin Cui. 2013. Terec: A temporal recommender system over tweet stream. *Proceedings of the VLDB Endowment* 6, 12 (2013), 1254–1257.
- [2] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & Deep Learning for Recommender Systems. In *DLRS@ RecSys*.
- [3] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS 2014 Workshop on Deep Learning, December 2014*.
- [4] Robin Devooght, Nicolas Kourtellis, and Amin Mantrach. 2015. Dynamic matrix factorization with priors on unknown values. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 189–198.
- [5] Ernesto Diaz-Aviles, Lucas Drumond, Lars Schmidt-Thieme, and Wolfgang Nejdl. 2012. Real-time top-n recommendation in social streams. In *Proceedings of the sixth ACM conference on Recommender systems*. 59–66.
- [6] Erzsébet Frigó, Róbert Pálóvics, Domokos Kelen, Levente Kocsis, and András Benczúr. 2017. Online ranking prediction in non-stationary environments. (2017).
- [7] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 1725–1731.
- [8] Lei Guo, Hongzhi Yin, Qinyong Wang, Tong Chen, Alexander Zhou, and Nguyen Quoc Viet Hung. 2019. Streaming session-based recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1569–1577.
- [9] Michael Jugovac, Dietmar Jannach, and Mozghan Karimi. 2018. Streamingrec: a framework for benchmarking stream-based news recommenders. In *Proceedings of the 12th ACM Conference on Recommender Systems*. 269–273.
- [10] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526.
- [11] Arun Mallya and Svetlana Lazebnik. 2018. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7765–7773.
- [12] Fei Mi and Boi Faltings. 2020. Memory Augmented Neural Model for Incremental Session-based Recommendation. *arXiv preprint arXiv:2005.01573* (2020).
- [13] Fei Mi, Xiaoyu Lin, and Boi Faltings. 2020. Ader: Adaptively distilled exemplar replay towards continual learning for session-based recommendation. In *Fourteenth ACM Conference on Recommender Systems*. 408–413.
- [14] Manos Papagelis, Ioannis Rousidis, Dimitris Plexousakis, and Elias Theoharopoulos. 2005. Incremental collaborative filtering for highly-scalable recommendation algorithms. In *International Symposium on Methodologies for Intelligent Systems*. Springer, 553–561.
- [15] Ruihong Qiu, Hongzhi Yin, Zi Huang, and Tong Chen. 2020. Gag: Global attributed graph neural network for streaming session-based recommendation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 669–678.
- [16] Sylvester-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. 2017. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 2001–2010.
- [17] Steffen Rendle and Lars Schmidt-Thieme. 2008. Online-updating regularized kernel matrix factorization models for large-scale recommender systems. In *Proceedings of the 2008 ACM conference on Recommender systems*. 251–258.
- [18] Anthony Robins. 1995. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science* 7, 2 (1995), 123–146.
- [19] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. 2016. Progressive neural networks. *arXiv preprint arXiv:1606.04671* (2016).
- [20] Ying Shan, T Ryan Hoens, Jian Jiao, Haijing Wang, Dong Yu, and JC Mao. 2016. Deep crossing: Web-scale modeling without manually crafted combinatorial features. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 255–262.
- [21] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. 2017. Continual learning with deep generative replay. In *Advances in neural information processing systems*. 2990–2999.
- [22] João Vinagre, Alípio Mário Jorge, and João Gama. 2014. Fast incremental matrix factorization for recommendation with positive-only feedback. In *International Conference on User Modeling, Adaptation, and Personalization*. Springer, 459–470.
- [23] JianGuo Wang, Joshua Zhexue Huang, Dingming Wu, Jiafeng Guo, and Yanyan Lan. 2016. An incremental model on search engine query recommendation. *Neurocomputing* 218 (2016), 423–431.
- [24] Weiqing Wang, Hongzhi Yin, Zi Huang, Qinyong Wang, Xingzhong Du, and Quoc Viet Hung Nguyen. 2018. Streaming ranking based recommender systems. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 525–534.
- [25] Yichao Wang, Huifeng Guo, Ruiming Tang, Zhirong Liu, and Xiuqiang He. 2020. A Practical Incremental Method to Train Deep CTR Models. *arXiv preprint arXiv:2009.02147* (2020).
- [26] Junfeng Wen, Yanshuai Cao, and Ruitong Huang. 2018. Few-shot self reminder to overcome catastrophic forgetting. *arXiv preprint arXiv:1812.00543* (2018).
- [27] Yishi Xu, Yingxue Zhang, Wei Guo, Huifeng Guo, Ruiming Tang, and Mark Coates. 2020. GraphSAIL: Graph Structure Aware Incremental Learning for Recommender Systems. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2861–2868.
- [28] Friedemann Zenke, Ben Poole, and Surya Ganguli. 2017. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*. PMLR, 3987–3995.
- [29] Yang Zhang, Fuli Feng, Chenxu Wang, Xiangnan He, Meng Wang, Yan Li, and Yongdong Zhang. 2020. How to retrain recommender system? A sequential meta-learning method. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1479–1488.
- [30] Yan Zhao, Shoujin Wang, Yan Wang, and Hongwei Liu. 2020. Stratified and time-aware sampling based adaptive ensemble learning for streaming recommendations. *Applied Intelligence* (2020), 1–21.
- [31] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1059–1068.

### A ILLUSTRATIONS OF VARIANTS

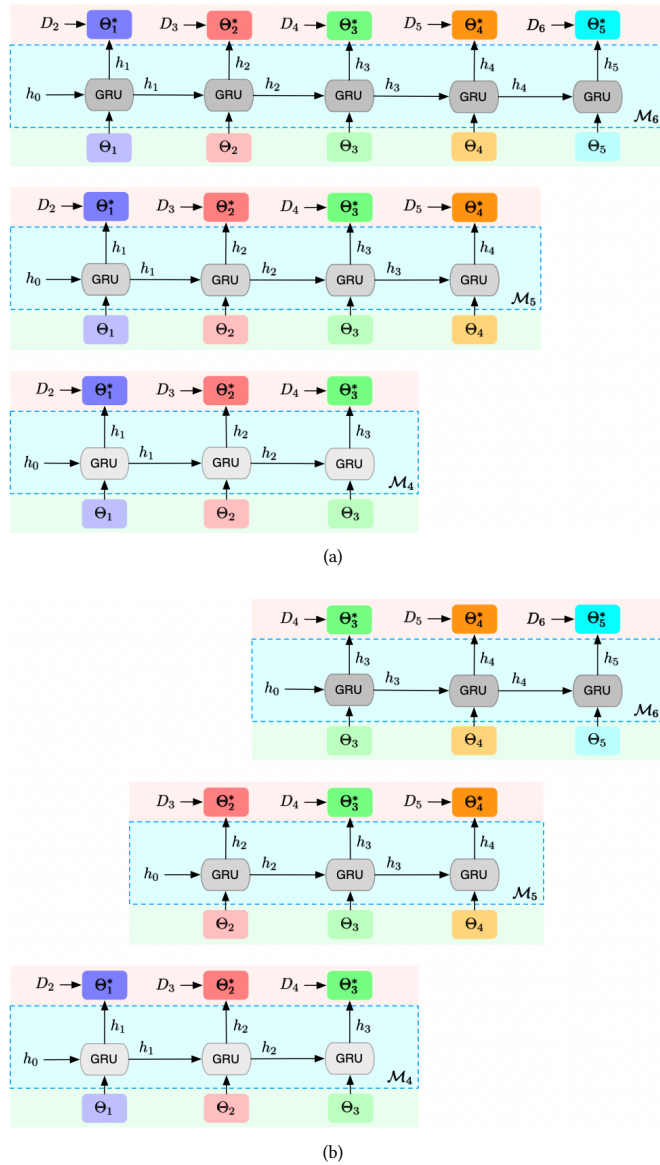


Figure 5: (a) ASMG-GRUfull trains GRU meta generator on the full sequence of historical models. (b) ASMG-GRUzero uses zero-initialized input hidden state for every period's training (i.e., same input hidden state  $h_0$  for the training of  $\mathcal{M}_4$ ,  $\mathcal{M}_5$  and  $\mathcal{M}_6$ ) instead of continuing on a previously learned hidden state.

### B SENSITIVITY ANALYSIS IN LOGLOSS

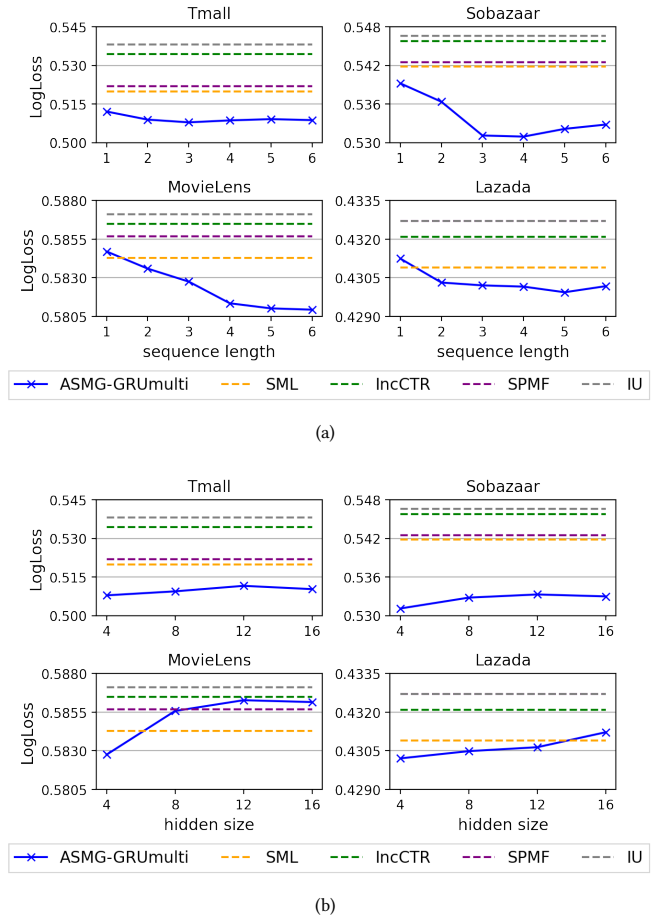


Figure 6: Prediction performance (average LogLoss over 7 test periods) of ASMG-GRUmulti w.r.t (a) different input sequence lengths and (b) different GRU hidden sizes.