

关于长序列用户行为建模预测点击率的实践

皮琦，卞伟杰，周国瑞，朱小强，盖坤*。

阿里巴巴集团 北

京，中国

{piqi.pq, weijie.bwj, guorui.xgr, xiaoqiang.zxq, jingshi.gk}@alibaba-inc.com

ABSTRACT

点击率（CTR）的预测对于推荐系统和在线广告等工业应用至关重要。实际上，通过从丰富的历史行为数据中挖掘用户兴趣，对这些应用中的点击率建模起着重要作用。在深度学习发展的推动下，人们提出了设计巧妙的用户兴趣建模架构的深度CTR模型，为模型性能带来了比常规指标更显著的改善。然而，要将这些复杂的模型部署到在线服务系统中进行实时推理，面对大量的流量请求，需要做出巨大的努力。当涉及到长序列的用户行为数据时，情况会变得更加困难，因为系统延迟和存储成本会随着用户行为序列的长度而线性增加。

在本文中，我们直接面对长序列用户行为建模的挑战，并介绍了我们为CTR预测任务共同设计机器学习算法和在线服务系统的实践。(i) 从服务系统来看，我们通过设计一个名为UIC（用户兴趣中心）的独立模块，将用户兴趣模型中最耗费资源的部分与整个模型解耦。UIC维护每个用户的最新兴趣状态，其更新取决于实时的用户行为触发器事件，而不是流量请求。因此，UIC对于实时CTR预测是无延迟的。(ii) 从机器学习算法的角度来看，我们提出了一个新的基于内存的架构，名为MIMN（多通道用户兴趣记忆网络），从长的连续行为数据中捕捉用户兴趣，实现了比最先进的模型更优越的性能。MIMN是以UIC模块的增量方式实现的。

理论上，UIC和MIMN的共同设计方案使我们能够处理具有无限长度的连续行为数据的用户兴趣建模。模型性能和系统效率之间的比较证明了所提方案的有效性。据我们所知，这是第一个能够处理长序列用户行为数据的工业解决方案之一，其长度可扩展到数千。它现在已经被部署在阿里巴巴的显示广告系统中。

*Q.Pi和W. Bian为共同第一作者。通讯作者为G. Zhou。

允许为个人或课堂使用本作品的全部或部分内容制作数字或硬拷贝，不收取任何费用，但拷贝不得以营利或商业利益为目的而制作或分发，且拷贝首页须注明本通知和完整的引文。必须尊重ACM以外的其他人拥有的本作品的版权。允许摘录并注明出处。以其他方式复制，或重新发表，张贴在服务器上或重新分发到名单上，需要事先获得特别许可和/或付费。请从permissions@acm.org申请许可。

KDD '19, 2019年8月4-8日，美国安克雷奇，AK，USA

© 2019年美国计算机协会。ACM ISBN 978-1-

4503-6201-6/19/08... \$15.00

<https://doi.org/10.1145/3292500.3330666>

CCS概念

• 信息系统 → 推荐系统; 在线广告.

关键字

点击率预测; 用户行为建模

ACM参考格式：

皮琦，卞伟杰，周国瑞，朱小强，盖坤. 2019. 长序列用户行为建模预测点击率的实践. In *The 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '19)*, August 4-8, 2019, Anchorage, AK, USA. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3292500.3330666>

1 简介

不断增长的互联网将我们带入一个具有个性化在线服务的数字世界。从在线系统中收集到的大量用户行为数据为我们提供了更好地了解用户偏好的绝佳机会。从技术上讲，从丰富的行为数据中捕捉用户的兴趣是非常重要的，因为它有助于显著改善典型的现实世界的应用，如推荐系统和在线广告[8, 30, 31]。在本文中，我们将自己限制在点击率（CTR）预测建模的任务上，这在在线服务中起着关键作用。这里讨论的解决方案也适用于许多相关的任务，如转换率预测和用户偏好建模。在深度学习的推动下，人们提出了具有巧妙设计的用户兴趣建模架构的深度CTR模型，达到了最先进的水平。这些模型可以大致分为两类：（1）基于池化的架构[4, 8, 31]，将用户的历史行为视为独立的信号，并应用sum/max/attention等池化操作来总结用户兴趣表示；（2）顺序建模架构[21, 30]，将用户行为视为顺序信号，应用LSTM/GRU操作来总结用户兴趣。然而，在工业应用中，需要付出巨大的努力，将这些复杂的模型部署到在线服务系统中进行实时推理，每天有数以亿计的用户访问该系统。当遇到极长的连续用户行为数据时，事情会变得更加困难，因为上述所有的模型都需要存储整个用户行为序列、

a.k.a. 特征，在在线服务系统中，并在严格的延迟范围内获取它们来计算兴趣表示。这里的“长”是指连续的用户行为的长度达到1000或更多。实际上，系统延迟和存储成本随着用户行为序列的长度而近似地线性增加。正如[30]中所报告的，部署顺序模型需要做很多工程工作，它只是处理具有最大限度的用户行为

长度为50。图1显示了在阿里巴巴的线上展示广告系统中, 用户行为序列的平均长度和相应的CTR模型性能。显然, 长序列用户行为建模的挑战是值得解决的。

在本文中, 我们介绍了我们在机器学习算法和在线服务系统的共同设计方面的实践。

我们将用户行为建模模块与整个点击率预测系统解耦, 并相应地设计出具体解决方案。

(i) **服务系统的观点**: 我们设计了一个独立的UIC (用户兴趣中心) 模块。UIC专注于用户行为建模的在线服务问题, 它维护每个用户的最新兴趣表示。UIC的一个关键点是其更新机制。用户状态的更新只取决于实时的用户行为触发事件, 而不是流量请求。也就是说, UIC对于实时CTR预测是没有延迟的。(ii) **机器学习算法观点**: 仅仅解耦UIC模块不能处理存储问题, 因为对于数以亿计的用户来说, 用户行为序列的长度扩展到1000, 存储和推理仍然相当困难。在这里, 我们借鉴了NTM[11]的记忆网络的理念, 提出了一个名为MIMN (多通道用户兴趣记忆网络) 的新型架构。MIMN以增量的方式工作, 可以很容易地用UIC模块实现。这有助于解决存储方面的挑战。此外, MIMN通过 *内存利用正则化* 和 *内存诱导单元* 的两种设计改进了传统的NTM, 使其在有限的存储空间下更有效地对用户行为序列进行建模, 并在模型性能上带来了显著的增益。

从理论上讲, 结合UIC和MIMN为我们提供了一个处理用户兴趣建模的方法, 即无限长度的连续行为数据。我们的实验表明, 所提出的解决方案在模型性能和系统效率方面都很出色。据我们所知, 这是第一个能够处理长度可达数千的连续用户行为数据的工业解决方案之一。

这项工作的主要贡献总结如下:

- 我们介绍了一个针对CTR预测任务的学习算法和服务系统的共同设计的实践。这个解决方案已经部署在一个世界领先的广告系统中, 为我们带来了处理长序列用户行为建模的能力。
- 我们设计了一个新的UIC模块, 它将繁重的用户兴趣计算与整个CTR预测过程相分离。UIC对于流量请求来说是无延迟的, 并允许在离线模式下进行复杂的模型计算和实时推理。
- 我们提出了一个新的MIMN模型, 它通过 *内存利用正则化* 和 *内存诱导单元* 的两种设计改进了原来的NTM架构, 使其更适合于用户兴趣学习。MIMN很容易与UIC服务器结合起来, UIC服务器可以增量更新用户的兴趣表示。
- 我们在公共数据集和从阿里巴巴的广告系统中收集的工业数据集上进行了仔细的实验。我们还详细分享了我们在部署所提出的解决方案的实际问题上的经验。我们相信这将有助于推动社区的发展。

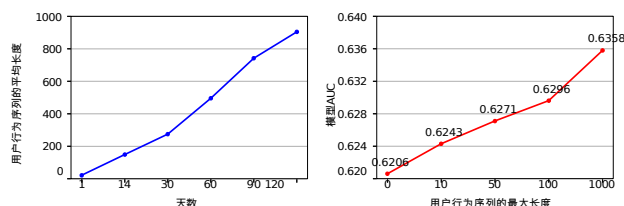


图1: 连续的用户行为数据的统计和相关的性能。在阿里巴巴的展示广告系统中的响应模型性能。

2 相关的工作

深度CTR模型。随着深度学习的快速发展, 我们已经在许多领域取得了进展, 如计算机视觉[15], 自然语言处理[2]。在这些成功的启发下, 一些基于深度学习的CTR预测方法[4, 8, 24, 29]已经被提出。这些方法不是传统方法中的特征工程, 而是利用神经网络来捕捉特征的相互作用。尽管想法似乎很简单, 但这些工作使CTR预测任务的发展向前迈进了一大步。此后, 业界更加关注模型架构设计, 而不是通过详尽的特征工程来提高性能。除了学习特征的相互作用, 越来越多的方法被提出来, 从丰富的历史行为数据中获取用户的洞察力。DIN[31]指出, 用户的兴趣是多种多样的, 并且随着项目的不同而变化。DIN中引入了关注机制来捕获用户的兴趣。DIEN[30]提出了一个辅助损失来捕捉来自具体行为的潜在兴趣, 并完善了GRU[6]来模拟兴趣的演变。

长期的用户兴趣。[17]认为长期兴趣是指一般的兴趣, 它回到了一个人的脑海中, 对个性化很重要。[19]提出对用户的长期兴趣分类进行建模。[5]逐步建立长期和短期的用户档案得分来表达用户的兴趣。所有这些方法都是通过特征工程来建立长期兴趣模型, 而不是通过自适应的端对端学习。TDSSM[25]提出对长期和短期用户兴趣进行联合建模以提高推荐质量。不幸的是, 这些基于深度学习的方法, 如TDSSM、DIN、DIEN很难被部署在面临极长的用户行为序列的实时预测服务器中。来自存储的压力, 计算的延迟将随着用户行为序列的长度而线性增长。在工业应用中, 行为序列的长度通常较小, 例如50个, 而在淘宝网上, 一个活跃的用户可能在两周内留下的行为, 如点击、转换等, 长度超过1000个。**记忆网络**。记忆网络[11, 26]已经被用来提取具有外部记忆成分的知识。这个想法已经被广泛地应用于NLP, 例如问题回答系统。一些工作[3, 10, 16, 27]利用记忆网络进行用户兴趣建模。然而, 这些方法忽略了长期利益模型和实际部署问题。

3 实时交易中心预测系统

在现实世界的推荐人或广告系统中, CTR预测模块作为一个关键的组成部分工作[8, 31]。通常情况下, 它接收一组候选人 (如项目或广告), 并相应地返回

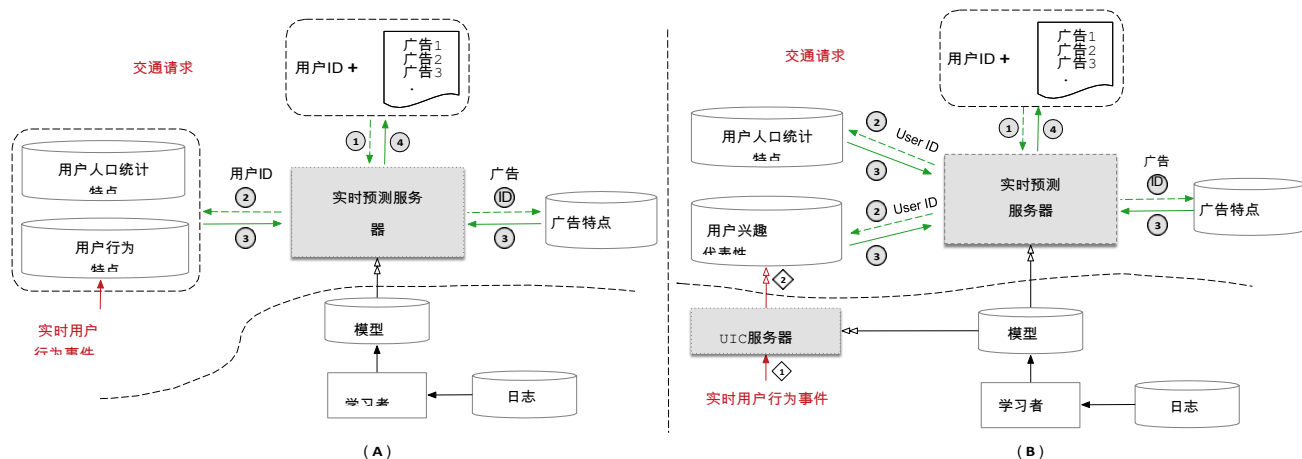


图2：用于CTR任务的实时预测（RTP）系统图解。通常，它由三个关键部分组成：特征管理模块、模型管理模块和预测服务器。(A)是我们的RTP系统的最后一个版本，(B)是带有提议的UIC服务器的更新版本。系统A和B的关键区别在于用户兴趣的计算：(i)在A中，它是在预测服务器内执行的，与请求有关。(ii)在B中，它是在UIC服务器中单独执行的。实时的用户行为事件。也就是说，它与流量请求脱钩，无延迟。

通过执行实时模型推断来预测概率分数。这个过程需要在严格的延迟限制下完成，实际中的典型值是10毫秒。

图2的A部分简要说明了我们的在线显示广告系统中CTR任务的实时**预测**（简称**RTP**）系统。为了便于读者理解，我们假设RTP的输入请求只包括用户和广告信息，而忽略了背景或其他因素。

3.1 用长序列的用户行为数据提供服务的挑战

在工业应用中，如电子商务行业的推荐系统[30, 31]，用户行为特征在特征集中贡献最大。例如，在我们的系统中，近90%的特征是用户行为特征，其余10%为用户人口统计特征和广告特征。这些行为数据包含丰富的信息，对用户兴趣建模很有价值[5, 17, 19, 25]。图1显示了我们的系统在不同日子里收集的用户行为序列的平均长度，以及用不同长度的用户行为特征训练的基本模型（Embedding&MLP[31]）的基本表现。在没有任何其他努力的情况下，与长度为100的序列相比，长度为1000的基本模型在AUC上有0.6%的提高。值得一提的是，0.3%的AUC改进对我们的业务来说是足够重要的。这一改进表明，利用这些长序列的用户行为数据是非常有价值的。

然而，利用长序列的行为数据带来了巨大的挑战。实际上，来自数以亿计的用户的行为特征是一个巨大的容量。为了保持推荐系统的低延迟和高吞吐量，行为特征通常被存储在一个额外的分布式内存存储系统中，例如我们系统中的TAIR[9]。这些特征会被提取到预测服务器上，并在系统中进行分析。

参与实时推理的计算, 当收到

流量请求。根据我们的实践经验，在我们的系统中实现DIEN[30]要花费大量的工程工作。在用户行为序列长度为150的情况下，延迟和吞吐量都已经达到了RTP服务器的性能边缘，更不用说长度为1000的情况了。直接涉及更多的用户行为数据是相当困难的，因为它面临着几个挑战，其中最关键的两个包括：

- **存储限制。**在我们的系统中，有超过6亿的用户。每个用户的行为序列的最大长度为150。这需要花费大约1TB的存储空间，不仅要存储产品ID，还要存储其他相关的特征ID，如商店ID、品牌ID等。当行为序列的长度达到1000时，将消耗6TB的存储空间，并且这个数字随着用户行为序列的长度线性增加。正如我们之前提到的，在我们的系统中使用高性能的存储，以保持低延迟和高吞吐量，而持有如此巨大的存储是非常昂贵的。巨大的存储量也会导致用户行为特征的对应计算和更新的成本足够高。因此，一个相当长的行为序列意味着一个不可接受的存储消耗。
- **延迟限制。**众所周知，使用连续的深度网络进行实时推理是非常具有挑战性的，特别是在我们有大量请求的情况下。DIEN[30]采用了一些技术，在每个工作者的QPS（每秒查询次数）为500的情况下，将我们系统中DIEN服务的延迟降低到14ms。然而，当用户行为的长度达到1000时，DIEN的延迟在500QPS时达到200ms。在我们的显示广告系统中，500QPS的延迟限制为30ms，这是很难承受的。因此，在目前的系统架构下，不可能获得长行为的好处。

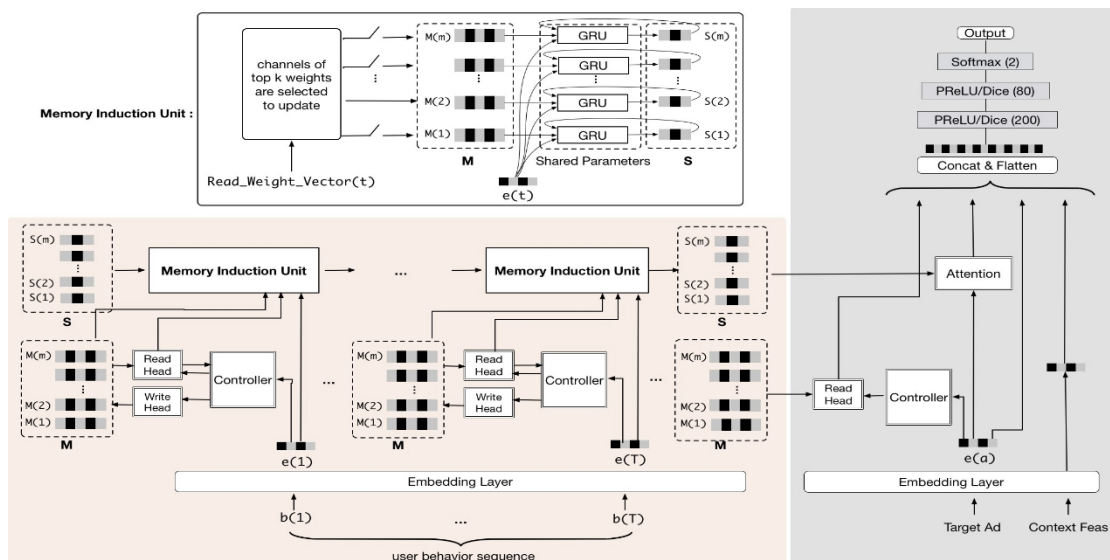


图3：提议的MIMN模型的网络结构。MIMN由两个主要部分组成：(i)左边的子网络，主要是利用顺序行为特征进行用户兴趣建模；(ii)右边的子网络遵循传统的Embedding&MLP范式，将左边子网络的输出和其他特征作为输入。MIMN的贡献在于左侧子网络，它是由NTM模型激发的，包含两个重要的存储架构：

a) 基本的NTM存储器单元，具有存储器读和存储器写的标准操作；b) 具有多通道的存储器感应单元。通道的GRU根据前面学到的NTM记忆来捕捉高阶信息。

3.2 用户兴趣中心

为了解决上述长序列用户行为建模的挑战，我们提出了一个机器学习算法和服务系统共同设计的解决方案。由于用户行为建模是CTR预测系统中最具挑战性的部分，我们设计了一个UIC（用户兴趣中心）模块来处理它。

图2的B部分说明了新设计的带有UIC服务器的RTP系统。系统A和B的区别在于用户兴趣表征的计算。在B中，UIC服务器维护每个用户的最新兴趣表示。UIC的一个关键点是它的更新机制。用户智能状态的更新只取决于实时的用户行为触发事件，而不是重新探索。也就是说，UIC对于实时CTR预测是没有延迟的。在我们的系统中，UIC可以将用户行为长度为1000的DIEN模型的延迟从200毫秒减少到19毫秒（500QPS）。

4 多通道用户兴趣记忆网络

在这一节中，我们将详细介绍我们的机器学习算法，用于长序列的用户行为建模。

4.1 从长序列的用户行为数据中学习 的挑战

众所周知，从长序列数据中学习是很困难的。简单的RNN（RNN[28], GRU[7], LSTM[14]）在面对相当长的序列时失败并不奇怪。引入注意力机制是为了

通过压缩所需的数据来提高模型的表达能力。

顺序数据的信息变成一个固定长度的张量[1, 31]。例如，在DIN[31]中，注意力机制通过软搜索与目标项目相关的部分隐藏状态或源行为序列来工作。为了进行实时推理，它需要存储所有的原始行为序列，这给在线系统带来了很大的存储压力。此外，注意力的计算成本随着行为序列的长度线性增长，这对于长序列的用户行为建模是不可接受的。实际上，RNNs中的隐藏状态并不是为了存储过去源序列的全部信息，而是为了更多的关注预测目标。因此，最后的隐藏状态可能会忘记长期信息。此外，存储所有的隐藏状态也是多余的。

最近，NTM（神经图灵机[11]）被提出来，从源序列中获取信息并将其存储在一个固定大小的外部存储器中，在许多长序列数据的建模任务中实现了比RNNs模型的显著改进。借用NTM的思想，本文提出了一个基于内存网络的模型，它为我们提供了一个处理长序列用户行为建模的新方案。我们把这个模型命名为MIMN（多通道用户兴趣记忆网络），如图3所示。UIC存储MIMN的外部记忆张量，并为用户的每一个新行为更新它。通过这种方式，UIC从用户的行为序列中逐步捕捉用户的兴趣。尽管UIC存储了一个固定长度的内存张量，而不是原始的行为序列，但考虑到存储压力，内存张量的尺寸必须受到限制。在本文中，我们提出了内存利用正则化，通过增加内存张量的维度来提高UIC的表达能力。

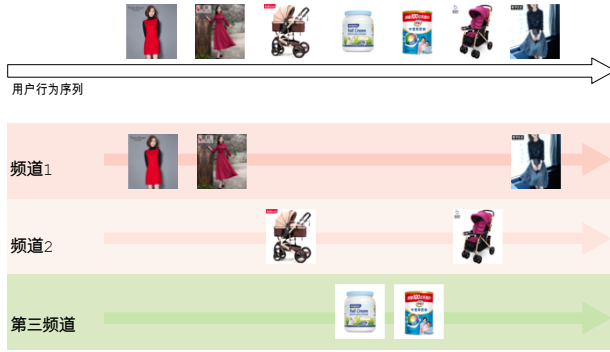


图4：多通道记忆的诱导过程。

内存的利用率。另一方面，由于用户的兴趣不同

以及随着时间的推移而演变，我们提出记忆诱导单元来帮助捕捉高阶信息。

4.2 神经图灵机

MIMN遵循传统的Embedding&MLP范式[8, 30, 31]，其中更多细节我们请读者参考[31]。MIMN的结构如图3所示。

标准的NTM通过一个存储网络来捕捉和存储顺序数据的信息。在 t 的时间步长中，存储器的参数表示为 M_t ，它由 m 个存储槽组成 $\{M_t(i)\}_{i=1}^m$ 。NTM的两个基本操作是内存读取和

内存写入，通过控制器与内存互动。

内存读取。输入有第 t 个行为嵌入向量，即控制器产生一个读键 k_t 来寻址内存。它首先跟踪与所有存储槽相比，产生一个权重向量 $w_t^r \propto \exp(K_{k_t, M_t}(i))$

$$w_t^r(i) = \frac{\exp(K_{k_t, M_t}(i))}{\sum_{j=1}^m \exp(K_{k_t, M_t}(j))}, \text{ for } i = 1, 2, \dots, m \quad (1)$$

其中

$$K_{k_t, M_t}(i) = \frac{t}{\|k_t\| \|M_t(i)\|}, \quad (2)$$

然后计算出一个加权的记忆总结作为输出 r_t ，

$$r_t = \sum_{i=1}^m w_t^r(i) M_t(i). \quad (3)$$

内存写入。内存写入时的权重向量 w_t^w 尽可能多地将源数据写入内存。

数据的生成类似于内存读取，对公式(1)进行运算。

控制器还产生了两个额外的键，即添加矢量 a_t 和擦除矢量 e_t ，它们控制存储器的更新。

其中 E_t 是擦除矩阵， A_t 是添加矩阵， $E_t = w_t^w \otimes e_t$ 和 $A_t = w_t^w \otimes a_t$ ， \otimes 为指点积和外产品分别。

4.3 内存利用率正规化

实际上，基本的NTM存在内存利用不平衡的问题，特别是在用户兴趣建模的情况下。也就是说，热门项目往往很容易出现在用户的序列中。

行为数据，并主导着内存的更新，使得使用

的内存是低效的。之前在NLP领域的工作[22, 23]提出使用LRU策略来平衡每个内存的利用率。因为LRU在处理过程中非常注意平衡每一个短暂的序列中的内存利用率，所以LRU几乎不会将信息写入相邻时间段的同一个槽中。然而，在我们的场景中，用户可能会与几个属于同一兴趣的行为进行交互，因此这些行为应该被写入同一槽位。LRU会扰乱内容寻址，不适合我们的任务。在本文中，我们提出了一个新的策略，名为内存利用率正则化，它被证明对用户兴趣建模是有效的。

内存利用率正规化。内存利用率正则化策略背后的想法是对内存利用率的方差进行正则化。在不同的内存插槽中写入重量，推动了内存

利用，以达到平衡。让 $g_t = \sum_{c=1}^C w_t^w$ 是累积的上等的权重，直到第 t 个时间步骤，其中 w_t^w 代表重新平衡的

在第 c 个时间步长中的写入权重。重新平衡的写入权重 w_t^w 可以表述为：

$$p_t = \text{softmax}(w \partial g_t) \quad (5)$$

$$w_t^w = \frac{w_t^w}{p_t}. \quad (6)$$

w_t^w 是4.2小节中介绍的原始写权重， w_t^w 代表内存更新时的新写权重。权重转移矩阵 p_t 取决于(i) g_t ，代表第 t 步每个存储槽的累积利用率；(ii) 参数矩阵 W_∂ ，通过正则化损失学习：

$$w_t^w = \frac{1}{t} \sum_{i=1}^t w_i^w, \quad (7)$$

$$L_{reg} = \lambda \sum_{i=1}^m \left(\frac{1}{m} \sum_{t=1}^T w_t^w(i) - \frac{1}{m} \sum_{t=1}^T w_t^w(i)^2 \right), \quad (8)$$

其中 m 是内存槽的数量。 L_{reg} 有助于减少不同内存插槽间更新权重的差异 w_t^w ，所有 m 个槽的更新率趋于均匀。在这种情况下

这样，所有的内存插槽的利用率都得到了提高，达到了平衡。利用率正规化可以帮助内存张量实现

从源行为数据中存储更多信息。

4.4 记忆诱导单元

NTM中的存储器被设计用来存储来自于以下方面的原始信息尽可能多的源数据。这里面的一个问题是，它可能会

错过了对一些高阶信息的捕捉，如不断变化的

兴趣的每一部分的过程。为了进一步提高用户兴趣提取的能力，MIMN设计了一个记忆诱导单元(MIU)。MIU也包含一个内部存储器，槽的数量为 S ，与NTM相同。这里我们指的是每个记忆槽作为一个兴趣通道。在第 t 个时间步，MIU：

(i)选择 k 个通道，通道索引在集合 $\{i : w_t^w(i) \in r\}$ 。 w_t^w 是内存读数的权重向量，即 $\text{top}_k(w_t^w)$ 。

前述的NTM，如公式(1)所示。(ii)对于第 j 个选择的通道，根据公式(9)更新 $S_t(i)$ 。

$$S_t(i) = \text{GRU}(S_{t-1}(i), M_t(i), e_t), \quad (9)$$

其中 $M_t(i)$ 是NTM的第 j 个内存插槽， e_t 是行为嵌入向量。公式(9)显示，MIU捕获的信息来自于

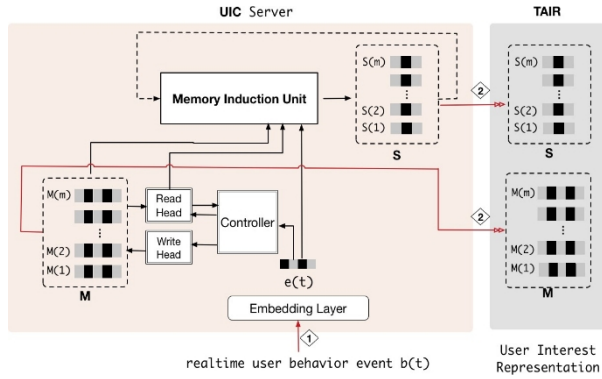


图5：在UIC服务器中用NTM和MIU实现用户间模型的子网络。

在NTM模块中记忆的原始行为输入和信息。这就像一个归纳过程，如图4所示。多通道记忆的GRU的参数是共享的，没有增加参数量。

4.5 在线服务的实施

与[30, 31]应用注意力机制来获得以候选人为中心的兴趣表示不同，MIMN学习捕捉并明确地将用户的不同兴趣存储在每个用户的外部存储器中。这种基于内存的架构不需要在候选者（例如我们系统中的目标广告，如图3所示）和用户行为序列之间进行交互式计算，并且可以增量执行，使其可以扩展到长序列的用户行为建模。

在线服务的MIMN的实现是直接的。正如第3.2节所介绍的，我们将整个模型拆分并在两个服务器中实现：如图5所示，在UIC服务器中实现用于用户间最繁重计算的左侧子网络，其余右侧子网络在RTP服务器中实现。图3清楚地说明了这种实现。

NTM和MIU模块都享有增量计算的好处。最新的内存状态代表了用户兴趣，并被更新到TAIR，用于实时CTR预测。当收到一个新的用户行为事件时，UIC计算并再次将用户兴趣代表更新到TAIR。通过这种方式，用户行为数据不需要被存储。在我们的系统中，巨大的长期用户行为量可以从6T减少到2.7T。

讨论。UIC服务器和MIMN算法的共同设计使我们能够处理长度可达数千的连续的用户行为数据。UIC的更新是独立于整个模型计算的，这使得它在实时CTR预测中没有延迟。MIMN建议以增量方式建立用户兴趣模型，而不需要像传统解决方案那样存储整个用户行为序列。此外，MIMN的设计采用了改进的内存架构，使模型性能更优越。然而，它并不适合所有的情况。我们建议将这个解决方案应用于

的应用：(i)丰富的用户行为数据，(ii)流量规模

表1：本文使用的数据集的统计数据。

| 数据集 | 用户 | 项目 ^a | 类别 | 实例 |
|---------|--------|-----------------|---------|--------|
| 亚马逊（图书） | 75053 | 358367 | 1583 | 150016 |
| 淘宝网。 | 987994 | 4162024 | 9439 | 987994 |
| 工业。 | 0.29亿 | 0.6亿 | 100,000 | 12.2亿 |

^a 对于工业数据集，项目指的是广告。

实时用户行为事件的数量不能明显超过实时CTR预测的请求。

5 实验

在本节中，实验分为两部分：（一）我们详细介绍了算法验证，包括数据集、实验设置、比较模型和相应的分析。公共数据集和实验代码都是可用的¹。

(ii) 我们讨论并分享了我们在阿里巴巴的展示广告系统中部署所提出的解决方案的实际经验。

5.1 数据集和实验设置

在两个公共数据集和一个从阿里巴巴的在线展示广告系统中收集的工业数据集上进行了模型比较。表1显示了所有数据集的统计数据。

亚马逊数据集²是由产品评论和元数据组成。

来自亚马逊的数据[20]。我们使用亚马逊数据集的书籍子集。对于这个数据集，我们将评论视为一种交互行为，并按时间对一个用户的评论进行排序。假设用户 u 有 T 个行为，我们的目的是利用 $T-1$ 个行为的预览来预测用户 u 是否会写 T 个评论中的评论。为了专注于长序列的用户行为预测，我们过滤这些行为序列长度短于20的样本，并在长度为100时截断行为序列。

淘宝数据集³是一个来自淘宝网的用户行为集合。

推荐系统[12]。该数据集包含几种类型的用户行为，包括点击、购买等。它包含了大约一百万用户的用户行为序列。我们将每个用户的点击行为按照时间排序，以构建行为序列。假设用户 u 有 T 个行为，我们使用前 $T-1$ 个点击的产品作为特征来预测用户是否会点击第 T 个产品。行为序列的长度被截断为200。

工业数据集收集自阿里巴巴的在线展示广告系统。样本来自印象记录，以“点击”或“不点击”作为标签。训练集由过去49天的样本组成，测试集由第二天的样本组成，是工业建模的经典设置。在这个数据集中，每天的样本中的用户行为特征包含前60天的历史行为序列，长度被截断为1000。

实验设置对于所有模型，我们使用Adam[18]求解器。我们采用指数衰减法，学习率从0.001开始，衰减率为0.9。FCN（全连接网络）的层数设定为200×80×2。嵌入维数为

¹ <https://github.com/UIC-Paper/MIMN>

² <http://jmcauley.ucsd.edu/data/amazon/>

³ <https://tianchi.aliyun.com/dataset/dataDetail?dataId=649&userId=1>

表2：公共数据集上的模型性能（AUC）。

| 模型 | 淘宝（平均值±std） | 亚马逊（平均值±std） |
|---------|------------------|------------------------|
| 嵌入&MLP | 0.8709 ± 0.00184 | 0.7367 ± 0.00043 |
| DIN | 0.8833 ± 0.00220 | 0.7419 ± 0.00049 |
| GRU4REC | 0.9006 ± 0.00094 | 0.7411 ± 0.00185 |
| ARNN | 0.9066 ± 0.00420 | 0.7420 ± 0.00029 |
| RUM | 0.9018 ± 0.00253 | 0.7428 ± 0.00041 |
| DIEN | 0.9081 ± 0.00221 | 0.7481 ± 0.00102 |
| MIMN | 0.9179 ± 0.00150 | 0.03250.7593 ± 0.00150 |

设置为16，这与内存插槽的尺寸相同。MIU中GRU的隐藏维数被设定为32。在NTM和MIU中，内存槽的数量是一个参数，在消融研究部分会仔细检查。我们把AUC作为衡量模型性能的指标。

5.2 竞争者

我们将MIMN与最先进的CTR预测模型在长序列用户行为建模的情况下进行比较。

- **嵌入&MLP**是CTR预测的基本深度学习模型。它采用和池操作来整合行为嵌入。
- **DIN**[31]是一项早期的用户行为建模工作，它提出了对候选人进行软搜索的用户行为。
- **GRU4Rec**[13]以RNN为基础，是第一个使用递归单元来模拟连续用户行为的工作。
- **ARNN**是GRU4Rec的一个变种，它使用注意力机制来对所有隐藏状态进行加权求和，以更好地表示用户序列。
- **RUM**[3]使用一个外部存储器来存储用户的行为特征。它还利用软写和注意力阅读机制与存储器进行交互。我们使用特征级的RUM来存储序列信息。
- **DIEN**[30]整合了GRU和以候选人为中心的注意力技巧，以捕捉用户兴趣的演变趋势，并取得了最先进的性能。为了便于比较，我们省略了DIEN中用于更好地嵌入学习的辅助损失技巧，否则所有上述模型都应该实现这一技巧。

5.3 公共数据集的结果

表2列出了所有比较模型和MIMN的结果。每个实验重复3次。

所有其他模型都击败了Embedding&MLP，这验证了网络架构设计对用户行为建模的有效性。MIMN击败了所有的模型，在AUC指标上有明显的提升。我们相信这是因为基于内存的架构容量巨大，适合于用户行为建模。正如在第4.1节中所讨论的那样，用户兴趣的背后是漫长的时间。

顺序行为数据是多样化的，并随着时间的推移而演变。MIMN通过多通道的记忆来学习捕捉用户的兴趣，在两个方面：(i) 在基本的NTM中的内存，平衡利用了

兴趣记忆，(ii) MIU中的记忆，进一步捕捉到了

表3：模型在不同槽位数上的表现

| 模型 | 淘宝（平均值±std） | 亚马逊（平均值±std） |
|---------|------------------|------------------|
| MIMN 4槽 | 0.9046 ± 0.00135 | 0.7522 ± 0.00231 |
| MIMN 6槽 | 0.9052 ± 0.00202 | 0.7503 ± 0.00120 |
| MIMN 8槽 | 0.9070 ± 0.00186 | 0.7486 ± 0.00071 |

通过诱导基于NTM记忆的利息顺序关系来获得高阶信息。

5.4 消融研究

在本节中，我们研究MIMN中不同模块的影响。**内存插槽的数量**。我们在MIMN上进行了不同数量的内存插槽的实验，这是一个手动设置。为了简化，我们只用基本的NTM架构来评估MIMN，省略了**内存利用率正则化**的设计。

和**记忆诱导单元**。表3显示了结果。从经验上看，插槽数会影响模型的性能。对于

我们的分析结果表明，这与数据集中用户行为序列的长度有关。每个内存槽都是随机初始化的。对于具有较长行为序列的数据集，例如淘宝数据集，内存有更多的机会学习并实现稳定的表示。在行为序列较短的情况下，例如亚马逊数据集，内存容量较大的模型性能会受到学习的影响。特别是当所有内存插槽的利用率不平衡时，部分内存向量可能没有得到充分的利用和更新，这意味着这些内存向量仍然保持在原始初始化附近。这将损害模型的性能。因此，我们提出了内存利用率正则化来缓解这个问题。

存储器利用正则化。由于每个用户的inter-est强度不均匀，以及内存的随机初始化，在基本的NTM模型中，存储的利用可能是不平衡的。这个问题会影响内存的学习，使其不能充分地利用有限的内存存储。我们采用了**内存利用规律化**的技巧来帮助解决这个问题。图6显示了内存利用率，这验证了所提出的正则化器的有效性。如表4所示，这种平衡效应也带来了模型性能的改善。

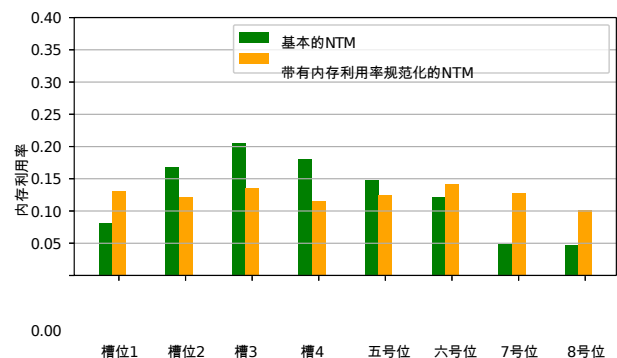


图6：NTM中不同槽位的内存利用率

表4：有/无内存利用正则化和内存诱导单元的MIMN模型性能（AUC）比较

| 模型 | 淘宝网（平均值±标准差） | 亚马逊（平均值±标准差） |
|-----------------------------|------------------------------|---------------------------------------|
| 不含MUR和MIU的MIMN ^a | 0.9070±0.001860.7486±0.00071 | 含MUR的MIMN0.9112±0.002670.7551±0.00121 |
| 含MUR和MIU的MIMN ^b | 0.9179±0.00208 | 0.7593±0.00296 |

^a MUR是指内存利用正则化。^b MIU是指记忆诱导单元。**表5：工业数据集的模型性能（AUC）**

| 模型 | AUC |
|--------------------|--------|
| DEIN | 0.6541 |
| 医学博士 | 0.6644 |
| 在不同步设置下的MIMN（一天内）。 | 0.6644 |
| 用大宗销售数据训练的MIMN | 0.6627 |

记忆诱导单元。通过从基本的NTM诱导记忆，带有记忆诱导单元的MIMN能够捕获高阶信息并带来更多的改进，如表4所示。它增强了用户兴趣提取的能力，有助于从长序列行为数据中建立用户兴趣模型。

5.5 工业数据集的结果

我们进一步对从阿里巴巴的线上广告系统收集的数据集进行了实验。我们将MIMN与DIEN的模型进行了比较。表5显示了结果。MIMN改进了DIEN，AUC增益为0.01，这对我们的业务是非常重要的。除了离线模型的性能，MIMN和DIEN模型在系统问题方面也存在很大差异。图7显示了真实的CTR预测系统在使用MIMN和DIEN的模型时的系统性能。与UIC服务器共同设计的MIMN在很大程度上击败了DIEN，它拥有恒定的延迟和吞吐量的特性。因此，MIMN已经准备好在我们的系统中利用长度可扩展到数千的长序列用户行为数据，并享受模型性能的改善。相反，使用DIEN的系统在延时和系统吞吐量方面都受到影响。由于系统的压力，DIEN作为我们最后的产品模型所利用的用户行为序列的长度只有50。这就验证了再次证明我们提出的解决方案的优越性。

在线A/B测试。我们已经在阿里巴巴的展示广告系统中部署了所提出的解决方案。从2019-03-30到2019-05-10，我们进行了严格的在线A/B测试实验来验证提出的MIMN模型。与DIEN（我们上一个产品模型）相比，MIMN实现了7.5%的点击率和6%的RPM（每米收入）收益。我们把这归功于从长的连续行为数据中挖掘出额外的内部形成，而拟议的协同设计解决方案使。

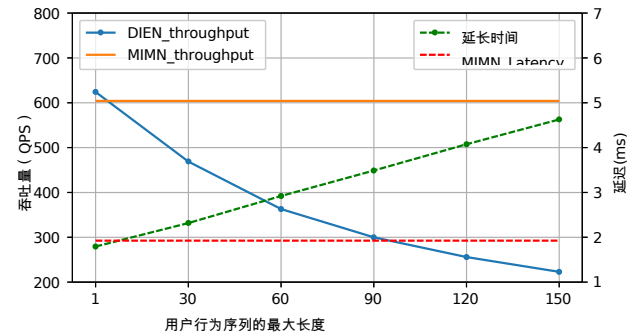


图7：用MIMN和DIEN模型服务的实时CTR预测系统在不同长度的用户行为序列中的系统性能。MIMN模型是通过UIC服务器的设计实现的。

5.6 部署的实践经验

在本节中，我们将讨论在我们的在线系统中部署UIC和MIMN的拟议解决方案的实际经验。

UIC服务器和RTP服务器的同步。如4.5节所述，MIMN是与UIC和RTP服务器一起实现的。因此，UIC和RTP服务器之间存在一个不同步的问题。两台服务器的参数更新不同步，可能会导致实际系统中的模型推断不正确，这是很危险的。我们进行了实验来模拟失同步的情况。表5显示了结果。MIMN在参数更新不同步的情况下，对模型性能的影响不大。注意在这个实验中，同步外更新时间的差距在一天之内，这是工业系统中的一个传统设置。实际上在我们的实际系统中，模型部署被设计为每小时执行一次，进一步降低了风险。我们相信这是由于MIMN对用户兴趣的稳定表示，导致MIMN有良好的泛化性能。

大甩卖数据的影响。如今，许多电子商务网站都采用大促销的方式来吸引顾客进行在线消费，例如中国著名的阿里巴巴11.11大促销。在这种极端的情况下，样本的分布以及用户的行为都与日常的情况有很大的不同。我们比较了MIMN在我们的系统中使用和不使用11.11大促销日收集的训练数据的性能。结果显示在表5中。我们发现，根据经验去除大甩卖的数据会更好。**热身策略。**尽管UIC的设计是为了增加更新，但从一开始就需要相当长的时间来进行稳定的积累。实际上，我们采用热身策略，用预先计算好的用户兴趣表征对UIC进行初始化。也就是说，我们收集每个用户过去120天的历史行为（用户行为序列的平均长度为1000），并在离线模式下用训练好的MIMN模型进行推理，然后将积累的记忆推送到UIC进行进一步的更新。这一策略使模型的性能得到了合理的提升。

当尽快部署拟议的解决方案。

回滚策略。在出现意外问题的情况下，例如大规模在线作弊对训练样本的破坏，在

UIC服务器的增量更新机制可能会受到很大的影响。A

麻烦的挑战是寻找异常情况的发生点。为了抵御这种风险,我们设计了一个回滚策略,在每天凌晨00:00存储所学到的用户兴趣表示的副本,并记录过去7天的情况。

6 结论

在本文中,我们专注于通过机器学习算法和在线服务系统的共同设计来利用长序列的用户行为数据。在计算方面,存储是从相当长的连续用户行为数据中捕获长期用户兴趣的主要瓶颈。我们介绍了我们对新的解决方案的实践—用于实时推断用户兴趣模型的解耦UIC服务器和基于内存的MIMN模型,该模型可以逐步实现,并优于其他最先进的模型。

值得一提的是,深度学习给我们带来了一个强大的工具包,可以为工业应用提供更多有价值的数据。我们相信这项工作通过对极长的连续用户行为数据进行建模开辟了新的空间。在未来,我们计划进一步推进研究,包括学习算法、训练系统以及在线服务系统。

鸣谢

作者要感谢马国强、沈振中、刘浩宇、王维钊、马驰、宋俊涛、易鹏涛,他们为在线系统的实施做了大量的工作。

参考文献

- [1] Dzmity Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. 通过联合学习对准和翻译的神经网络翻译。 *arXiv预印本 arXiv:1409.0473* (2014)。
- [2] Dzmity Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. 通过联合学习对齐和翻译的神经机器翻译。在 *第三届“学习表征”国际会议的论文集中*。
- [3] 陈旭, 徐宏腾, 张永峰, 唐家喜, 曹一心, 秦正, 和查宏远. 2018. 带有用户记忆网络的顺序推荐。在 *第十一届ACM网络搜索和数据挖掘国际会议论文集*。ACM, 108-116。
- [4] Heng-Tze Cheng 和 Levent Koc. 2016. 用于推荐系统的广泛和深度学习。In *Proceedings of the 1st Workshop on Deep Learning for Recommended Systems*. ACM, 7-10。
- [5] Christina Yip Chung, Abhinav Gupta, Joshua M Koran, Long-Ji Lin, and Hongfeng Yin. 2011. 行为定位系统中长期和短期用户资料分数的递增更新。美国专利 7,904,448。
- [6] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2014. 门控递归神经网络在序列建模上的经验评估。 *arXiv预印本 arXiv:1412.3555* (2014)。
- [7] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2014. 门控递归神经网络在序列建模上的经验评估。 *arXiv预印本 arXiv:1412.3555* (2014)。
- [8] Paul Covington, Jay Adams, and Emre Sargin. 2016. 用于YouTube推荐的深度神经网络。在 *第十届ACM会议上, Recommender Systems的论文集*。ACM, 191-198。
- [9] duolong, daoan, and fanggang. 2017. TAIR, 一个由阿里巴巴集团开发的分布式键值存储系统。 <https://github.com/alibaba/tair>。
- [10] Travis Ebesu, Bin Shen, and Yi Fang. 2018. 用于推荐系统的协作记忆网络。 *arXiv预印本 arXiv:1804.10862* (2018)。
- [11] Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. 神经图灵机。 *arXiv预印本 arXiv:1410.5401* (2014)。
- [12] Zhu Han, Pengye Zhang, Guozheng Li, He Jie, and Kun Gai. 2018. 学习基于树的深度模型用于推荐系统。(2018), 1079-1088。
- [13] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. *arXiv preprint arXiv:1511.06939* (2015). 基于会话的推荐与递归神经网络。
- [14] Sepp Hochreiter and Jürgen Schmidhuber. 1997. 长短期记忆。 *神经计算* 9, 8 (1997), 1735-1780。
- [15] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. [n.d.]. 密集连接的卷积网络。
- [16] Jin Huang, Wayne Xin Zhao, Hongjian Dou, Ji-Rong Wen, and Edward Y Chang. 2018. 用知识增强的记忆网络改进顺序性推荐。In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. ACM, 505-514。
- [17] Hyoungh R Kim and Philip K Chan. 2003. 学习隐含的用户兴趣层次, 用于个性化的背景。In *Proceedings of the 8th international conference on Intelligent user interfaces*. ACM, 101-108。
- [18] Diederik P Kingma 和 Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)。
- [19] 刘洪哲和MS Zamanian. 2007. 基于短期和长期相结合的 用户行为兴趣, 在网络上选择和提供广告的框架。美国专利申请。11/225,238。
- [20] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. 基于图像的风格和替代物推荐。在 *论文集 of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 43-52。
- [21] Massimo Quadrana, Alexandros Karatzoglou, Balázs Hidasi, and Paolo Cremonesi. 2017. 用分层递归神经网络实现基于会话的个性化推荐。在 *第11届ACM推荐人会议论文集系统*。ACM, 130-137。
- [22] Jack Rae, Jonathan J Hunt, Ivo Danihelka, Timothy Harley, Andrew W Senior, Gregory Wayne, Alex Graves, and Timothy Lillicrap. 2016. 用稀疏的读和写扩展记忆增强的神经网络。In *Advances in Neural Information Processing Systems*. 3621-3629。
- [23] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. *arXiv preprint arXiv:1605.06065* (2016). 用记忆增强的神经网络进行一次性学习。
- [24] Ying Shan, T Ryan Hoens, Jian Jiao, Haijing Wang, Dong Yu, and JC Mao. [n.d.]. 深度交叉: 网络规模的建模, 无需手工制作组合特征。
- [25] Yang Song, Ali Mamdouh Elkahky, and Xiaodong He. 2016. 用于时间性推荐的多速率深度学习。在 *第39届ACM SIGIR国际信息检索研究与发展会议*上。ACM, 909-912。
- [26] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. 2015. 端对端记忆网络。在 *神经信息处理系统的进展中*。2440-2448。
- [27] 王庆勇, 尹鸿志, 胡志亭, 连德福, 王浩, 和金子. 2018. 具有对抗性训练的神经记忆流推荐者网络。In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2467-2475。
- [28] 罗纳德-J-威廉姆斯和大卫-齐普赛尔. 1989. 一种用于持续, 运行完全递归神经网络的学习算法。 *神经计算* (1989), 270-280。
- [29] 翟双飞, 张景浩, 张若飞, 和张中飞. 2016. Deepintent: 用递归神经网络学习在线广告的注意力。In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. acm, 1295-1304。
- [30] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*. Honolulu, USA。
- [31] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yangghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. 用于点击率预测的深度兴趣网络。在 *第24届ACM SIGKDD知识发现与数据挖掘国际会议论文集*。ACM, 1059-1068。