# Continuous Personalized Knowledge Tracing: Modeling Long-Term Learning in Online Environment

Chunpai Wang, Sherry Sahebi

University at Albany - SUNY





State University of New York



## Problem



Continuous Personalized Sequence Modeling

# Problem



- Continuous Personalized Sequence Modeling
- Online Education Systems
  - abundant of practice problems with large topic varieties
  - students can have a long continuous learning experiences
- Knowledge Tracing (KT)
  - model students' knowledge level over time , to
    - predict student performance
    - create a study plan
    - recommend learning materials, etc.



# **Knowledge Tracing Approaches**



- Traditional KT approaches
  - probabilistic and hidden Markov models (e.g. BKT [Corbett et al.])
  - logistic regression models (e.g., IRT [Frederic et al.], PFA [Philip et al.])
  - work well on small datasets
- Deep KT approaches
  - work well on large datasets (e.g., DKT [Piech et al.], DKVMN [Zhang et al], SAKT [Pandey et al.])

# **Knowledge Tracing Approaches**



- Traditional KT approaches
  - probabilistic and hidden Markov models (e.g. BKT [Corbett et al.])
  - logistic regression models (e.g., IRT [Frederic et al.], PFA [Philip et al.])
  - work well on small datasets
- Deep KT approaches
  - work well on large datasets (e.g., DKT [Piech et al.], DKVMN [Zhang et al], SAKT [Pandey et al.])
  - but, aim to learn **global** patterns as opposed to **personalized** ones
    - learn a shared set of parameters for all students
  - fail to represent student knowledge in the long run
    - truncate long learning sequences, and train on shorter batches
    - lose dependence between student sequence batches
    - difficult on personalization

# **Our Proposed Solution: CPKT**



- Continuous Personalized Knowledge Tracing (CPKT):
  - personalized
    - personalized memory slots
    - personalized knowledge acquisition and forgetting patterns
  - continuous
    - online learning and inference paradigm
    - without truncating learning trajectories
  - generalizable
    - Transition-Aware Stochastic Shared Embedding (TA-SSE)
    - avoid overfitting

#### **CPKT: Overall Architecture**





- Dynamic key-value memory networks with personalized parameters
  - global static latent concept features  $M^k$
  - personalized dynamic value matrix  $M_{s,t}^{v}$  (student mastery state)
  - student embeddings  $u_s$



- Personalized knowledge acquisition  $(a_s)$  and forgetting  $(f_s)$  rates to update student mastery over time
  - $\boldsymbol{f}_s = \sigma(w_3^T[\boldsymbol{v}_t, \boldsymbol{u}_s] + \boldsymbol{b}_3)$
  - $M_{s,t}^{v}(i) = M_{s,t-1}^{v}(i) \otimes [\mathbf{1}^{d_{h}} w_{q}(i)\mathbf{f}_{s}]$

• Personalized knowledge level  $(r_{s,q,t})$  and ability-knowledge summary  $(x_{s,q,t})$ 



# **CPKT: Continuous Knowledge Tracing**

• Update model continuously over time with incremental data feeding

- moving window to extract and train on the most recent historical records
- avoid catastrophic interference problem



**UNIVERSITY AT ALBANY** State University of New York

# **CPKT: Continuous Knowledge Tracing**

• Update model continuously over time with incremental data feeding

- moving window to extract and train on the most recent historical records
- avoid catastrophic interference problem
- Caveat:
  - many parameters
  - easily overfitted



**UNIVERSITY AT ALBANY** State University of New York

## **CPKT: Stochastic Shared Embedding to Rescue**

- Data-driven regularization of embedding layers
- Stochastically swapping embeddings during the training
- Which embeddings to swap?
  - Transition-Aware Stochastic Shared Embedding (TA-SSE) for problems
  - random for students

#### **CPKT: Transition-Aware Stochastic Shared Embedding**

- Adds generalizability and avoids overfitting
  - generate item transition matrix *T* according to learning trajectories :

$$T_{i,j} = prob(j|i) = \frac{|i \to j|}{i}$$

• stochastically swap problems in the learning trajectory based on T

## **Continuous Learning Algorithm Summary**

#### Generate item transition matrix T

For each test item index t do

For each problem response pair  $(q_i^s, a_i^s)$  in each student s's last t-H to t-1 responses With probability  $\rho$ 

replace  $\boldsymbol{q}_i^s$  's embedding and response with a random problem according to T

For each student s with embedding  $\boldsymbol{u}_s$ 

With probability  $\boldsymbol{\rho}$ 

replace  $u_s$  with random student embedding  $u_z$ 

Update the personalized parameters by forward and backward passes

Predict the target student response at t

Collect the target student response as new data

Update the transition matrix T

Increase t by 1

# **Three Research Questions**



- RQ1. How is the model performance compared with state of-the-art baselines?
- RQ2. How do different proposed components affect its prediction performance?
- RQ3. How does the model perform on the users with different lengths of learning trajectories?

# **RQ1: CPKT Outperforms Baselines**



16

UNIVERSITY AT ALBANY State University of New York

# **RQ2: Ablation Study**



• Both personalization and TA-SSE Help

Mathada	MORF	ASSIST2015	EdNet	Junyi	
ivietnoas	RMSE	AUC	AUC	AUC	
CPKT-W/O-Pers	0.1919±0.0092	0.7202±0.0040	0.5745±0.0036	0.8546±0.0078	
CPKT-W/O-SSE	0.1895±0.0067	0.7092±0.0049	0.5753±0.0058	0.8333±0.0051	
СРКТ	0.1752±0.0081	0.7274±0.0032	0.6558±0.0072	0.8802±0.0072	

## **RQ3: Better Performance for Longer Trajectories**

• Arrange all students based on their test sequence lengths (corresponding to trajectory lengths) into three groups with roughly equal sizes: short, medium, long.

Group	#Llcorc	Range of Test		Mean AUC	P-Value of CPKT vs.		
	#USETS	Length	DKT	DKVMN	СРКТ	DKT	DKVMN
Short	199	[0,110]	\$0.64	0.6427	0.6321	p=0.8421	p=0.4767
Medium	210	[110,380]	0.6422	0.6363	0.644	p=0.7375	p=0.1398
Long	234	[380,900]	0.6413	0.6315	0.6475	p=0.0745	p=1.51e-05

Group	#1.10.010	Range of Test		Mean BCE	P-Value of CPKT vs.		
	#USEIS	Length	DKT	DKVMN	СРКТ	DKT	DKVMN
Short	199	[0,110]	0.655	0.6884	0.653	p=0.8583	p=0.0154
Medium	210	[110,380]	0.6474	0.675	0.6375	p=0.0475	p=2.87e-09
Long	234	[380,900]	0.6562	0.6942	0.6389	p=4.71e-05	p=1.70e-19

# Conclusions



Continuous Personalized Knowledge Tracing (CPKT):

- ★ could track personalized student knowledge
- over long learning trajectories using an online learning and prediction paradigm
- using a transition-aware stochastic shared embedding regularization method that could resolve the overfitting issues



Source code of CPKT: <u>https://github.com/persai-lab/CIKM2023-CPKT</u> Contact: Chunpai Wang <u>chunpaiwang@gmail.com</u>

This paper is based upon work supported by the National Science Foundation under Grant No. 2047500.







State University of New York

# **Related Work**



- Traditional Knowledge Tracing
  - probabilistic models (e.g. BKT) and logistic models (e.g. IRT, AFM)
  - works well on small dataset
  - allow personalization (learn student specific parameters)
- Deep Knowledge Tracing (KT)
  - works well on large dataset (e.g. DKT, DKVMN, SAKT)
  - aim to learn global patterns (learn a shared set of parameters for all students)
    - truncate long learning sequences, and train on shorter batches
    - lose dependence between student sequence batches
    - difficult on personalization
    - fail to represent student knowledge in the long term

# **Motivation**



#### **Online Education System**

- abundant of practice problems with large topic varieties
- student can have a "lifelong" or continuous learning experience
  - can practice same topics many times
  - can learn a variety of topics over long periods of time
  - demand of personalized long-term learning experience





# СРКТ



- We propose Continuous Personalized Knowledge Tracing (CPKT):
  - personalized KT
    - track individualized student knowledge and predict student performance
    - personalized memory slots to maintain learner's knowledge in a lifelong manner
    - personalized knowledge acquisition and forgetting patterns
  - continuous KT by online learning and inference paradigm
    - mimic the real world long-term continuous learning scenario
    - without truncating learning trajectories to train the model
  - incorporation of personalization and continuous model learning
    - transition-aware stochastic shared embedding
    - avoid overfitting

#### • Personalized knowledge tracing:

- given  $\{(q_1^s, a_1^s), (q_2^s, a_2^s), \dots, (q_{t-1}^s, a_{t-1}^s), q_t^s\}$ , predict  $a_t^s$ .
- the general deep learning-based models omit the superscript s in the context and do not differentiate distinct students' historical records.



- Personalized Memory Slots
  - maintain a dynamic value matrix  $\mathbf{M}_{s,t}^{v}$  for each student



- Personalized knowledge acquisition  $(a_{s,t})$  and forgetting  $(f_{s,t})$  rates to update student mastery over time
  - $\boldsymbol{f}_{s,t} = \sigma(w_3^T[v_t, u_s] + b_3)$
  - $M_{s,t}^{v}(i) = M_{s,t-1}^{v}(i) \otimes [\mathbf{1}^{d_{h}} w_{q}(i)\mathbf{f}_{s,t}]$

• Personalized knowledge level  $(r_{s,q,t})$  and ability-knowledge summary  $(x_{s,q,t})$ 



Personalized Knowledge Acquisition and Forgetting



## **CPKT: Stochastic Shared Embedding to Rescue**

- Stochastic Shared Embedding (SSE) [1]
  - data-driven regularization of embedding layers
  - stochastically swapping similar item embeddings during the training
  - implicitly adds exponentially many distinct reordering layers above the embedding layer and leads to exponentially many models trained at the same time
  - the loss landscape with SSE regularization becomes smoother and leads to better generalization.
  - requires an auxiliary knowledge graph to compute the switching probability distribution

[1] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James L Sharpnack. 2019. Stochastic Shared Embeddings: Data-driven Regularization of Embedding Layers. Advances in Neural Information Processing Systems 32 (2019).

#### **CPKT: Transition-Aware Stochastic Shared Embedding**

- Adds generalizability and avoids overfitting
  - propose to generate and use a transition matrix based on the student learning trajectories.
  - learning trajectories typically contain some information on the learning material similarities.
  - better generalization.
  - Q by Q transition matrix T:

$$T_{ij} = \operatorname{prob}(j \mid i) = \frac{|i \rightarrow j|}{|i|}$$

# Algorithm

#### Algorithm 1: CPKT

- **Input:** Observed student responses  $\Omega_{obs}$ , including other students' and the target student's historical responses. Hyperparameter  $\rho$  denotes by SSE threshold and H denotes by time window size.
- 1 Sort each student's responses by the timestamp.
- <sup>2</sup> Generate the learning transition matrix **T** based on all observed learning trajectories.
- <sup>3</sup> Generate a dictionary D that stores each problem as key and a list of corresponding observed score or correctness from all students as value.

4 <b>f</b>	or each testing time index t do
5	Extract each student's responses between time index
	$t - H$ and $t - 1$ , denotes by $\Omega_{t-H}^{t-1} = \{(q_i^s, a_i^s)\}_{i=t-H}^{i=t-1}$ .
6	Feed $\Omega_{t-H}^{t-1}$ along with student IDs into model.
7	<b>for</b> each problem $q_i^s$ and corresponding interaction
	$(q_i^s, a_i^s) \in \Omega_{t-H}^{t-1}$ do
8	Identify the problem embedding $\mathbf{k}_i$ and interaction
	embedding $\mathbf{v}_i$ .
9	Generate a random number $\gamma \in [0, 1]$ .
10	if $\gamma < \rho$ then
11	Replace $\mathbf{k}_i$ with $\mathbf{k}_j$ , where $j \sim T_{ij} = prob(j \mid i)$
12	Randomly sample a response $a_j$ for problem $q_j$
	from $\mathcal{D}$ .
13	Identify the interaction embedding $\mathbf{v}_j$ for
	$(q_j, a_j)$ , and replace $\mathbf{v}_i$ with $\mathbf{v}_j$ .
14	end
15	end

# Algorithm

16	for each students do								
16									
17	Identify the student embedding <b>u</b> <sub>s</sub> .								
18	Generate a random number $\gamma \in [0, 1]$ .								
19	if $\gamma < \rho$ then								
20	Randomly sample a student <i>z</i> from all students.								
21	Identify the student embedding $\mathbf{u}_{z}$ .								
22	Replace $\mathbf{u}_s$ with $\mathbf{u}_z$ .								
23	end								
24	end								
25	Forward and backward pass with the new embeddings								
	to train the model by minimizing the training loss.								
26	Predict the target student's response at time <i>t</i> .								
27	Collect the target student's new response into $\Omega_{obs}$ .								
28	Update the transition matrix <b>T</b> as well as $\mathcal{D}$ .								
29	Increase the testing time index by 1.								

30 end

• RMSE (H denotes the sliding window size)

$$\ell_{RMSE} = \sqrt{\frac{\sum_{s} \sum_{t}^{t+H} (a_t^s - p_t^s)^2}{n}}$$

• Binary Cross Entropy

$$\ell_{BCE} = -\sum_{s}\sum_{t}^{t+H} \left(a_t^s \log p_t^s + \left(1 - a_t^s\right) \log \left(1 - p_t^s\right)\right)$$

#### **Experimental Setup**



#### 5-Fold User Stratified Cross Validation

- Observe all historical records for all training users.
- Observe all records before a preset threshold time index for all testing users.
- Predict test users' performance after the threshold time index.
- The testing records will be uncovered over time, along with the rolling prediction.
- Set the threshold index at roughly 10% of the maximum sequence lengths for each data.

Descriptive Statistics of 4 Real World Datasets.									
Dataset	Users	Questions	Question Record	Mean Question Responses	STD Question Responses	Correct Question Responses	Incorrect Question Responses	Max Sequence Length	
MORF	686	10	12031	0.7763	0.2507	N/A	N/A	46	
ASSIST2015	19840	100	683801	N/A	N/A	500379	183,422	1000	
EdNet	1000	11249	200931	N/A	N/A	118767	82,184	1,000	
Junyi	1564	142	120984	N/A	N/A	86654	34328	1000	

#### **Experimental Setup**



#### Table 4: Hyperparameters of CPKT.

Dataset	$d_h$	$d_u$	Ν	H
MORF	32	2	11	20
ASSIST2015	128	8	6	150
EdNet	128	8	41	150
Junyi	128	4	7	200

Descriptive Statistics of 4 Real World Datasets.									
Dataset	Users	Questions	Question Record	Mean Question Responses	STD Question Responses	Correct Question Responses	Incorrect Question Responses	Max Sequence Length	
MORF	686	10	12031	0.7763	0.2507	N/A	N/A	46	
ASSIST2015	19840	100	683801	N/A	N/A	500379	183,422	1000	
EdNet	1000	11249	200931	N/A	N/A	118767	82,184	1,000	
Junyi	1564	142	120984	N/A	N/A	86654	34328	1000	