

Oversea Internship Final Report

SIDA GAO

Supervised by Prof. Emma Brunskill

Reinforcement Learning and Education Lab, CMU CSD

August 29, 2016

Abstract

This is the final report of my summer internship in CMU. It is *NOT* a formal academic paper, since all of the projects pictured in this report are still works in progress. Therefore, some parts of the write-up won't be very formal, especially the reference part, which would only appear as footnotes (and works that are not that relevant would not be listed at all). This report assumes the readers have a general understanding of traditional AI, machine learning, deep learning and decision making, but will provide some backgrounds in reinforcement learning (RL), and educational data mining (EDM).

The main project for this summer is Deep Batch RL, in which I work with a master student, aiming to compare the generalization ability of model-based and model-free RL methods which both adopt deep learning models. Our original hypothesis believes that model-free methods should outperform the model-based ones, however, our empirical results suggest the opposite. The side project, in which I take the lead, focuses on extracting prerequisite relationships in educational data. This project also adopts deep learning models, which hold an advantage of discovering complicated latent relations in the data over traditional methods. The last project I take part in is Adaptive Fraction, which aims to find robust tutoring policy with offline educational data collected from Fraction, an online tutor system. A robust policy is one we can expect it to be at least as good as our evaluations with simulated students. This project is currently at a rather preliminary phase and we are still working on designing and implementing the experiments.

1 Introduction

Artificial intelligents could be roughly divided into two classes: those that are trying to understand the world (classification, clustering, etc.) and those that are trying to improve the world with actions (decision making, planning, etc.). Reinforcement learning aims to do both, since it's trying to make optimal decisions in an unknown world. By doing some episodes of trials in the world, an RL agent collects trajectories of states (i.e. observations), rewards and actions. With these trajectories, the RL agent can pick the best action in any observed state, which would gain him the most reward. RL differs from decision making methods like Markov decision processes, because MDPs have a perfect understanding and observation of the world, while RL agents know nothing about the transition model of states or the reward model (i.e., the rewards related to each state). In real life situations, the actual states are often not observable, and the agents would also need to learn to infer the states with observations.

Education, or to be more specific, intelligent tutor, is a perfect case where RL methods can be applied. We can easily draw some parallels between tutor systems and RL agents. The actions are the various tutoring materials that could be presented to students, while the rewards could be students' performance in tests. The world is also unknown: the transition model of students' knowledge states

(i.e. how students would digest the materials) is unknown, and the reward model is also unknown (i.e. how student with a certain knowledge state would perform in a test). In education, we are also dealing with a partially observable environments, since we won't be able to directly see the knowledge state for students, and we can only have a belief state inferred from observations of the students.

Education also poses an even larger challenge than other RL tasks. In conventional RL tasks like robotics, the robot can run trials in the real environments as much as we want, since in most cases the worst we can expect is merely a broken mechanical arm resulted from tripping over a stone. In education, it won't be ethical to collect online data on real students with an agent that's still learning from trials and errors. A crappy work-in-progress would do irreversible harms to real kids in this practice, and the things that could be broken would be the kids' future. Therefore, doing offline policy evaluation with trajectories collected from real world tutor systems (in which runs policies either designed by expert humans or well-tested AI agents) is an important track of research in educational RL.

Prof. Emma Brunskill's group, in which I interned for the summer, lies its interest in RL and education. We have shown that these two areas are intrinsically related. Our group covers a wide spectrum of research areas, from RL theories (machine learning) and traditional AI, all the way to more application oriented tracks like data mining and human computer interaction (in education). Given the special, yet not unique challenge of education described previously, our research projects are mostly focused on offline policy evaluation or planning. Educational data mining, which focuses on data collected from intelligent tutoring systems, is one of the may side interests of our group. Several distinctive characteristics of EDM include dealing with sequential data, and a special focus on discovering the latent hierarchical structure of data items (i.e. prerequisite structure of knowledge components). Deep learning, while has been widely adopted in various areas like computer vision, is still an emerging methods in RL and EDM, where only only a handful of previous works could be found. It is also a new area for the group.

In this section, we have briefly described the backgrounds and motivations of the projects. The following sections are organised as: section 2 described a deep learning model which my work in all 3 projects are based on; section 3, section 4 and section 5 discusses the big pictures of the three projects, where we are and what part I play in them. At last, section 6 jots down some other aspects of the internship besides work, and summarizes the 10-week internship.

2 Deep Knowledge Tracing

2.1 Knowledge Tracing Going Deep

Knowledge tracing is a common task for educational data mining, with a goal towards accurately model the students' knowledge acquisition process. The conventional method, Bayesian Knowledge Tracing (BKT), keeps a binary state for each skill (i.e. knowledge component, or concept), either known or unknown. When a skill is not learned, by doing exercises on the skill, the student may have some probability to transition to the learned state. When answering a question on a certain skill, the student could get it wrong even when he already knows the skill (i.e. there exists a slip rate). And reversely, a student can get a question correct with luck, without knowing the skill, which corresponds to a guess rate. There have been various extentions to BKT, like modeling the forgetting process, and collapsing similar skills into one (since BKT tracks each skill separately, this kind of practice is trying to address the relations between skills).

There has been a quite influential work published in 2015, *Deep Knowledge Tracing*¹ (DKT), which adopts an LSTM to model students' knowledge acquisition. Retrospectively speaking, LSTM is basically an upgraded version of BKT. Instead of having one binary state representation for each skill, LSTM uses a long vector of real numbers (which is normally more than the number of skills) to represent the knowledge state. By doing this, LSTM could capture complicated (and unfortunately, not very interpretable) relationships between different skills. The four gate layers in LSTM also closely resemble the key parameters in a BKT. The output gate, which reads out the hidden state vector to get an output (i.e. belief state on each skill), is a counterpart of slip and guess rate. The update gate, which scales the update value from the input, is like the acquisition rate. And the forget gate is literally modeling the forgetting of a knowledge component. In a nutshell, LSTM just upgrades the model of these processes from a single real number, to a layer of fully connected neural network. Therefore, we wouldn't be surprised to see DKT vastly outperforms BKT.

DKT can be simplified as a function approximator which takes a student's practice history (sequence of exercises he did, and whether he got each exercise correct or not) as input, and outputs the predictions (belief states) of the probability of getting each skill correct if presented as the next question.

2.2 Implementing DKT

DKT is a good modeling method for the student, i.e., a learned transition model of hidden knowledge states. This could be useful in various parts in our group's projects. Therefore, the first thing to do is implementing the model.

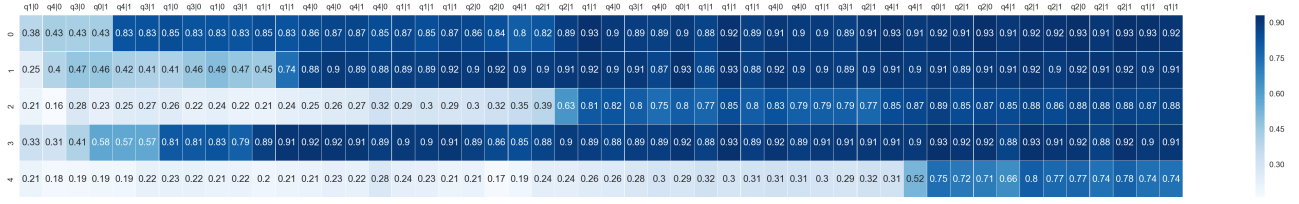
Although there has already been several published DKT implementations, besides being a great opportunity to sort out all the subtlety like sequence padding and dropouts in LSTM, there are several other advantages of having our own implementation. First of all, our implementation is in Python, under the Tensorflow framework. Python is a friendly language for machine learning and data mining practices with its various powerful packages, and Tensorflow is a more flexible deep learning framework, both of which provide convenience not presented in previous implementations. More importantly, our ultimate goal of DKT doesn't stop at getting a prediction metric on the test dataset, but to adopt it to do planning and data mining, thus calling an off-the-shelf LSTM cell from the framework won't work for us, since we won't be having access to the inner state of the LSTM that way. We need to implement the LSTM from scratch to better accustom to our goals. At last, since we have the full access to the model's structure, we made a small step forward of DKT by connecting the readout layer to the trainable initial state vector, which enables our DKT to predict the first answer without any previous input (which is a feature not presented in the original DKT implementation).

2.3 Sanity Checks for DKT

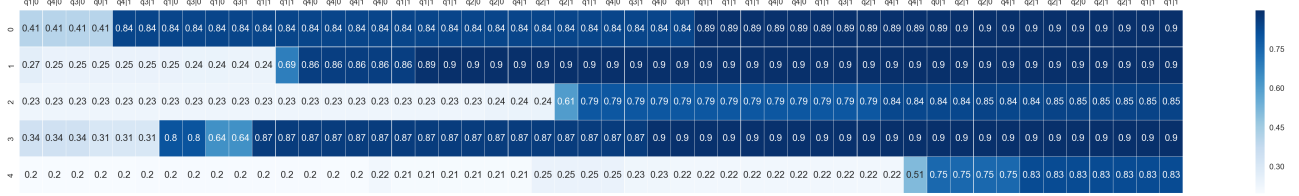
The biggest challenge of implementing a deep learning model is verifying the model. It requires little tuning (which is a great characteristic of using Adam optimizer in stochastic gradient descent) to get an AUC prediction performance comparable with previous works on the test set, however, we need more than a single AUC number to convince us that the model can be trusted.

Besides conventional checks like printing out or plotting several crucial variables/metrics, we did a more convincing sanity check. The goal is to test if DKT can fully recover a BKT model. First, we did a BKT simulation with 5 skills, each has its unique set of BKT parameters. As an extension to vanilla

¹*Deep Knowledge Tracing*, Chris Piech et al, NIPS 2015.



(a) DKT result on student 9999.



(b) Oracle BKT result on student 9999.

Figure 1: An example of the knowledge tracing visualization of DKT and BKT on one of the simulated student trajectory. The row labels are student ids. The column labels, $qi|j$, are student actions in each time step (with the leftmost action as the earliest): the question is on skill i , $j = 1$ when the student answers correctly, otherwise $j = 0$. Each column is the prediction for the action in the column label.

BKT, we add several prerequisite relations into the simulation: skill 0 and 1 are prerequisites for skill 2, while skill 2 and 3 are prerequisites for skill 4. We simulate such relations in a rather simplified manner: the student’s acquisition rate of a skill would stay 0 until he learns all the prerequisites.

BKT model is actually a two-state HMM model, with acquisition rate corresponds to transition model, slip and guess rate corresponds to sensor model. Therefore, we can easily construct an oracle BKT inference model, which knows all the hidden parameters of BKT, and can infer the student’s hidden knowledge state with the simulated trajectories as observations. A small tweak is needed to accustom to the prerequisite relations: the acquisition rate is an expected probability, conditioned on the knowledge states of prerequisites.

DKT can achieve the same AUC, 0.81, with this oracle BKT inference model. Again, only a AUC number won’t be enough (there are millions of ways to accidentally get a high AUC). We computed the MSE (mean squared error) of the step-by-step belief states predicted by both models on the test set and got 0.003, which is a small number, but we still don’t know if it’s small enough. Therefore, we plot the knowledge tracing process given by the two models and compare the trends. From 1 we can see that DKT gives nearly identical predictions with oracle BKT. Combined with all the previous checks, we can safely say that our DKT model can be trusted and adopted in the projects.

In addition, we did a more novel sanity check, in which we reorder the train or test set in various ways and observe the consequent AUC performance. This yields some interesting results and end up being a new project, which is elaborated in section 5.

3 Deep Batch RL (for Education)

In this project, we want to compare two methods, the model-free method (DRQN) and the model-based method (expectimax search with DKT), both based on deep learning models, on off-policy learning for offline (i.e. a batch of) educational data. I focused on the model-based method, including its implementation and evaluation. The model-free method is done by Li Zhou, a graduate student

in our group, and I was involved in the discussions of it.

3.1 Model-based and Model-free RL

In RL, our goal is to find the optimal policy in an unknown environment. There are mainly three ways of doing RL (the latter two are model-free methods, while we only focus on the comparison of the first two):

- First explicitly learn the model for the environment, i.e., the reward and transition model, then find the optimal policy using MDP (Markov Decision Process) method based on those.
- Learn the Q values, i.e. the expected long-term utility of doing action a in state s , $Q(s, a)$, without learning the environment model. With Q values, we can easily derive the optimal policy: $\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$.
- Directly learn the policy, i.e., we are directly mapping a state s to an action a , without explicitly computing any model or utility.

Our hypothesis is, model-free method should hold an advantage over model-based method in various aspects. First, a more apparent one, model-free method is more computational feasible to scale-up, while doing MDP with a world model might involve tree searches, which is exponential. The more relevant advantage of a model-free method is that it's more sample efficient: you need a lot of samples to train an accurate world model (a slightly inaccurate model might suffer from an error explosion in the decision making process, thus result in a bad policy), but way less sample to simply know which action is better given the current state (i.e., no need to know the exact long-term utility).

3.2 RL Going Deep

3.2.1 Deep World Model

DKT is a good world model in education. It's a simple step forward to adopt DKT to do planning. With a DKT, we know the knowledge belief state for the student given any practice history, so we can evaluate any policy based on this.

To be more specific, we do a expectimax search with DKT. An expectimax search tree has interleaving layers of max and expect nodes. A max node pick the children with the maximum value, while an expect node's value is an expectation of the value of its children. In our case, a max node will have children expect nodes each corresponding to an action (i.e. which skill to teach). An expect node will have children max nodes each corresponding to a possible probability for each outcome of the action (i.e., correct or incorrect). DKT can assign a probability to each outcome based on the observed history, thus an expectation can be computed over the children max nodes.

Unlike other tree search like minmax, expectimax tree could not be trimmed without a loss on the accuracy (since it's computing an expectation, and all the branch values are needed). Therefore, expectimax is expensive to compute (as is the case for most model-based algorithms).

Besides conventional speed-up coding tricks like minimize the function calls and the use of abstract classes, we specifically made an optimization for DKT. Since DKT does several matrix multiplications for each input vector, we can stack several vectors into an input matrix and compute the output simultaneously. This trick significantly reduces the computation time to the one sixth of the unstacked version.

3.2.2 Deep Q-Learning

Q-Learning, a model-free RL algorithm, is off-policy, which means that the policy we pick in updating and learning Q-values isn't necessarily the actual policy we run. The Q-values should satisfy equation 1, the Bellman function.

$$Q(s, a) = R + \gamma \max_{a'} Q(s', a') \quad (1)$$

where R is the reward of state s , γ is a discount factor for future's rewards, and other notations consistent with the description in 3.1. From this equation we can get a better understanding of off-policy: the a' we pick is not the action that actually run, but the action than can maximize the Q value in state s' .

For discrete states we can simply put the Q values in a table and do value iterations until the Bellman functions are satisfied (close to iterative methods to solve linear equations). Yet we can't do that for continuous states as in our case, so we want to learn a parameterized function to approximate the Q function. Previous works have shown that this function approximation is not guaranteed to converge when it's not linear, especially when it's a neural network.

The Deep Q-Network² (DQN) adopts several ticks to resolve this diverging issue. DQN is an online (but still off-policy) algorithm, which does online exploration using $\epsilon - greedy$ search during training. Q-learning requires each update to be independant, thus when training the network, DQN employs an experience replay mechanism which feed randomized (s, a, r, s') experience minibatches sampled from the pool of exploration history to the neural network. The loss that the neural network is minimizing is the difference of the left-hand and the right-hand size of equation 1.

The second trick is having a target network, which is identical to DQN but updated with a lower frequency. Since we are doing $\epsilon - greedy$ exploration, the distribution of data acquired during training is actually shifting, thus making the training very unstable (both sides of the Bellman equation is changing). If the two sides of the Bellman equation is provided by the same DQN, it would be like a dog chasing its own tail. The target network, however, can serve as a stable right-hand side of equation 1, thus resolve the diverging issue.

3.2.3 Deep Recurrent Q-Learning

The DQN deals with Atari games, where the state is fully observable from the consecutive 4 frames of the game. However, in our case, we are dealing with a partially observable environment, since we can only observe whether a student gets a problem correct, rather than their actual knowledge status. Therefore, we need to adopt an LSTM to map an action history (observations) to the Q-value. Previous work³ has already added this extention to DQN, resulting in Deep Recurrent Q-Network (DRQN). However, in their work, the training of DRQN is still an online process.

Although DKT and DRQN both adopt an LSTM, both take the trajectory of students' actions as inputs, and even have the same number of outputs (one for each skill), they are very different function approximators (i.e. completely different loss functions). DKT's output is an approximate of the student's knowledge status, i.e. the probability of getting each skill correct, thus a supervised learning model (the result of the next question serves as the label). The output of DRQN is the long-term utility (i.e. rewards) of picking each skill as next action (i.e. question to give). Unlike DKT, the DRQN can be directly used to pick a policy.

²Human-level control through deep reinforcement learning, Google DeepMind, 2016.

³Deep Recurrent Q-Learning for Partially Observable MDPs, Matthew Hausknecht and Peter Stone, 2015.

For real world data, besides being a comparison model of DRQN, another DKT instance, which is trained on a hold-out dataset unseen by DRQN or the expectimax-DKT, should serve as the "ground truth" for the evaluation of the policy picked by both DRQN and expectimax-DKT. However, we start with a comparison on BKT simulation data, where the ground truth BKT model is known, thus a fair comparison between DRQN and expectimax-DKT.

3.3 Challenges in Batch RL

3.3.1 Reward Function Design

In any RL practice, we'll always face the challenge of how to properly define a reward function, since in most cases, the rewards would not be explicitly presented in this data. A myopic or misspecified reward function might result in very undesirable outcomes⁴.

The reward function should be based on the trajectories in the data. We are not looking for myopic definitions like getting a reward whenever question is correctly answered (which would result in the agent constantly giving the same question to hack a high reward). Currently we are using a one-time-only reward associated to each skill, which is awarded once the student get 3 correct in a row.

In our DKT-expectimax framework, any trajectory based reward function can be swapped in. There are more variants of reward functions. A more natural one is to give a decaying reward for questions on the same skill. If we take the applicability into consideration, we probably should also give negative rewards for too many exercises on the same skill since we could run out of exercises in our pool. To make the training of DRQN more stable, we could also try giving awards only when trial ends, where some regularization may be needed. And we could also use discounted rewards, since we want to get rewards in the near future (the sooner the better).

3.3.2 Constraints on Exploration

In RL, when we talk about batch we mean offline (since we have a huge batch of fixed data to train on). For online RL, its single data entry (since we are constantly getting new data through greedy search) or mini-batch. We are doing an offline RL task, thus batch RL.

Several obstacles were encountered after we altered the DRQN into an offline version (the only difference is that we have a fixed pool for experience replay, while in the online version there are always new data coming into the pool). First, the model starts diverging before converging to the minimum loss. Second, though the model can achieve its best performance (which is significantly higher than a random policy) when the loss is at its minimum, it cannot achieve a comparable performance with DKT-expectimax. Currently, the experiments of DKT are only restricted to the BKT simulated data (where we have the ground truth BKT). Table 1 shows the comparison of different models and settings.

As we have shown in section 2, DKT can fully recover the BKT, so we are not surprised to see that a fully-trained DKT (with 5000 instances) can gain an optimal policy. However, in real world, we shouldn't expect DKT to fully recover the underlying model of real students. To better simulate this situation, we train DKT with little data (merely 100 instances). With this DKT, both the AUC on test set and the performance of its policy suffer a huge drop (which is consistent with our prior hypothesis, that DKT's performance is vulnerable to data available).

However, even an inaccurate DKT like that still outperforms DRQN, and in this case DRQN is trained on way more data. This greatly conflicts our prior hypothesis, since DRQN should be more

⁴*Concrete Problems in AI Safety*, Dario Amodei, Chris Olah, John Schulman, Jacob Steinhardt, Paul Christiano, Dan Mane, 2016.

Algorithm	Average rewards after a 20 problem episode
BKT	2.6
DRQN #ins=5000	2.1
DKT #ins=5000	2.6
DKT #ins=100	2.28
Random	0.94

Table 1: Planning performances comparison. BKT and DKT models are used in a 5-step expectimax search for planning. The number of training instances is also listed. BKT-expectimax, with the BKT being the ground truth model, is the perfect policy. Random policy randomly picks a question at each time step, and serves as a baseline. (The 0.8 to 1.0 standard deviations are omitted in the table.)

sample efficient and unless the myopic expectimax which can only look 5 steps ahead (when keeping the computation feasible), DRQN’s Q-value estimation is infinite horizon (i.e. looks infinite steps ahead).

We’ve tried several tricks and fine tuning for DRQN. For example, for offline training, the data distribution is fixed, so maybe we don’t need a target network like the online DQN. However, empirically, the tricks we have tried don’t make a difference.

The breakthrough finding came when we modify the algorithm to online (which is only possible for simulation since we need the ground truth simulator) it can perform on par with the optimal policy. When we do greedy search, it converges fast (because we have a lot of near-optimal sequences), while when we do completely random search, the performance drops after a short boost (just like the offline performance), then come back to the optimal performance eventually, which could be mainly due to the massive (around 10^5 episodes) random dataset (in which there are enough good/near-optimal examples) that has already added in the pool. Greedy search is basically changing the training dataset, so we can easily accept the fact that it’s way better.

The current drawback of offline DRQN is clear: unable to do exploration, thus require a massive amount (10^5 episodes) of random data to see enough near-optimal trajectories. On the contrary, DKT can easily pick up a vague prerequisite structure even under a small dataset (100 episodes), by observing that some skills are easier to learn in the beginning while others are not. Based on this little knowledge, the expectimax can get us an optimal policy. In this current setting, model-based method seems to win on planning performance (especially when model-free method is not that sample efficient as we thought it should be).

3.4 Future Work

Experimental results so far have shown that deep batch Q-learning possibly couldn’t work due to the inability of doing exploration in the environment. It turns out that the generalization ability of deep neural network is not that powerful as we expected. However, we are still trying to fine tune the DRQN to see if its performance could go up. A more promising direction is to integrate return based off-policy RL algorithms⁵ into deep RL. This has not been done by any previous work before thus we might be on to something really innovative.

There are several concerns on the current experiment setting on simulation though. We only have 20 trials for each episode, i.e. our horizon is fairly short. A method as shallow as DKT-expectimax (5 steps look ahead) can perform as near-perfect because we can only expect the students to learn those

⁵*Safe and Efficient Off-Policy Reinforcement Learning*, Remi Munos, Tom Stepleton, Anna Harutyunyan, Marc G. Bellemare, 2016.

skills without prerequisites. If we have a setting where a shallow sight might be catastrophic (like a longer horizon and a more complicated skill structure, with loose prerequisite relations as in the real world), we might see DRQN to win, since it can look infinite steps. However, none of these could account for the low performance of DRQN on the simpler simulation data (it should be able to get the optimal policy).

A minor technical issue we need to point out here. Currently the model-based method is technically not a real RL algorithm, since the reward model is not learned but directly coded in the tree search algorithm. We can use something like DKT to learn the model, but in this case it would be a regression task rather than a classification one. However, this won't happen in application since the reward model in education is defined by human. The fix won't be easy and is not going to help with our main concern (why DRQN is underperforming) either.

There are more interesting directions that we have in mind, including planning with multi-skill questions and making sense of the necessary horizon length for getting an optimal policy (in extreme cases like all the skills are independent, pick any skill wouldn't make a difference, thus zero horizon is as good as infinite horizon). As for application, we should also account for a certain dropout rate when searching for policy. If a student drops out in the middle, a short horizon might be better, since the long-term benefits we planned for the student might not pay off. These thoughts are even more off-topic, and could be discussed in future projects.

4 Adaptive Fraction

This project, Adaptive Fraction, is led by another graduate student in our group, Shayan Doroudi. Our group has some collaboration with an HCI group at CMU, and we run an online intelligent tutor system, Fraction, which tutors grade school level math (on fraction). With this system, we can deploy experiments on real students. However, as we have mentioned before, we want to make sure that our policy would be robust, which means that when we deploy it on our system, we should be sure that the result won't be worse than what we evaluated in (offline) experiments.

4.1 Robust Policy

There are many methods to model a student's knowledge acquisition process, DKT being one of them. Our assumption is, some of these models (or at least some part of it) is correct, and the result can be partially trusted. However, if we only adopt one student model, we might end up with an overfitted policy on this model, and might have negative consequences on real students. Therefore, we need to see what the best or worst our policy could do on various student models.

We will have a matrix to show how different policies (which could be built on one of the student models or worlds) can perform on different worlds. DKT could serve as one world, and expectimax is one way to learn a policy on this world. The matrix would have world and policy as its two axes, and performance as its values. We expect to see that the values on diagonal be higher (since diagonal values correspond to the performance of policies tested on the same world that generated the policy in the first place), and hope to find a policy that's good on all world models. Shayan has already implemented a framework of testing policies on different worlds. My job is to port my DKT implementation, and its simulation code into his framework. DKT has a fairly good AUC, 0.81, on the Fraction dataset.

4.2 Future Work

It's possible to incorporate expectimax with DKT as a policy as well. However, the Fraction dataset's structure is more complicated than simulation or other real datasets: instead of having a single latent skill, each problem consists of multiple steps; students will only meet the same problem once, but different problems can share same steps; different steps within the same problem can have different results; each step is associated with a skill. Therefore, by taking an action (picking a problem), there are 2^{num_steps} possible outcomes, thus expect node would have more children. And the possibility associated with each max node (as a child of expect node) need to come after a multiple step simulation of DKT. We are not sure if we are going to need DKT-expectimax in the framework though, since it would be prohibitive to even look ahead 3 steps (thus very myopic).

In Fraction, after exercises (a tutor session), the student will do a post test. Post test score is our reward. Besides student models, we also use the data to train a reward model. We first feed the tutor session trajectory to a student model, then use the final belief state to do a LASSO regression (linear regression with L1 regularization) on post test score, which gives us the reward model. However, this method gives nearly identical predictions for all policy-world combinations. So one of the major next steps we are thinking is to combine the student and reward model as one (so that we can directly predict the posttest performance given the within-tutor trajectory). This probably would require some tweak to the DKT.

As an extension to the matrix, we are also planning to train student models on a subset of students (above or below average students). Moreover, my implementation of the loose-prerequisite BKT model (with different sets of parameters for different numbers of mastered prerequisites) could be integrated in the framework as well. The current BKT inference model didn't fully leverage the prerequisite structure, since it doesn't update the children or ancestor's belief when encounters a skill. There are still some math to work out for this upgrade.

5 Influence of Skill Orders

This project came up from the investigation of the dataset, Assistments, used in the original DKT paper. There are several data ordering issues in Assistments and preliminary results suggest that the structure inside the sequential data could be discovered by investigating the DKTs trained on reordered data. I'm mainly responsible for the project, and Shayan, as an EDM veteran, would also advice me in this project.

5.1 Data Quality Issues for Assistments Dataset

The original DKT paper reports a 0.86 AUC on the Assistments dataset, which is the largest public educational dataset so far. However, by closely investigating this dataset, several data quality issues emerge. I independently discovered some of those, while a recent EDM paper⁶ has a more comprehensive coverage. Such issues greatly inflate the AUC result of the DKT, which make the reported AUC completely lose its credibility.

Besides trivial misprocessings like keeping the scaffolding problems, there are three major data quality issues. First, a great percentage of the dataset is duplicated, where same record with the same timestamp could be repeated for more than 200 times. This would obviously boost the prediction performance, since the model would observe 200 identical records in a row.

⁶Going Deeper with Deep Knowledge Tracing, Xiaolu Xiong, Siyuan Zhao, Eric G. Van Inwegen, Joseph E. Beck, EDM 2016.

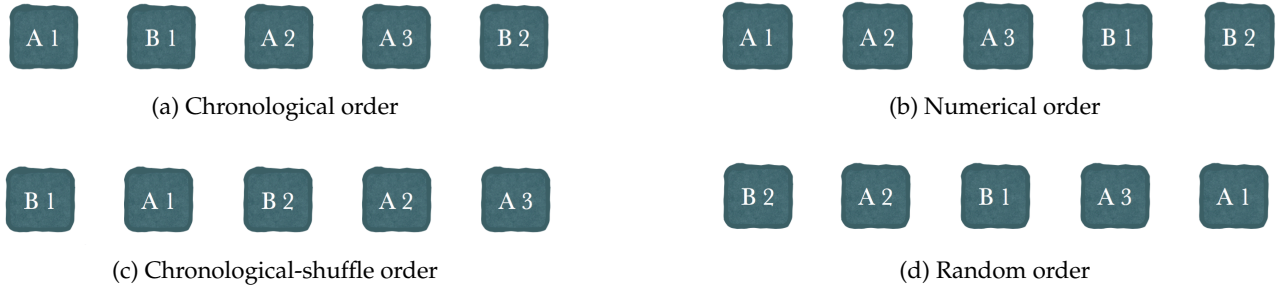


Figure 2: Four order schemes. Skill A has three records, while skill B has two. The numbers stand for the original(chronological) order within the same skill.

The other two issues, which direct us to launching this project, are more subtle and would be described in the following two sections.

5.2 Reordering the Skills

In the original Assistments dataset, records are not presented in chronological order. Instead, the trajectory is organized in what we call a numerical order, where records of the same skill are grouped together, and the chronological order of results (i.e. correct or incorrect) within a skill is preserved. To better understand the influence of the numerical order, we train on data with one of the order schemes, and test on data with another ordering. As comparison, we add two more reordering schemes. Chronological-shuffle means that we shuffle the order of the skill encountered by the student, yet still keep the chronological order of results within the same skill (like what we did in numerical ordering). The baseline order is random order, which completely shuffles the data without preserving any chronological order. The four order schemes are showcased in figure 2.

Figure 3a is the initial result we got on the Assistments data, with the performance of training and testing on chronological order significantly outperforms other order schemes. We can at least draw three conclusions from this result:

- There is a pattern of knowledge acquisition over time for our model to learn (thus chronological > random).
- There is a strong correlation between the acquisition of different skills (thus chronological > numerical).
- The structure of skill acquisition preserves a bit in numerical order, but would be completely broken down in chronological-shuffle (thus numerical > chronological-shuffle).

In a nutshell, the fact that AUC drops significantly after shuffling the skill orders (which simultaneously break the inter-skill structure presented in the data) suggests that we probably could measure skill relations with AUC drops after reordering. This is the decisive result that make us think this project might be promising.

5.3 Multi-skill Questions

However, a closer look at the data reveals another data quality issue that could severely inflate the prediction AUC. Some problems in the Assistments dataset are associated with multiple skills. Such problems are presented as multiple consecutive records, each associated with a single skill. Since they are actually one problem, the correctness (result) is same for the records. An example of this

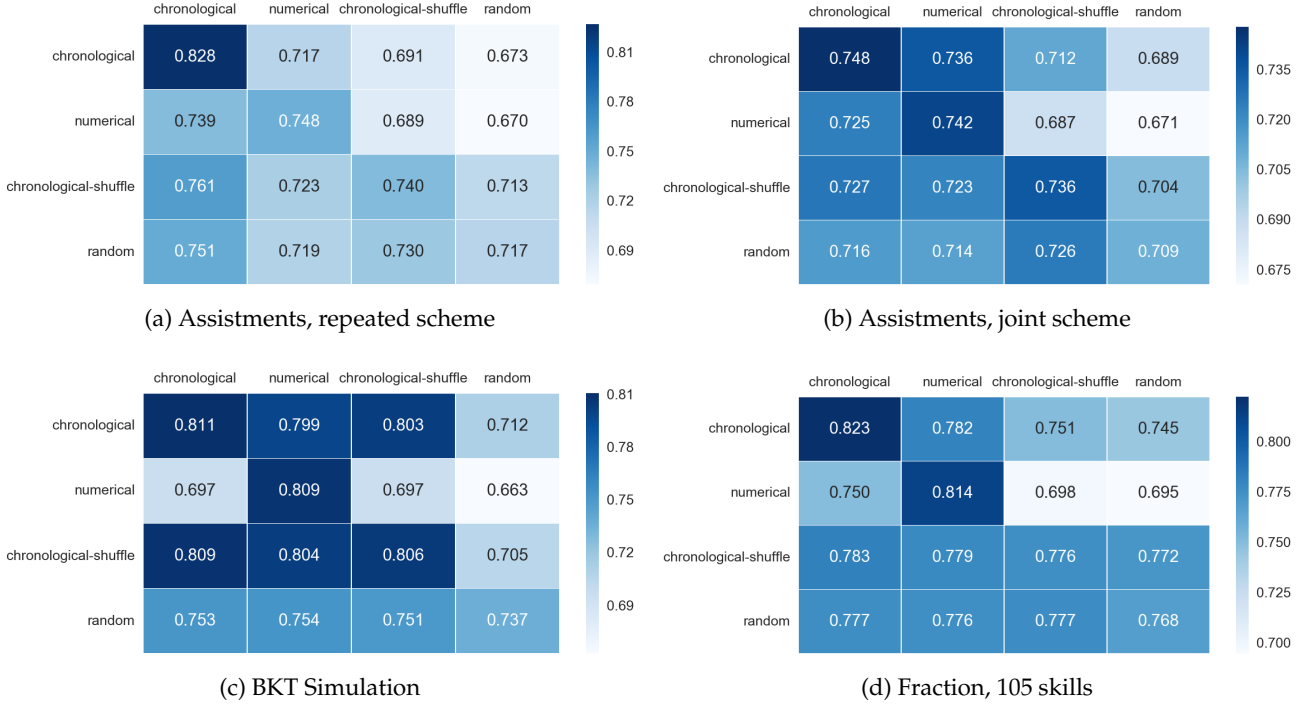


Figure 3: Reordering test results (AUC of prediction) on four datasets. Row labels are the order of the train set, and column labels are the order of the test set.

Time Stamp	Problem	Skill	Correct
5678	100	1	0
5678	100	2	0
5678	100	3	0
5679	101	1	1
5679	101	2	1
5679	101	3	1

Table 2: Repeated scheme: repeated records for different skill tags of the same problem.

repeated scheme can be found in table 2. Such scheme would hint the DKT with the ground truth when DKT takes the first record as an input. For the consecutive records that are actually of the same problem, the result won’t change, thus would definitely be correctly predicted. To confirm this, we separately compute the AUC for the repeated records in the test set, and the AUC is over 0.9999, which is basically a perfect prediction.

To resolve this issue, we need to keep only one record for each problem. We can make every skill combination of multi-skill problems a new skill, which is assigned a new skill id, as shown in table 3. Under this scheme, we repeat the reorder test and the results are shown in figure 3b. We can see that most AUC results don’t change much, yet the AUC of train on chronological test on chronological suffers a huge drop, since in this case there won’t be a hint to leverage as in the repeated the scheme.

Joint scheme has its own deficiencies, since it completely loses the representation of its components. DKT won’t know (at least not explicitly) the sub-skills of a joint-skill, thus might perform worse on predicting other stand-alone records of the sub-skills. We can simply sample one of the skills in the multi-skill question and discard all other records to preserve some of the components of a joint-skill. However, this yields an even worse AUC, since it’s only an approximation of the

Time Stamp	Problem	Skill	Correct
5678	100	145	0
5679	101	145	1

Table 3: Joint Scheme: one record per problem, with a new skill id.

problem.

We made a novel modification to DKT, aiming to resolve such dilemma. The new DKT takes a multi-hot encoding of an action as input, where multiple elements in the vector are 1, each corresponds to a sub-skill (while the original DKT takes one-hot encoded inputs). And when making predictions, the DKT multiplies the belief states of sub-skills based on the assumption that a student has to know all the sub-skills to correctly answer a multi-skill question. To our surprise, this still results in the same AUC as the joint scheme. By taking a closer look at the LSTM, we notice that a multi-hot vector won't be treated as a combination of several one-hot vectors in the computation, but more close to a completely new vector. Therefore it's basically equivalent to the joint scheme.

5.4 Extracting Prerequisite Relations with Reordering

Though not as significant as figure 3a, we can still observe that training on chronological data holds a clear advantage over other order schemes when tested on chronological data. In order to get a sense of the scale of AUC difference we should be expecting, we conducted the same reordering experiments on the BKT simulation dataset. However, as shown in figure 3c, the result is quite surprising: training on numerical order is way worse, while training on chronological-shuffle order is on par with chronological.

On Fraction dataset we observed a descent AUC difference between chronological-shuffle and chronological, as shown in figure 3d. We took a closer look by computing the AUC for each skill separately when trained on chronological-shuffle and tested on chronological. In this case, we won't get a correct DKT model, since the relative order of skills is messed-up in the train set. Therefore, only independent skills would keep the same AUC, while skills with either ancestors or children would suffer a big AUC drop. About 20 out of 105 skills suffer a huge AUC drop. This result gives us more confidence on our method.

We tried another way of reordering with the simulation dataset. We train DKT on the chronological ordered data, but when we test it, we completely shuffle the order of one of the skill while preserving the order of other skills. When we investigate the per skill AUC in the test set, obviously the shuffled skill would have the biggest drop. According to figure 4, the closer the skill in the hierarchical structure is, the bigger the AUC drop. However, we are expecting AUC drops with bigger values.

5.5 Future Work

We decided to write up our findings as a paper. Therefore, besides keeping progressing on the experiments, we also need to start literature search.

As for the experiments, the priority would be to work out the detailed logistics of extracting knowledge components' relations with reordering. There are several concerns: shuffle a single skill or a group of skills; shuffle the train test or the test set; what threshold to set for filtering the noise. To better investigate these issues, we need to train an interpretable model, i.e. a BKT, so that we can check the learned model by comparing it with the ground truth model used in simulation. Based on

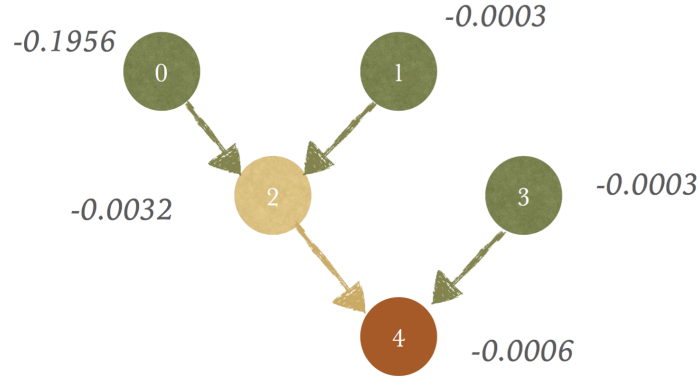


Figure 4: The per skill AUC drop on the BKT simulation test set when skill 0's order is shuffled. Skill 2 takes 0 and 1 as prerequisites, while skill 4 takes 2 and 3 as prerequisites.

this, we will move on to DKT and real world data.

The bizarre results observed in figure 3c also need to be properly explained, since that may lead us to a deeper understanding of sequential data orders and DKT. We have already tried different simulated data with additional attributes like interleaving questions or shifting question distributions to figure out when AUC stays the same after shuffling. However, no conclusion has been reached.

The preliminary literature search suggests several possible advantages of our methods over proposed EDM methods. First, DKT as a deep learning model could capture more subtle skill relations, and unlike many other methods, no domain knowledge is required. As for the only proposed method for extracting skill relations with DKT, it suffers from several deficiencies. This method uses the belief state change of skill j after correctly answering a question on skill i as a metric of the relation between the two skills. Figure 5 shows the tracing results given by DKT when input a trajectory of getting one of the skills correct 20 times in a row. We can see that when a skill's belief state goes up, its prerequisites' belief states would also go up. And after extreme cases like 20 corrects in a row, they would all end up in a near perfect belief.

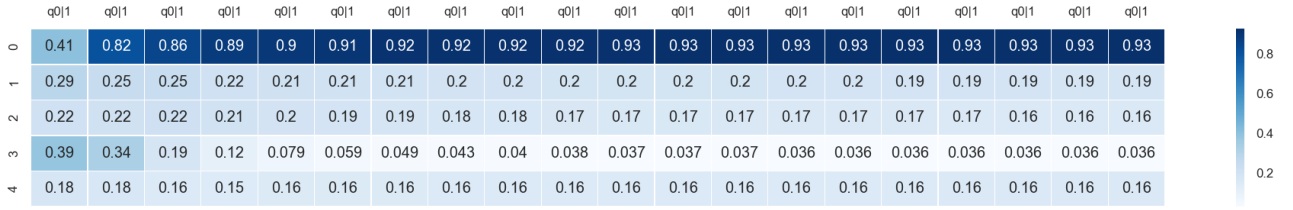
However, in real world data, such strict prerequisite relations won't exist, thus we couldn't tell if a AUC rise is truly resulted from skill relations. In figure 5b, the belief for skill 0 goes up as well, but it's independent with skill 1. In real world dataset, such noise would harm the validity of our conclusions. Moreover, the relations between skills might show up after several observations, as shown in figure 5e. Thus basing our conclusion on only one observation like the proposed method would be too myopic, even for simulated data.

More literature search need to be done to fully motivate and justify our method.

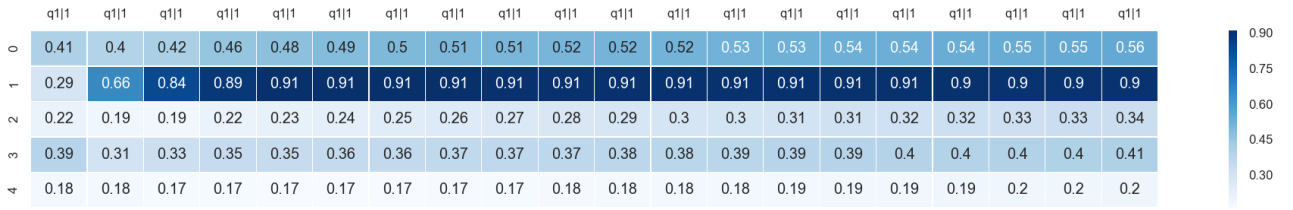
6 Summary

During the past ten weeks, my project started out as repeating a published model and ended up with some new exciting directions to explore. Besides my hands-on experiments, I also participated in a lot paper reading and discussion on RL in the group, which got me acquainted to a whole new area and a group of interesting people.

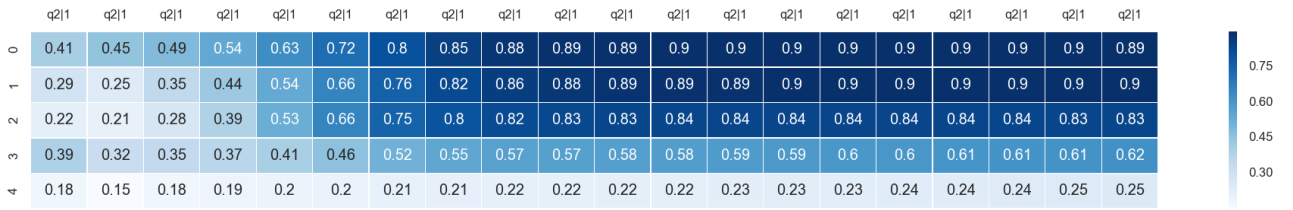
This ten-week internship has been deeply rewarding. As a college-senior-to-be, I've never been more capable of doing research and I've never gotten the opportunity of fully devoting myself to research for such a long time span. This experience has been exhausting and mad fun at the same time, and it further established my mind to be a researcher, and a truth seeker.



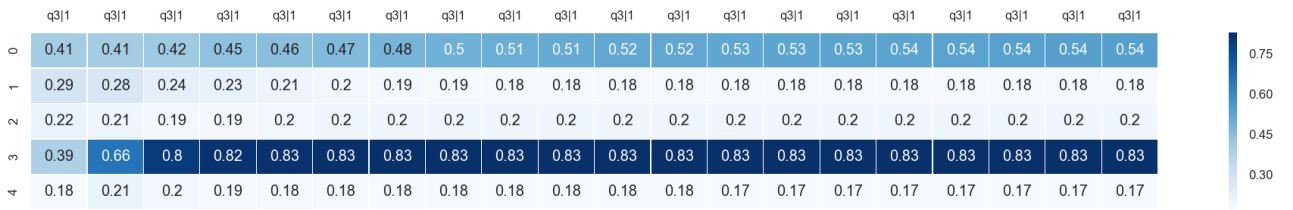
(a) Skill 0



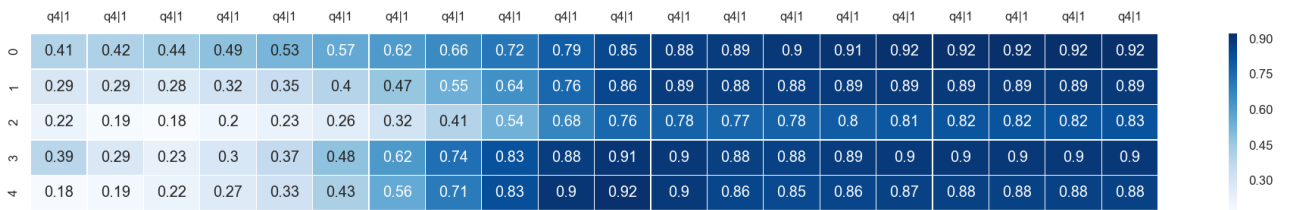
(b) Skill 1



(c) Skill 2



(d) Skill 3



(e) Skill 4

Figure 5: DKT's predictions when getting the same skill correct 20 times in a row. Each column is the belief state for a time step. Each row is the belief states of a skill over time.