# Assisting Learning through Student Mastery Metrics

A Major Qualifying Project Report

Submitted to the Faculty of

Worcester Polytechnic Institute

by

Justin Weintraub

Noah Goodman

Date: 3/1/2024

Project Advisor:

Professor Neil Hefferman

Worcester Polytechnic Institute

# Abstract

When it comes to learning a subject matter or skill, each student can be stuck in different areas and have differing needs. Our work on this MQP was done to provide a means to help out with this variability, researching into methods to assist student learning on the ASSISTments platform, a platform used by students all across the world. We looked into knowledge tracing, a way to predict student mastery in relation to problems. We then used that work to implement a service that uses reinforcement learning to provide tailored hint messages to help students with their varying needs.

# Acknowledgments

Firstly, we would like to thank our advisor, Neil Hefferman, for his continued support and advice for the project. We'd also like to thank our graduate research assistant, Morgan Lee. Morgan gave us tasks throughout our MQP and was a great help with any questions we had. Also, we'd like to thank our teammate, Leo Riesenbach, an IQP student who joined us a term after we started. A substantial amount of the work mentioned in the paper was contributed to by him. Joseph St. Pierre and Ryan Emberling, members of ASSISTments, were both a great help in advising us and helping us out with our tasks, so we'd like to thank them. Lastly, we'd like to thank Ethan Prihar for making the Reinforcement Learning Service that we built off of. We are very thankful for all the support we got from everyone over the course of this project.

Table of Contents

# 1 Introduction

ASSISTments is a free learning platform for students, where students are able to learn subject matter and do homework problems off of it. It is used by many school systems in the US, covering the school years K-12 with a focus on middle school mathematics. As of late, there's been an effort to individualize student learning. And by that, I mean tailoring the problems, hints, and the like for each student. This is where our work comes in. The Reinforcement Learning Service (RLS) is a service that chooses what hints to give students based on certain criteria. Our work, as described later, was to get this service up and running. The criteria that this service was based on was another part of the work we've done. We also did work on the Mastery Service. We learned about the Mastery Service and its underlying algorithms, as well as helped to get it up and running. This mastery value is one such metric now being used by the RLS service. This paper will go into more detail about our work in these two services in a mostly chronological order.

# 2 Mastery Service

## 2.1 Background

Our first major form of work for our MQP was on the "Mastery Service", a service that uses knowledge tracing to determine student knowledge of skills. With a skill being something like the "Pythagorean theorem" or "Basic addition". To determine student knowledge, knowledge tracing was used, which is a section of algorithms that determine student mastery over tasks. There are multiple types of knowledge tracing algorithms that were researched, as discussed in the next section. Our main work in knowledge tracing was to help get the service up and running by doing tasks such as creating a knowledge tracing app, creating displays, and demoing said deliverables so the knowledge tracing app could get implemented. We also helped to get data for a paper that will compare the different forms of knowledge tracing on ASSISTments data over the years.

## 2.2 Types of Knowledge Tracing

Part of our work was researching different forms of knowledge tracing, both for understanding them and for implementation.

### 2.2.1 Bayesian Knowledge Tracing

Bayesian Knowledge Tracing (BKT) is one of the simpler forms of knowledge tracing, and the one being used by the Mastery Service. BKT is formulaic, and its formula (Catherine) is as described by Code 1.

```
function bkt_prediction(initial_mastery, correct, learns, forgets, guesses, slips) {
  mastery_posterior = 0;
  if(correct == 1) { // user got problem correct
    mastery_posterior = (initial_mastery * (1-slips)) / (initial_mastery * (1-slips) + (1-initial_mastery) * guesses)
  } else {
    mastery_posterior = (initial_mastery * slips) / (initial_mastery * slips + (1-initial_mastery) * (1-guesses))
  }

  return (1-forgets) * mastery_posterior + learns * (1-mastery_posterior)
}
```

Code 1. The formula for calculating mastery.

The mastery equation has multiple components based on the skill's characteristics. The initial

mastery is a variable between 0 and 1. "Learns" is the chance of the user learning how to do the

problem upon getting it right, "Forgets" is the chance of the user forgetting how to do a problem,

"Guesses" is the chance of the user not understanding the problem but getting it right, and

"Slips" is the chance of getting something wrong despite the student already knowing the

problem. A new mastery value is calculated from this equation, with the mastery value

determining the chance the student would get the problem correct. Despite BKT being relatively

simple, it's been found to be effective for the Mastery Service.

## 2.2.2 Deep Knowledge Tracing

We also researched other forms of knowledge tracing that weren't used by the Mastery Service.

Deep Knowledge Tracing (DKT) is a form of knowledge tracing that uses a deep neural network

to determine mastery. Specifically, DKT uses long short-term memory and a recurrent neural

network to provide a flexible model/one with a high area under the curve measurement

performance (Mohammad). DKT is passed in time series data, with each data point being a tuple

that consists of the exercise answered and whether or not they got the exercise correct. In large

feature spaces,  dimension reduction of the input is performed. Its loss function is binary

cross-entropy loss, and its output, like BKT, is the probability of a student answering the problem

correctly (or in other words the mastery value). It's argued as to whether or not DKT's higher complexity leads to higher performance or not than that of BKT.

### 2.2.3 Performance Factor Analysis (PFA)

While not directly knowledge tracing, Performance Factor Analysis (PFA) is another method used to predict student performance. It uses logistic regression, a type of machine learning model (Philip). The variables it uses related to student performance are variable, typically using information from past successes or failures for predictions. PFA allows for fast optimization and is good for finding new parameters, but has trouble giving predictions with high certainty (not returning 0 and 1 predictions for mastery, instead being somewhere in between). We researched this as an alternative method of prediction.

### 2.2.3 Other Types

There are other forms of knowledge tracing, such as with factorization machines and performance factors analysis, that don't directly use knowledge tracing. While we researched less into these types, they will potentially be looked into for the paper comparing different types.

## 2.3 Analyzing Knowledge Tracing (Google Collab)

To better understand BKT and logistic regression we ran code on a pandas dataframe to test the model and better understand it. We first ran a dataset through the model using the pyBKT python library and evaluated the model for its accuracy, area under curve, and root mean squared error. We tested the data with and without forgetting, a metric that determines if a student has lost the

ability to do a skill.  We then tested the model with differing minimum sequence lengths for the

amount of times a student has done a problem. A length of 0 means there is no limit.

```
[ ]  flist = [0,1,5,10]

  ▶   for f in flist:
          print( "filter:", f)
          filtered_df = grouped_df.filter(lambda x: len(x)>f )
          for x in [False, True]:
            model.fit(data=filtered_df, defaults= defaults, forgets = x)
            auc = model.evaluate(data=sample_df, metric='auc')
            acc = model.evaluate(data=sample_df, metric="accuracy")
            rmse = model.evaluate(data=sample_df)
            print("forgets:",x,"auc:",auc,"acc:",acc,"rmse:" ,rmse)

  ⮑   filter: 0
      forgets: False auc: 0.8255084379088086 acc: 0.7925759487826649 rmse: 0.3838536774471284
      forgets: True auc: 0.9061435450026515 acc: 0.8359718332016648 rmse: 0.33092437854939954
      filter: 1
      forgets: False auc: 0.8293774270748058 acc: 0.793168144992557 rmse: 0.3851687716710273
      forgets: True auc: 0.9027610168036746 acc: 0.83373121222723 rmse: 0.33356504748465043
      filter: 5
      forgets: False auc: 0.813767889909291 acc: 0.7905690866253101 rmse: 0.38746772985214223
      forgets: True auc: 0.8977969699644766 acc: 0.8399148440334029 rmse: 0.3334078002490736
      filter: 10
      forgets: False auc: 0.821833570054722 acc: 0.8076357990166206 rmse: 0.3861534229464671
      forgets: True auc: 0.911246878494844 acc: 0.8547762352764398 rmse: 0.32281163888558767
```

Figure 1. Our Tests for Differing Amounts of Data with and without Forgetting and Filters.

```
from sklearn.linear_model import LogisticRegression
clf = LogisticRegression(random_state=0).fit(x, y)
clf.predict(x)
pd =clf.predict(x)
print(clf.predict_proba(x))
# get auc and rmse
print(clf.score(x, y))
print(pd[np.where(pd==0)])
print(pd)

[[0.28973664 0.71026336]
 [0.28692081 0.71307919]
 [0.28412139 0.71587861]
 ...
 [0.27857229 0.72142771]
 [0.27582282 0.72417718]
 [0.27309022 0.72690978]]
0.7686714851144382
[0. 0. 0. ... 0. 0. 0.]
[1. 1. 1. ... 1. 1. 1.]
```

Figure 2. Shows difference in area under curve with just linear regression

As shown in Figure1, we found that including forgetting made the model better at predicting

mastery. The amount of filtering didn't have much of an effect on the results. We also ended up

testing the model with logistic regression using the scikit-learn python package and got a worse

value for accuracy than with BKT as shown in figure 2.

## 2.4 Displaying Mastery

One of our tasks was to show how the mastery value would be used in ASSISTments. We also received experience with elements of design by using Figma to create a mockup of the UI.



Figure 3. Mockups made in Figma for the Mastery Progress Bar

In Figure 3, we placed a mastery progress bar in the top right corner displaying the current mastery of the student in this specific skill or problems. There is also an information button that can be used to display more information such as the specific percentage of mastery or how much mastery they had gained in that session. Color would indicate how close the student was to mastery.

## 2.5 Creating a Knowledge Tracing Application

One of our major tasks was creating a demo for knowledge tracing, that would showcase a student's mastery being updated in real-time for different skills on an assignment. The app could select an assignment, and when doing so it would give a display such as that in Figure 4. Every skill related to an assignment gets its own mastery bar and other features.
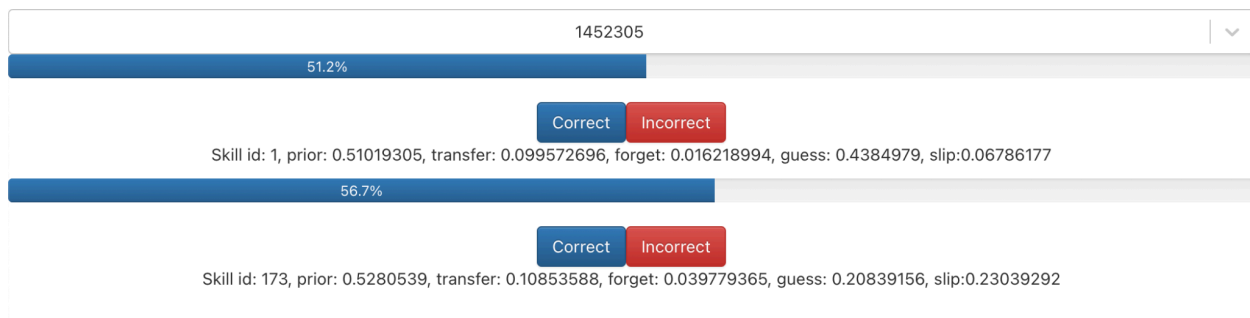


Figure 4. The Knowledge Tracing Application's View for Assignment 1452305.

Depending on the skill's fitted parameters, pressing Correct or Incorrect raises the mastery bar by a certain amount, as listed in Code 1. The initial mastery value is shown on the progress bar for the skill. This demo assumes only one user is doing the assignment, the one interacting with the demo.

To get the data for the assignments and skills, SQL queries on the ASSISTments database were performed. We joined the tables for assignments, skills, student mastery, and any other tables necessary to do the linkage. We stored that information in a local Postgres database. The app made for the demo had a Postgres database, a Node.js backend, and a React frontend. The database's mastery value for a skill would get updated every time correct or incorrect is pressed. After this app was completed, we presented our work and it spurred the work on putting mastery into ASSISTments.

## 2.6 Creating Student Datasets

Part of our work was getting logs from students' problem attempts sorted year by year, to be used for a paper previously mentioned on testing different forms of knowledge tracing on student data. Our work on getting logs was done to compare/evaluate the differences between the data every year (with every year being the last 4 school years) for the paper. This work is building off of a previous one made by Morgan Lee et al. that looked into mastery in previous school years, but didn't look into investigating different forms of knowledge tracing and dealt with oddities in data as the pandemic was starting. To do our work, this required SQL queries the linked data between the problem logs and assignment logs tables. We wanted to get data from assignments that were assigned in that school year, and the performance of the students on said assignments, as shown by Code 2.

```sql
create or replace view assign19_20 as
select * from core.assignments where assign_date between '30-May-20 23:00:00.000' and
'31-May-21 00:00:00.000' ;
select pl.discrete_score as discrete_score, al.user_xid as user_xid, a.sequence_id as assignments,
metadata_value_id as skill_id
from (select * from student_data.problem_logs pl
where pl.start_time between '30-May-20 23:00:00.000' and '30-May-21 00:00:00.000') pl
inner join student_data.assignment_logs al
on pl.assignment_log_id = al.id
inner join assign19_20 a
on al.assignment_xid = a.id
inner join cas_core.legacy.assistment_to_sequence_associations a2s
```

```
on a.sequence_id = a2s.sequence_id

inner join cas_core.legacy.problems p

on a2s.assistment_id = p.assistment_id

inner join

(select * from cas_core.legacy.metadata_taggings where metadata_target_id = 1) mt

on p.id = mt.object_id;
```

Code 2. The query for a specific school year.

Doing this method of getting all the logs from each year proved to be ineffective, as the data

gathered was simply too large and the queries took too long. We applied several filters to limit

the data to what was the most important for gathering inferences. First off, we limited

assignments to only be for school years 6-7, as that's where the majority of data was from. We

selected the assignments that were the most frequently used, being the top 500 from each grade.

These steps significantly reduced the file sizes, the final step for creating the student datasets.

# 3 Reinforcement Learning Service (RLS)

## 3.1 Background

RLS or Reinforcement Learning Service is a program to be used in ASSISTments to recommend specific help messages based on the student's mastery of the specific skill being tested. This is so that it can better help the student based on how proficient they are at the topic. For example, in the case of a problem on fractions, a student could be given a hint about how division works, or one related to how to add two fractions together.

RLS, as the name mentions, works through reinforcement learning. Reinforcement learning is a means of a model taking actions through getting or not getting rewards for their past actions. In this case, if a user gets the next problem correct after solving a problem with a hint given by RLS, that hint will be more likely to be used again under the same circumstances for another student. There is also a set chance for RLS to select a random hint to give, as a means of discovering potentially new valuable options. As RLS makes decisions and analyzes the outcomes, it will get better at helping students.

We had multiple tasks related to RLS. Our first major task was linking the Mastery Service into RLS. How RLS used to work was that it used the fraction of problem correctness to determine what hints to give to the student. However, if a student is learning something new, that isn't a great metric to use due to a lack of data. So instead, we got to work at using mastery as a parameter, due to that being better suited for the prediction task. An empirical study in the future will be made to determine the effectiveness of using mastery in RLS. Our next major task was getting RLS up and running in ASSISTments. In addition to making it so it could run on ASSISTments, we had to make sure it'd be safe to run in the production stack of

ASSISTments (in comparison to the testing stack). To do that, we had to do testing and code cleanup. RLS is not yet up and running, and will be in the future.

## 3.2 Linking the Mastery Service to RLS

To do the task of linking the two services, we needed to update the RLS codebase to get the data from the Mastery Service. Some of the work we did to get data from the Mastery Service can be seen in Code 3.

```python
def get_mastery(conn_dict: dict[Connection], user_id: str, skills: pd.DataFrame) -> (int | None, np.float64 | None):
    '''
    Gets the mastery values for a student, given a user ID and the relevant skill(s).

    If there are multiple skills, then it takes the mastery value for the one with
    the highest n_observations (most data points used for determining mastery).
    '''

    if skills.empty or pd.isnull(user_id):
        return None

    mastery_query = f'''
                SELECT sm.id, sm.mastery, er.target_id as user_id, sm.skill_id
                FROM student_data.student_mastery sm
                INNER JOIN (
                    SELECT * FROM core.external_references WHERE xref = \'{user_id}\' AND xref_type_id = 1
                ) er ON sm.user_xid = er.target_id
                WHERE sm.skill_id IN ({", ".join(str(i) for i in skills["skill_id"])})
                ORDER BY sm.n_observations DESC
                LIMIT 1
                '''

    query_result = sqlio.read_sql_query(sql_text(mastery_query), conn_dict['dev'])
    if query_result.empty:
        return None, None
    row = query_result.iloc[0]
    return int(row['id']), row['mastery']
```

Code 3. Getting Data from the Mastery Service

As can be seen by Code 3, to get mastery for a specific skill, we needed to do a query between the two tables and a connection table. "N_observations" is used to choose the mastery value that is the most recent, as there is a mastery value for every instance a student has done a problem in the database. After getting the mastery data, we updated the preexisting RLS code to use the

mastery value instead of user correct whenever possible (if there exists a mastery value). With that done and some test cases, the linkage was complete.

## 3.3 Containerizing RLS

One of the major tasks we had was containerizing RLS, or in other words, making it so all of the service's dependencies are installed in one go in an environment that is accessible anywhere. This was done so we could easily install the RLS service onto ASSISTments's EC2 instances (computing machines). To do that, we had to first figure out what dependencies needed to be installed. We did a clean run-through of installing all the dependencies from scratch on a computer, specifically the pip dependencies necessary for our Python application to work. With all of the requirements stored in a requirements.txt file, we then created a pyproject.toml file, which packaged the RLS service and its dependencies. This would allow us to simply import RLS and its functionality. With that done, we passed on this work to Joe St. Pierre so he could create a pipeline for making Docker images that when containerized would run the packaged RLS we generated in a simulated environment. This Docker image would be put on ASSISTment's ECR, a storage location running on AWS.

## 3.4 Code Cleanup

To make the code more readable and efficient we went through a code review. We first went through the codebase and made TODO comments on what needed to be changed and then worked to improve it. One such problem was that the previous code was taking in parameters to functions through the function's arguments and it would not tell you what connection you would need until you looked at the function itself. We turned the connections into a dictionary as an

improvement. This dictionary would be passed to the functions and they would use specific features from the dictionary. Since the code was written in Python, we implemented type hinting for more readability as it describes what type each of the arguments are as well as what the return type of the functions is. For an example of us improving the clarity of the code, we simplified a variable declaration by making repeated values their own variables and passing them into the original declaration as shown in Code 4.

```
22    -         # TODO: Pull len(...['option']) out into a local for clarity
23    -         d = len(self.feature_names_dict['user']) + \
24    -             len(self.feature_names_dict['problem']) + \
25    -             len(self.feature_names_dict['option']) + \
26    -             len(self.feature_names_dict['user']) * \
27    -             len(self.feature_names_dict['option']) + \
28    -             len(self.feature_names_dict['problem']) * \
29    -             len(self.feature_names_dict['option'])
         23  +  num_users = len(self.feature_names_dict['user'])
         24  +  num_problems = len(self.feature_names_dict['problem'])
         25  +  num_options = len(self.feature_names_dict['option'])
         26  +  d = num_users + num_problems + num_options + (num_users * num_options) + (num_problems * num_options)
```

Code 4. How we Restructured the Code to be More Readable (Red is old Code and Green is new).

We also created more comments and documentation for the code to make it more understandable including adding docstrings to functions. This is also being continued after we leave by our teammate Leo.

## 3.5 Unit Testing

We created test cases for many functions in the program. This testing included queries, functions, and model tests. The model used was Dynamic Linear Epsilon Greedy, which was used for choosing the help message. For the model tests we simulated situations with random choice, a single option, and with two options weighted and unweighted as seen below in Code 5.

```python
    def test_two_options_uninitialized(self):
        '''
        Tests that if the weights and biases are the default values,
        exploitation choice will choose the first option.
        '''
        options_dict = self.tutor_strategies_dict
        exploitation_choice = self.model.exploitation_choice(self.user_features,self.problem_features,options_dict)
        self.assertEqual(exploitation_choice, 1)

    def test_two_options_initialized(self):
        '''
        Tests that if the weights are the default values and biases are equal to all 1,
        exploitation choice will choose the second option due to higher values for the second option.
        '''
        self.model.b += 1 # Add 1 to all biases so that different options have different reward values
        options_dict = self.tutor_strategies_dict
        exploitation_choice = self.model.exploitation_choice(self.user_features,self.problem_features,options_dict)
        self.assertEqual(exploitation_choice, 2)
```

Code 5. Test Cases for Testing the Decision Model.

Another test case tests various edge cases for functions that get data from the database including

missing skill and mastery data as well as checking whether the correct values are gotten from the

get functions as shown in Code 6.

```python
def test_mastery_of_problem_with_no_skills(self):
    '''
    This test case checks to see if mastery is not given skill \
    tags, then no mastery values are gotten
    '''
    mastery = get_mastery(conn_dict, get_testing_student_xref(), pd.DataFrame())
    self.assertIsNone(mastery)

def test_mastery_and_skills_gotten_when_valid_fields(self):
    '''
    This test checks to see if given an user and problem pair \
    that has skill values and mastery values, both are \
    correctly gotten from the code
    '''
    ids = get_testing_data_ids()

    # first checking to see if skills are gotten
    skills = get_skills(conn_dict, ids['problem_id'])
    self.assertIsNotNone(skills)

    # now checks to see if mastery is gotten
    mastery = get_mastery(conn_dict, ids['user_id'], skills)
    self.assertIsNotNone(mastery)

def test_insert_prior_knowledge_data(self):
    '''
    Checks the different cases of the function insert_prior_knowledge_data, \
    checking to see where mastery is gotten from in each test
    '''
    user_features = { 'user_correct': 'test' }
    user_features_no_mastery = { 'user_correct': 'test', 'prior_knowledge': 'test' }

    # first check to see if user_correct is used if no mastery_features
    res, used_mastery = insert_prior_knowledge_data(deepcopy(user_features), None)
    self.assertEqual(res, user_features_no_mastery)
    self.assertFalse(used_mastery)

    user_features_with_mastery = { 'user_correct': 'test', 'prior_knowledge': 0.1}
    res, used_mastery = insert_prior_knowledge_data(deepcopy(user_features), 0.1)
    self.assertEqual(res, user_features_with_mastery)
    self.assertTrue(used_mastery)
```

Code 6. Shows the test cases for getting values from the database.

## 3.6 Performance Testing

To make sure that RLS could go up and running without issues, we needed to perform a substantial amount of performance tests on the various functionality of RLS.

### 3.6.1 Running Standard Scripts

Our first main performance test was checking to see how the setup and update script ran in a new environment for RLS. The setup script initialized the RLS in a specified environment. So for setup, we purged information about RLS from the testing RDS (the host for the data on AWS), and then ran the script on the testing EC2. After some debugging of various small errors that occurred (such as needing to rename a column in a database), the script ran for approximately 5 hours without a hitch, initializing RLS on the testing EC2. The testing EC2 gets dumped on with up to date data from ASSISTments every week. So when the new week came around, we had a new batch of data to work with for the update script. After a similar debugging process, we ran the update script, and after a few hours, it finished normally.

### 3.6.2 Generating Synthetic Recommendation Logs

Another type of performance test done was to see how RLS would process newly made recommendation logs, as the update script wouldn't simply be testing this task. To do this test, we had to generate fake unprocessed recommendation logs. Even though they were fake, they were based on real data. To create these recommendation logs, we gathered information from the assignment logs (a table that shows what assignments each student has done), problem logs (what problems the student did during the assignment), and tutor strategies table (a table that has a list of possible hints a student could be given for a problem). We got data from 1000 random assignment logs and got their associated problems' data. For every problem, we got the tutor strategies associated with it. After that, we created rows in the user features, problem features, tutor strategy features, and finally, recommendation logs. These tables were all accordingly given

fake data and linkage to the associated prior tables. This way, we'd have a large set of random

features and recommendation logs to be used for the processing.

# 4 Future Work

As previously alluded to, there's still a lot of work left to do for the Mastery Service and RLS. With the Mastery Service now implemented, there's a paper being made about comparing different forms of knowledge tracing on ASSISTments data over the years. And as of when this paper was written, RLS hasn't gone live, being still in the works. There is an ongoing process of testing, cleanup, and implementation for this. RLS needs to be handled with care, due to the potential of giving unhelpful hint messages. We hope our work can help with student learning for ASSISTments users nationwide.

# Bibliography

Mohammad Khajah, et al. "How Deep Is Knowledge Tracing?" arXiv [Cs.AI], 2016,

        http://arxiv.org/abs/1604.02416. arXiv.

Morgan Lee, et al. "Knowledge Tracing over Time: A Longitudinal Analysis". Proceedings of

the 16th International Conference on Educational Data Mining, International Educational Data

Mining Society, 2023, pp. 296--301, doi:10.5281/zenodo.8115788.

Philip Pavlik Jr, et. al. "Performance Factors Analysis – A New Alternative to Knowledge

        Tracing", Carnegie Mellon University,

        https://pact.cs.cmu.edu/koedinger/pubs/AIED%202009%20final%20Pavlik%20Cen%20

        Keodinger%20corrected.pdf.

Catherine Yeh, et. al. "Predicting What Students Know". Computer Science Wlliams.

        https://www.cs.williams.edu/~iris/res/bkt-balloon/index.html.