

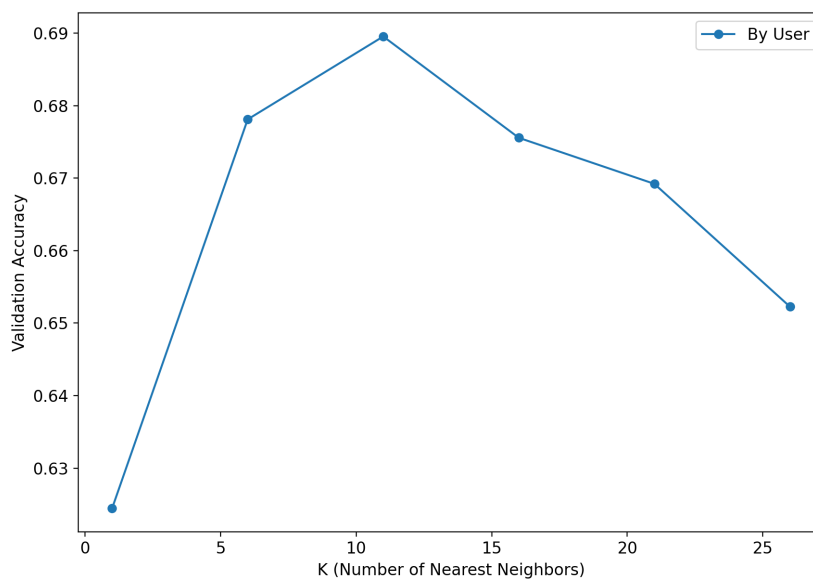
# CSC311 Final Project

Xinyu Kang  
Xinyan He  
Xin Peng

December, 2020

## 1 Q1(Xinyan He)

(a) Plot for the accuracy on the validation data:



Report the accuracy on the validation data:

When k is 1, validation accuracy: 0.6244707874682472  
When k is 6, validation accuracy: 0.6780976573525261  
When k is 11, validation accuracy: 0.6895286480383855  
When k is 16, validation accuracy: 0.6755574372001129  
When k is 21, validation accuracy: 0.6692068868190799  
When k is 26, validation accuracy: 0.6522720858029918

(b) For the item-based collaborative filtering methods:

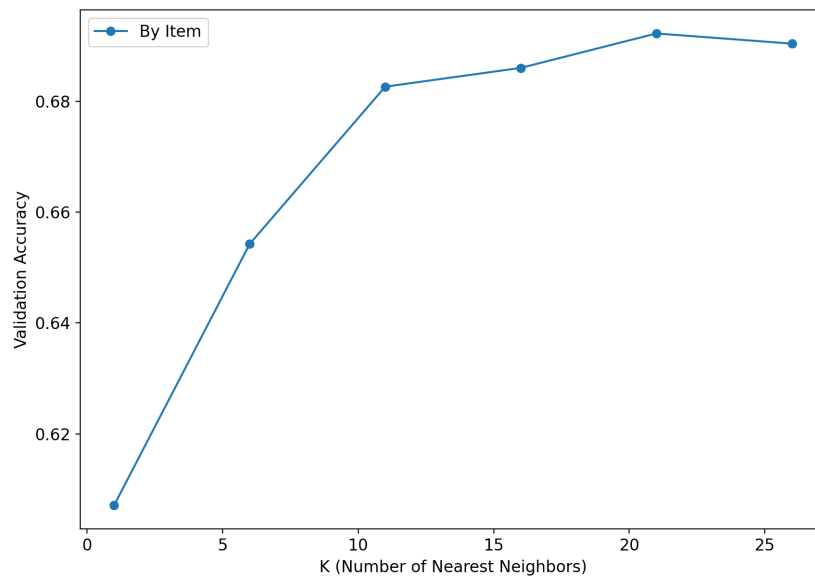
Report the chosen  $k^*$  and the final test accuracy:

$k^*$  is 11 with test accuracy 0.6841659610499576

(c) Report the accuracy on the validation data:

When k is 1, validation accuracy: 0.607112616426757  
When k is 6, validation accuracy: 0.6542478125882021  
When k is 11, validation accuracy: 0.6826136042901496  
When k is 16, validation accuracy: 0.6860005644933672  
When k is 21, validation accuracy: 0.6922099915325995  
When k is 26, validation accuracy: 0.69037538808919

**Plot for the accuracy on the validation data:**



**Report the chosen  $k^*$  and the final test accuracy:**

$k^*$  is 21 with test accuracy 0.6816257408975445

(d) We could see that the user-based collaborative filtering have a better performance on test data.

(e) **List at least two potential limitations of kNN for the task given:**

(1) **kNN is very sensitive to noisy data.**

Suppose here we are using the user-based collaborative filtering, we are predicting whether a student can answer a question correctly or not based on users that are closet to it. When the nearby data is not so accurate or has some mistake, this would directly lead an incorrect prediction of whether the target student can answer correctly or not.

(2) **kNN has a slow running time with large data.**

As kNN is based on distance calculation. As in our given task, the sparse matrix we are used is (542, 1774) which is quite huge. So here running time would be quite slow especially when we are choosing a large k.

## 2 Q2(Xinyu Kang)

- (a) Derive the log-likelihood  $\log p(C|\theta, \beta)$  for all students and questions.

From

$$p(c_{ij} = 1|\theta_i, \beta_j) = \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}$$

We could also get

$$p(c_{ij} = 0|\theta_i, \beta_j) = 1 - \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}$$

$$p(C|\theta, \beta) = \prod_{i=1}^{542} \prod_{j=1}^{1774} \left( \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)} \right)^{c_{ij}} \left( 1 - \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)} \right)^{1-c_{ij}}$$

Then we calculate the  $l(\theta, \beta) = \log(p(C|\theta, \beta))$

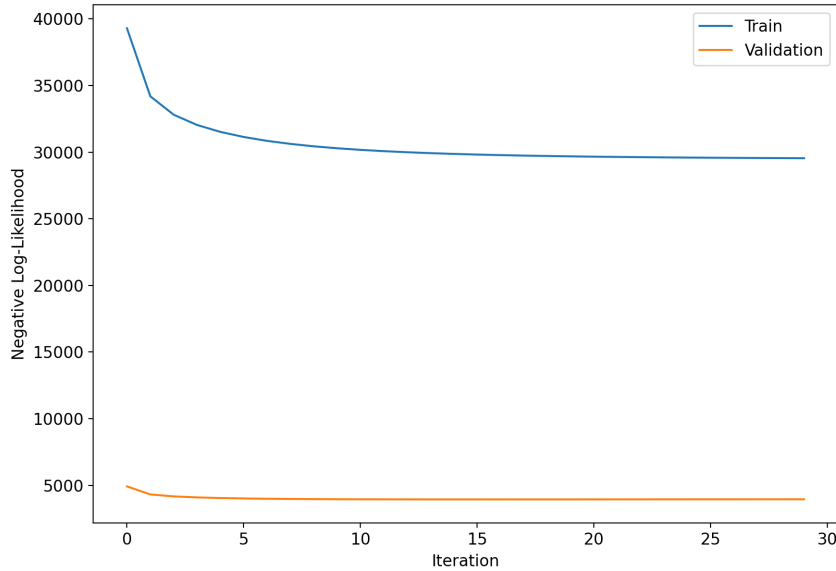
$$\begin{aligned} l(\theta, \beta) &= \sum_{i=1}^{542} \sum_{j=1}^{1774} c_{ij} \log\left(\frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}\right) + (1 - c_{ij}) \log\left(1 - \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}\right) \\ &= \sum_{i=1}^{542} \sum_{j=1}^{1774} \left\{ c_{ij} [(\theta_i - \beta_j) - \log(1 + \exp(\theta_i - \beta_j))] + (1 - c_{ij}) \log\left(\frac{1}{1 + \exp(\theta_i - \beta_j)}\right) \right\} \\ &= \sum_{i=1}^{542} \sum_{j=1}^{1774} \left\{ c_{ij} [(\theta_i - \beta_j) - \log(1 + \exp(\theta_i - \beta_j))] - (1 - c_{ij}) \log(1 + \exp(\theta_i - \beta_j)) \right\} \\ &= \sum_{i=1}^{542} \sum_{j=1}^{1774} \left\{ c_{ij} (\theta_i - \beta_j) - c_{ij} \log(1 + \exp(\theta_i - \beta_j)) - \log(1 + \exp(\theta_i - \beta_j)) + c_{ij} \log(1 + \exp(\theta_i - \beta_j)) \right\} \\ &= \sum_{i=1}^{542} \sum_{j=1}^{1774} \left\{ c_{ij} (\theta_i - \beta_j) - \log(1 + \exp(\theta_i - \beta_j)) \right\} \end{aligned}$$

Derivative of the log-likelihood with respect to  $\theta_i$  and  $\beta_j$

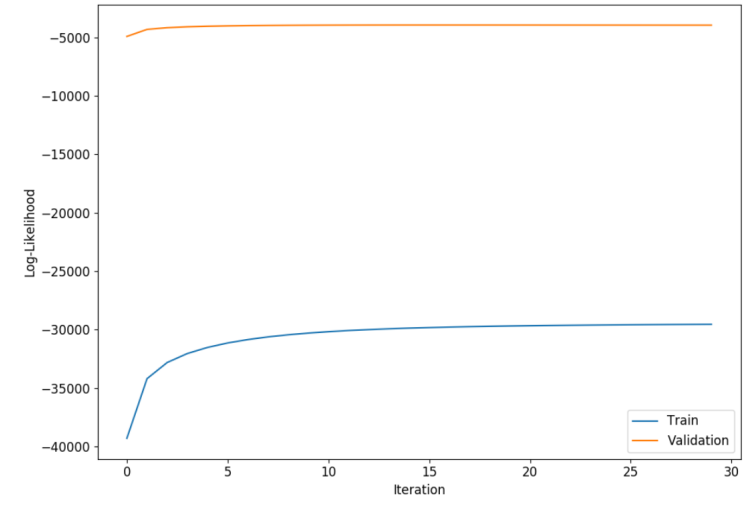
$$\frac{\partial l}{\partial \theta_i} = \sum_{j=1}^{1774} \left( c_{ij} - \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)} \right) \quad (1)$$

$$\frac{\partial l}{\partial \beta_j} = \sum_{i=1}^{542} \left( -c_{ij} + \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)} \right) \quad (2)$$

- (b) We have chosen learning rate to be 0.02 with 30 iterations. And we initialize theta and beta to zero vectors. The plot of negative log-likelihood as a function of iterations is as follows.

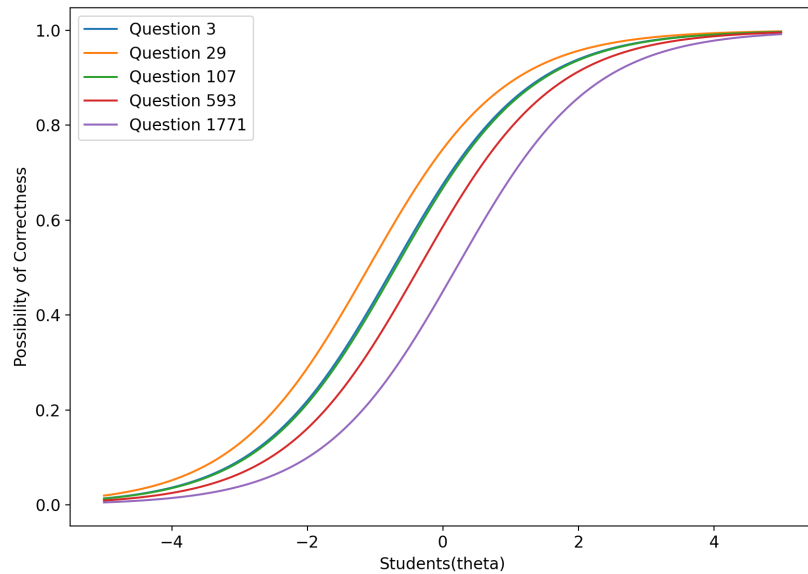


The plot of log-likelihood as a function of iterations is as follows.



(c) Validation Accuracy: 0.7074513124470787  
 Test Accuracy: 0.707874682472481

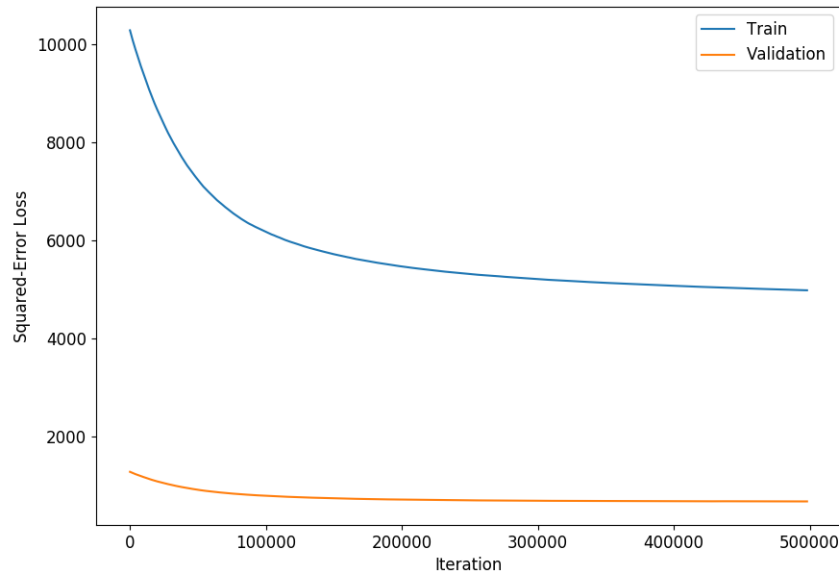
(d) The five questions we've selected are 3, 29, 107, 593, and 1771. The plot of probability of the correct response as a function of  $\theta$  given a question is as follows.



Each line represents the probability of correctness for a certain question as a function of the ability of students. We can see that the shape of lines are very similar to the sigmoid function, the probability of solving the same questions correctly are increasing when student ability increases. Also, for most of the students, the orange line is on the top, and the purple line is on the bottom. This implies that question 29 is the easiest and question 1771 is the hardest among the five questions.

### 3 Q3 - Option (i)(Xinyu Kang)

- (a) We tried out all  $k$ 's from the list  $[5, 6, 7, 8, 9, 10, 11, 12]$ , and found that  $k = 9$  is the best  $k$  with validation accuracy 0.6613039796782387 and test accuracy 0.6587637595258256.
- (b) We filled the missing entries with mean values, which would introduce noises to the dataset. It ignores the correlation between features, so the imputed values might not make sense. Also, adding mean values might reduce variance of the data and change what the top  $k$  eigenvalues should be.
- (c) See `matrix_factorization.py`
- (d) We have chosen learning rate to be 0.01 with 500000 iterations. We tried  $k$  values 10, 30, 45, 50, 55, 60, and found that  $k = 50$  is the best  $k$  value.
- (e) The plot of training and validation squared-error losses as a function of iteration is as follows.



The final validation accuracy is 0.7050522156364663 and the test accuracy is 0.7058989556872707 (for the latest run), but the accuracy varies for each run because of the randomness of `update_u_z`.

- (f) If we were to train the model as a binary classification problem, we could make the targets be from  $\{-1, 1\}$ . Suppose the target is  $t$  and our output is  $y$ , we can calculate the loss as

$$\mathcal{L}(y) = 1 - t \cdot y$$

So if  $t$  and  $y$  have the same sign, the loss is 0; otherwise, the loss is 1.

## 4 Q4(Xinyu Kang)

### 4.1 Ensemble Process

We use three ALS algorithms as our base models. For bootstrapping, we draw data points randomly from the training set with replacement, and our samples have exactly the size of the training set. After we run ALS three times on three different samples, we average them first before making binary predictions. That is, for validation or testing, we average the float possibilities from entries of the three matrices and then compare it to the 0.5 threshold to make the final predictions.

### 4.2 Results

Validation Accuracy: 0.7027942421676545

Test Accuracy: 0.7011007620660458

The accuracy varies (within 0.003) for each run, but on average the bagging ensemble does not have a better performance than the base model. We think this is because bagging can only reduce variance, but ALS does not overfit much to the training data. Also, since we are bootstrapping for only three samples, it is likely that none of the base models have seen some of the data at all. Repeated data in the samples also does not provide any advantages for ALS, because with 500000 iterations, any seen data might be trained for enough times.

## 5 Part B

In this part, We modified one-parameter IRT implemented in Part A. Since from IRT in Part A we could see there is no big difference between validation accuracy and test accuracy, indicating that we are not overfitting. So here we are trying to improve the model by adding more features to our model to increase the test accuracy. Our algorithm extension include two parts:

- (1) adding discrimination parameter  $\alpha_j$  for each question  $j$
- (2) initialize  $\theta$  by gender group.

### 5.1 Formal Description(Xin Peng, Xinyan He)

#### (1) Adding discrimination parameter $\alpha_j$

Problem in Part A:

In Part A, we have used one-parameter IRT model with  $\beta_j$  representing the difficulty of the question  $j$ , and  $\theta_i$  representing the  $i$ -th students ability.  $\beta_j$  is determined at the point of median point, and the question discrimination, i.e., the rate at which the probability of endorsing a correct question changes given ability levels, is fixed for all questions. However, with same discrimination for all questions, we cannot detect subtle differences in the ability of the respondents, which is not ideal. Since for some questions, they have a higher discrimination, and students with a higher ability can do much better than students with a lower ability. While for other questions, the discrimination is small, and students with different ability will perform similarly on these questions. For example, in figure 1, we have fitted the 1-parameter model using Part A Q2, and plot  $p(c_i = 1|\theta_i) = \frac{\exp(\theta_i - \beta)}{1 + \exp(\theta_i - \beta)}$  using  $\beta$  for each question. And each point in the figure represents the probability that the question is correctly answered given the student's ability using the data in train\_data.csv. From the figure, it is clear that the discrimination for question 1 is small, since when students' ability is small(-1), the probability that it can be answered correctly is around 0.8. However, even when the ability increases, the probability does not increase and it even decreases. While for Question 45, the discrimination would be quite big, since from the plot, it is very clear that probability of answering the question correctly is almost 1 for ability greater than 0, and almost 0 when ability is lower than -1. But since we are using 1-parameter model, the discrimination is the same for all questions, so from Figure 1, curves for these 3 different questions have the same shape, and the only difference is that the position of the median point is not different, and it does not fit the data very well.

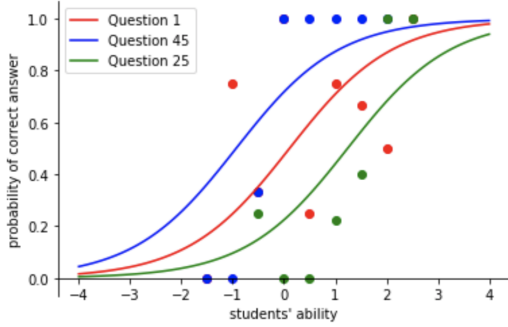


Figure 1: 1-Parameter Model

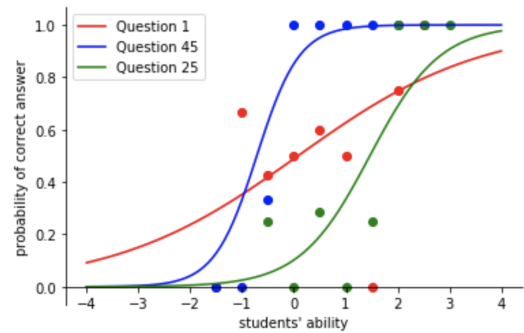


Figure 2: 2-Parameter Model

Extension:

So we decided to introduce discrimination  $\alpha_i$  for each question  $i$  in our new model. So different questions will have different slopes, and the steeper the slope, the higher the discrimination, and this question can detect subtle differences in respondents' ability. And our new model uses  $p(c_i = 1|\theta_i) = \frac{\exp(\alpha_i(\theta_i - \beta))}{1 + \exp(\alpha_i(\theta_i - \beta))}$ . And this new model can fit the data better. From Figure 2, using 2-parameter model, the curve for Question 1 is relatively flat, since as we have argued above that Question 1 cannot discriminate the probability of answering this question based on students' ability very well. And slope for Question 45 is very large, since Question 45 can greatly detect difference in probability based on students' ability. In comparison to Figure 1, the model in Figure 2 fits the data much better.

The way we implemented the 2-parameter IRT model is similar to what we have done in Part A, i.e., performing alternating gradient descent on  $\theta$ ,  $\beta$ ,  $\alpha$  to maximize the log-likelihood. And the comparison between 1-parameter model and 2-parameter model is listed below.

	1-parameter IRT	2-parameter IRT
probability that question $j$ is correctly answered	$p(c_{ij} = 1 \theta_i, \beta_j) = \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}$	$p(c_{ij} = 1 \theta_i, \beta_j, \alpha_j) = \frac{\exp[\alpha_j(\theta_i - \beta_j)]}{1 + \exp[\alpha_j(\theta_i - \beta_j)]}$
log-likelihood for all students and questions	$\log(\mathbf{C} \theta, \beta) = \sum_{i=1}^{542} \sum_{j=1}^{1774} [c_{ij} \log(\frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}) + (1 - c_{ij}) \log(1 - \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)})]$	$\log(\mathbf{C} \theta, \beta, \alpha) = \sum_{i=1}^{542} \sum_{j=1}^{1774} [c_{ij} \log(\frac{\exp(\alpha_j \theta_i - \alpha_j \beta_j)}{1 + \exp(\alpha_j \theta_i - \alpha_j \beta_j)}) + (1 - c_{ij}) \log(1 - \frac{\exp(\alpha_j \theta_i - \alpha_j \beta_j)}{1 + \exp(\alpha_j \theta_i - \alpha_j \beta_j)})]$
derivative with respect to $\theta_i$	$\sum_{j=1}^{1774} (c_{ij} - \frac{e^{\theta_i - \beta_j}}{1 + e^{\theta_i - \beta_j}})$	$\sum_{j=1}^{1774} (\alpha_j c_{ij} - \alpha_j \frac{e^{\alpha_j \theta_i - \alpha_j \beta_j}}{1 + e^{\alpha_j \theta_i - \alpha_j \beta_j}})$
derivative with respect to $\beta_j$	$\sum_{i=1}^{542} (-c_{ij} + \frac{e^{\theta_i - \beta_j}}{1 + e^{\theta_i - \beta_j}})$	$\sum_{i=1}^{542} (-\alpha_j c_{ij} + \alpha_j \frac{e^{\alpha_j \theta_i - \alpha_j \beta_j}}{1 + e^{\alpha_j \theta_i - \alpha_j \beta_j}})$
derivative with respect to $\alpha_j$		$\sum_{i=1}^{542} (c_{ij} \theta_i - c_{ij} \beta_j - (\theta_i - \beta_j) \frac{e^{\alpha_j \theta_i - \alpha_j \beta_j}}{1 + e^{\alpha_j \theta_i - \alpha_j \beta_j}})$

## (2) Initialize theta by gender group

Problem in Part A: In part A, we initialize every student ability( $\theta$ ) to zero which we are treating every student's ability equally at first. Our intuition here is that different gender of student may lead to different ability, like some people are tend to believe boys may perform better than girls on STEM subjects. So here we are trying to initialize the  $\theta$  according to student's gender. The basic process is to first separate the training data into 3 sets according to gender, and train each set to gain the approximation of that gender's  $\theta$ . Finally training the model using the whole training data and initialize the theta using the approximation we calculated for each gender.

### Algorithm box:

**Step 1:** improve the irt model by adding a new parameter.

```

1 def update_theta_beta_alpha(data, lr, theta, beta, alpha):
2     theta = theta - lr * gradient of theta
3     beta = beta - lr * gradient of beta
4     alpha = alpha - lr * gradient of alpha
5 return theta, beta, alpha

```

**Step 2:** split the training data according to the user's gender

```

1 def split_data_gender(training_data):
2     # we are going to split data according to their gender
3     for data in training_data:
4         gender = find_gender(data["user_id"])
5         if gender is 0:
6             put data in gender0
7         else if gender is 1:
8             put data in gender1
9         else:
10            put data in gender2
11 return gender0, gender1, gender2

```

**Step 3:** train each gender data using the improved irt model and get the corresponding theta for each gender group.

```

1 def theta_initializer(gender0, gender1, gender2):
2     # this list stores the respective initializer for each gender group
3     theta_initializer = []
4     # we are going to split data according to their gender
5     for i in 0, 1, 2:
6         theta_list = improved_irt(gender_i)
7         add average value of theta_list to theta_initializer
8 return theta_initializer

```

**Step 4:** This is our **final model**, which created based on improved irt model while here each theta is initialized according to their gender.

```

1 def final_irt(train_data, theta_initializer):
2     # initialize the theta according to gender
3     for data in train_data:
4         # identify gender of the user
5         gender = find_gender(data["user_id"])
6         if gender is 0:

```



```

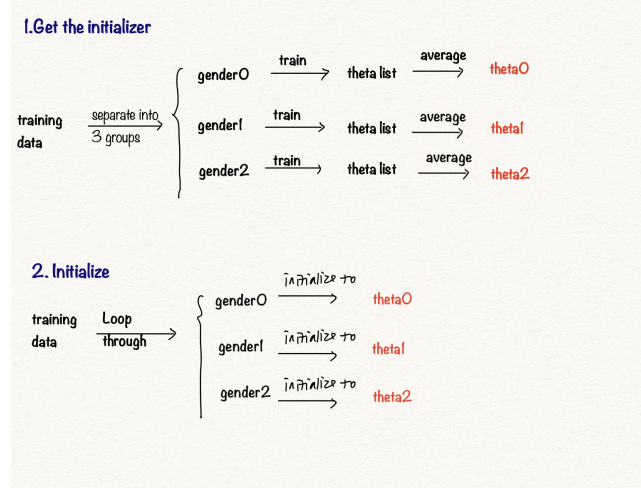
7         initialize it to theta_initializer[0]
8     else if gender is 1:
9         initialize it to theta_initializer[1]
10    else:
11        initialize it to theta_initializer[2]
12    #other things remain the same as the improved_irt

```

---

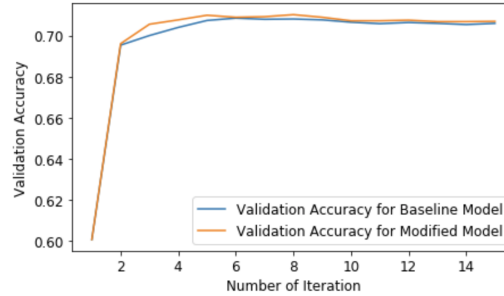
## 5.2 Diagram(Xinyan He)

Simple model for initialize theta by gender group.



Please refer to partb.py and the algorithm box above for the detailed implementation.

## 5.3 Comparison(Xin Peng)



	Baseline Model	Modified Model
Final Test Accuracy	0.7058989556872707	0.7106971493084956
Final Test Accuracy on Gender0(Unspecified)	0.6836734693877551	0.6938775510204082
Final Test Accuracy on Gender1(Female)	0.7073007367716009	0.7079705291359679
Final Test Accuracy on Gender2(Male)	0.7155425219941349	0.7221407624633431

The test accuracy using the modified model is 0.7106971493084956, and there is a 0.4% improvement on overall test accuracy using the modified model.

The overall purpose of using the modified model is to

1. Improve overall performance by using different discrimination for different questions.
2. Improve performance for unspecified, female, male group by using different initializations for each group.

From the table above, it is obvious that the overall test accuracy has improved.

And to further see whether accuracy on each group has improved, we further calculate the test accuracy for each group separately. Since test accuracy on each group has all increased, our model has indeed improved performance on overall model, as well as on each group, which is consistent with our hypothesis.

## 5.4 Limitations and Extension(Xin Peng, Xinyan He)

- In 2-parameter IRT, we have not taken group variables such as gender, question subject into account, so if the  $\alpha$ ,  $\beta$  are quite different for each group, our model would not perform well. Although we have tried

to separate genders into three groups: 0(unspecified), 1(female), 2(male), and fit three different  $\beta$ ,  $\alpha$  for each group, it does not work well using the given data.

One possible reason is that the training data is not large enough, since if we separate training data into three groups, there are some cases where no student in a group have answered a specific question, thus the  $\alpha$ ,  $\beta$  for this question and group will remain unchanged from initialization.

A possible extension is to use more data to train the model or to perform a leave-one-out cross validation to solve the problem of lacking data, and then to fit different  $\alpha$ ,  $\beta$ ,  $\theta$  for each group.

- We didn't make use of the question\_meta to improve our algorithm. So our algorithm may approach poorly when a student ability is quite unbalanced. For example a student is very good at solving algebra questions but are very bad at solving problems related to graph. If majority training data we get for this student is about the algebra, we might think the student a high theta. So when the test data is about graph, we may probably make a prediction that it would answer the question correctly. However in this case we would probably make a mistake.

A possible extension would be make several groups of questions according the subjects they are belong to. Then get a approximation of parameters  $\beta$  and  $\alpha$  for each group and make a different initialization during the training process.