

Part B

formal description

The performance of the IRT algorithm in part A is not satisfactory. We analyzed that the main reason is that a single decision tree is a high-variance model. We believe that the decision tree in part A overfits the training data, which makes the model's generalization ability poor.

Therefore, we decided to reduce the variance by averaging the predictions of multiple decision trees using a random forest. Each tree is trained on a different random subset of the training data. The results are aggregated at the end to make the model more robust by smoothing the data.

Algorithm Box

```

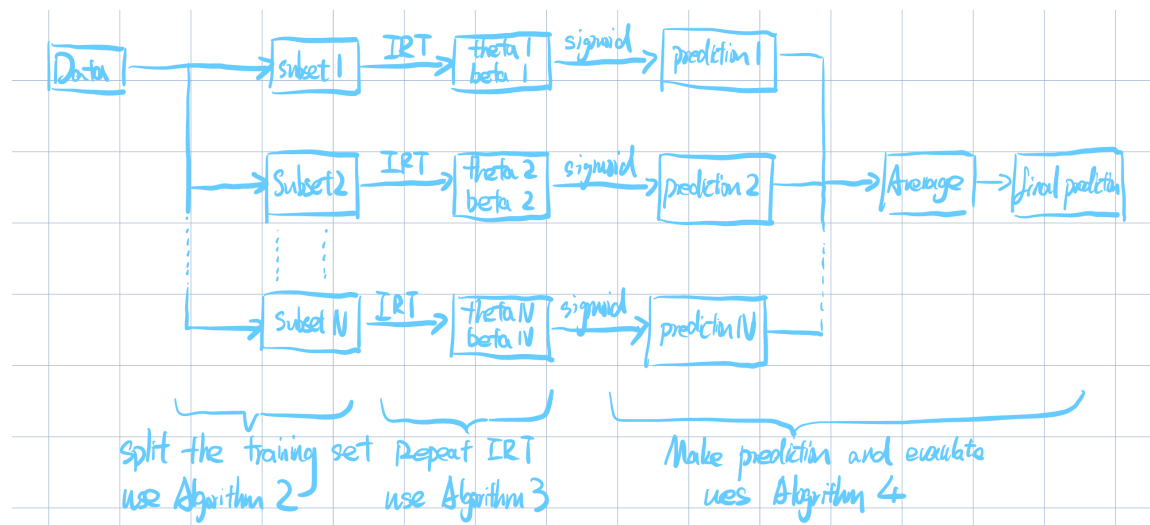
1           Algorithm 1: regular irt algorithm
2   Result: train a single irt model
3   Output: thetas, betas
4
5   # regular irt algorithm, no changes needed
6
7
8           Algorithm 2: split the training data
9   Result: split the training data into subsets
10  Output: a subset of the training data
11
12  # randomly split the training data into subsets
13
14
15           Algorithm 3: irt ensemble algorithm
16  Result: train an ensemble of irt models
17  Output: list of thetas, list of betas
18
19  theta_lst = []
20  beta_lst = []
21  int repeat
22  for i in range(repeat):
23      # split the training data
24      train_data = split_data(data)
25      # train a single irt model
26      thetas, betas = irt(train_data)
27      # store the thetas and betas in a list
28      theta_lst.append(thetas)
29      beta_lst.append(betas)
30
31  return theta_lst, beta_lst
32
33
34           Algorithm 4: evaluate the ensemble
35  Result: evaluate the ensemble of irt models
36  Output: the accuracy of the ensemble model
37
38  total_correct = 0
39  for i in range(len(data["user_id"])):
40      predict_lst = []
41      for j in range(len(theta_lst)):
42          # make prediction for each model and store it in a list
43          # take the average of the predictions
44
45      # then calculate the accuracy
46  return total_correct / len(data["is_correct"])
47

```

By training multiple decision trees on different data subsets in the form of random forests,

the overfitting problem of the model can be reduced.

Idea Diagram



Comparison or Demonstration

For comparison, when using the single irt algorithm, we obtained the following statistics. We use this group of data as the baseline models:

```
Final Validation Accuracy: 0.7063223257126728
Final Test Accuracy: 0.707310189105278
```

When using a random forest consisting of two decision trees, we get the following statistics:

```
Ensemble Validation Accuracy: 0.7022297488004516
Ensemble Test Accuracy: 0.7044877222692634
```

When using a random forest consisting of three decision trees, we get the following statistics:

```
Ensemble Validation Accuracy: 0.7054755856618685
Ensemble Test Accuracy: 0.7044877222692634
```

When using a random forest consisting of four decision trees, we get the following statistics:

```
Ensemble Validation Accuracy: 0.705193338978267
Ensemble Test Accuracy: 0.703076488851256
```

After comparison, our model does not significantly improve the accuracy.

experiment to test our hypothesis

We use the accuracy of the training data set and the accuracy of the validation data set to determine whether the model has signs of overfitting. If the training accuracy is significantly higher than the validation accuracy, it means that the model is too sensitive to the training data set and has signs of overfitting. On the contrary, if the training accuracy and validation accuracy are close, it means that the model does not have overfitting.

The test accuracy of the original model is as follows:

```
Final train Accuracy: 0.7398744002257973
```

Obviously, the training accuracy of the original model is significantly higher than the validation accuracy, indicating that the original model may be overfitting.

The random forests consisting of 2, 3, and 4 decision trees have the following statistics:

```
Ensemble Training Accuracy: 0.7379339542760373
```

```
Ensemble Training Accuracy: 0.7391158622636184
```

```
Ensemble Training Accuracy: 0.7374576629974597
```

Unfortunately, Random Forest only slightly reduces the gap between training accuracy and validation accuracy, and the overfitting problem still exists.

Limitations

As mentioned above, implementing Random Forest algorithm did not improve our original model significantly. This could be due to the following reasons:

The given dataset is not large enough or diverse enough. Random Forest is known to perform well on large-scale datasets with many features. In our improved model, we split the dataset into many subsets and trained a Random Forest model on each subset. However, the dataset may not be large enough to benefit from this approach.

Hence, changing the dataset may improve the performance of the Random Forest algorithm.

In our bagging implementation, we assume all predictions from the base models are equally important. This may not be the case in practice. Some models may perform better on certain types of data, and simply averaging their predictions may not be the best approach.

We can use a weighted average of the predictions to give more importance to the better performing models, but this requires tuning the weights, which can be time-consuming.

The Random Forest algorithm is computationally expensive. It requires a lot of time to tune the hyperparameters and train the model. We may not have been able to find the best hyperparameters in the time we had.

IRT model itself may be limited in predicting the student's performance. If we integrate more complex models in the ensemble, we may improve the performance of the model.