

FedIRT: An R package and shiny app for estimating federated item response theory models

Biying Zhou¹ and Feng Ji¹

¹ Department of Applied Psychology & Human Development, University of Toronto, Toronto, Canada
Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

We developed an R package FedIRT, to estimate traditional IRT models, including 2PL and the graded response models with additional privacy, allowing parameter estimation in a distributed manner without compromising estimation accuracy. Numerical experiments demonstrate that Federated IRT estimation achieves comparable statistical performance to mainstream IRT packages in R, with the benefits of privacy preservation and minimal communication costs. The R package also includes a user-friendly Shiny app that allows clients (e.g., individual schools) and servers (e.g., school boards) to apply our proposed method in a user-friendly manner.

Statement of Need

IRT ([Embretson & Reise, 2013](#)) is a statistical modeling framework grounded in modern test theory, frequently used in the educational, social, and behavioral sciences to measure latent constructs through multivariate human responses. Traditional IRT estimation mandates the centralization of all individual raw response data in one location, thereby potentially compromising the privacy of the data and participants ([Lemons, 2014](#)).

Federated learning has emerged as a field addressing data privacy issues and techniques for parameter estimation in a decentralized, distributed manner. However, there is currently no package available in psychometrics, especially in the context of IRT, that integrates federated learning with IRT model estimation.

Mainstream IRT packages in R, such as `mirt` ([Chalmers, 2012](#)) and `ltm` ([Rizopoulos, 2007](#)) require storing and computing all data in a single location, which can potentially lead to violations of privacy policies when dealing with highly sensitive data (e.g., high-stakes student assessments).

We have therefore developed a specialized R package, FedIRT, to integrate federated learning with IRT. We have also developed an accompanying Shiny app to recognize real-world challenges and aim to reduce the burden of learning R programming for applying this package. This app implements the method in a user-friendly and accessible manner.

Method

Here we briefly introduce the key idea behind integrating federated learning with IRT. For details, please refer to our methodological discussions on Federated IRT ([Zhou & Ji, 2023, 2024, In submission](#)).

37 Model formulation

38 The two-parameter logistic (2PL) IRT model is often considered the most popular IRT model.
 39 In 2PL, the response by person i for item j is often binary: $X_{ij} \in \{0, 1\}$, and the probability
 40 of person i answering item j with discrimination α_j and difficulty β_j correctly:

$$P(X_{ij} = 1|\theta_i) = \frac{e^{\alpha_j(\theta_i - \beta_j)}}{1 + e^{\alpha_j(\theta_i - \beta_j)}}$$

41 To make our package available for polytomous response, we also developed a federated learning
 42 estimation algorithm for the Generalized Partial Credit Model (GPCM) in which the probability
 43 of a person with the ability θ_i obtaining x scores in item j is:

$$P^{\text{GPCM}}(X_{ij} = x|\theta_i) = \frac{e^{\sum_{h=1}^x \alpha_j(\theta_i - \beta_{jh})}}{\sum_{c=0}^m e^{\sum_{h=1}^c \alpha_j(\theta_i - \beta_{jh})}}$$

44 In this function, β_{jh} is the difficulty of scoring level h for item j , and for each item j , all
 45 difficulty levels have the same discrimination α_j . m_j is the maximum score of item j .

46 Model estimation

47 In both 2PL and GPCM, often we assume the ability follows a standard normal distribution,
 48 thus we can apply MMLE.

49 We use a combination of traditional MMLE with federated average (FedAvg) and federated
 50 stochastic gradient descent (FedSGD) (McMahan et al., 2017). In our case, the log-likelihood
 51 and partial gradients are sent from the clients to the server. Then, the server uses FedSGD to
 52 update the item parameters and send them back to clients. By iterations, the model converges
 53 and displays the estimates on the interface.

54 Taking the 2PL model as an example, which has a marginal log-likelihood function l for each
 55 school k that can be approximated using Gaussian-Hermite quadrature with q (by default,
 56 $q = 21$) equally-spaced levels, and let $V(n)$ to be the ability value of level n , and $A(n)$ is the
 57 weight of level n .

$$l_k \approx \sum_{i=1}^{N_k} \sum_{j=1}^J X_{ijk} \times \log\left[\sum_{n=1}^q P_j(V(n))A(n)\right] + (1 - X_{ijk}) \times \log\left[\sum_{n=1}^q Q_j(V(n))A(n)\right]$$

58 By applying FedAvg, the server collects the log-likelihood values from all k schools and then
 59 sums up all the likelihood values to get the overall log-likelihood value: $l = \sum_{k=1}^K l_k$.

60 The server collects a log-likelihood value l_k and all derivatives $\frac{l_k}{\partial \alpha_j}$ and $\frac{l_k}{\partial \beta_j}$ from all clients,
 61 then observe that $\frac{\partial l}{\partial \alpha_j} = \sum_{k=1}^K \frac{l_k}{\partial \alpha_j}$ and $\frac{\partial l}{\partial \beta_j} = \sum_{k=1}^K \frac{l_k}{\partial \beta_j}$ by FedSGD, the server sums up all
 62 log-likelihood values and derivative values.

63 Also, we provided an alternative solution, Federated Median, which uses the median of the
 64 likelihood values to replace the sum of likelihood values in Fed-MLE (Liu et al., 2020). It is
 65 more robust when there are outliers in input data.

66 With estimates of α_j and β_j in 2PL or β_{jh} in GPCM, empirical Bayesian estimates of students'
 67 ability can be obtained (Bock & Aitkin, 1981).

68 Comparison with existing packages

69 We showcase that our package could generate the same result as traditional IRT packages,
70 for example, mirt (Chalmers, 2012). Take 2PL as an example, we use a synthesized dataset
71 with 160 students and 10 items. %For traditional packages, the whole dataset is used. For our
72 package, the dataset was separated into two parts, which contain 81 and 79 students.

73 Figure 1 and Figure 2 show the comparison of the discrimination and difficulty parameters
74 between mirt and FedIRT based on example_data_2PL in our package.

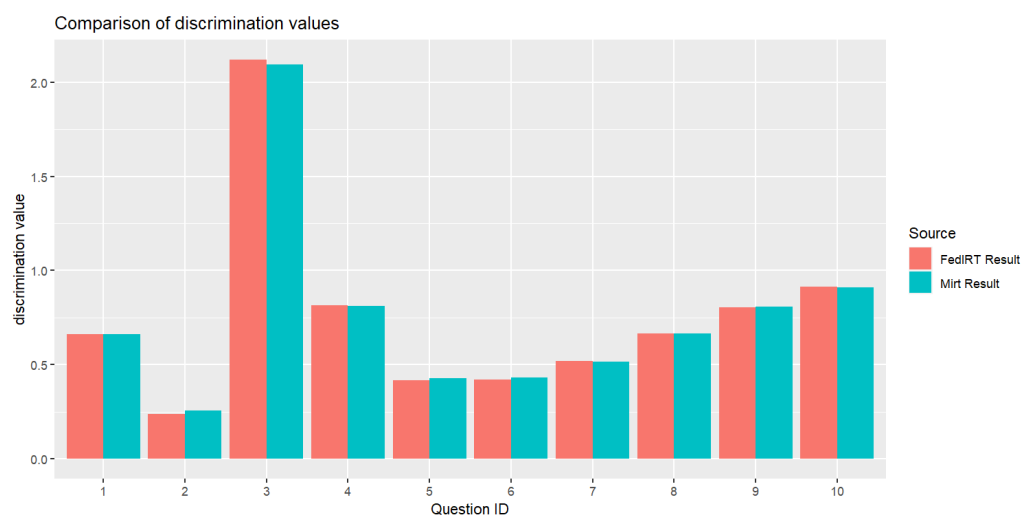


Figure 1: Discrimination parameter estimates comparison

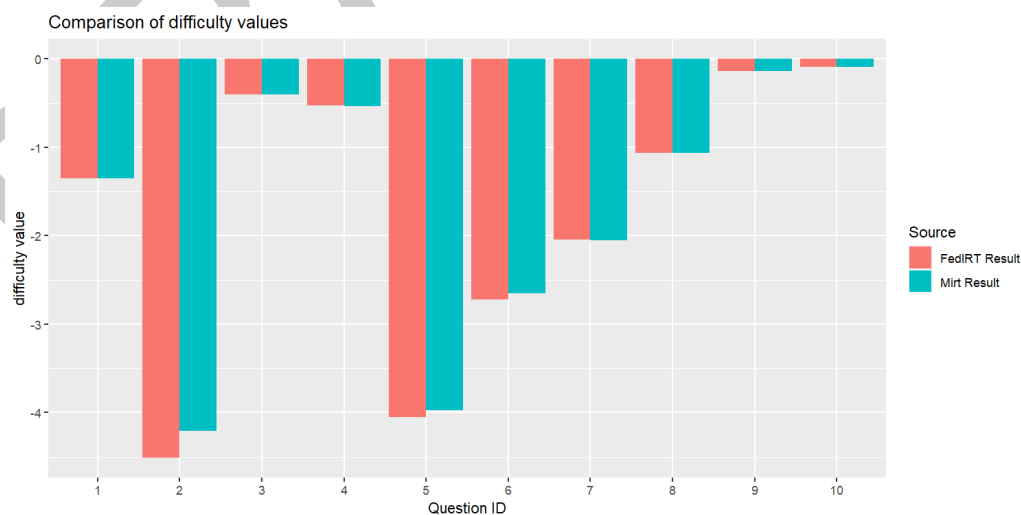


Figure 2: Difficulty parameter estimates comparison

75 Availability

76 The R package FedIRT is publicly available on [Github](#). It could be installed and run by using
77 the following commands:

```
devtools::install_github("Feng-Ji-Lab/FedIRT")
library(FedIRT)
```

Sample of the integrated function

We provide a function `fedirt` in the package, and the detailed usage of the function is shown in the user manual. We demonstrate a sample here.

Suppose we have a dataset called `dataset.csv`, and the head of this dataset is shown below:

site	X1	X2	X3	X4	X5
10	1	0	0	0	0
7	0	0	1	0	0
9	0	0	1	1	1
1	1	0	1	1	1
2	1	0	0	0	0

First, we need to split the dataset by different sites. The index of each site is indicated in the column `site`.

```
# split the dataset by sites
data <- read.csv("dataset.csv", header = TRUE)
data_list <- split(data[, -1], data$site)
```

Then, we change every sites' data into a list of matrices in R.

```
# change the list into matrices
inputdata <- lapply(data_list, as.matrix)
```

Then, we call the function `FedIRT::fedirt()` to get the result. It returns a list of item discriminations, item difficulties, and each sites' effect and each students' abilities. Call `summary()` to see the returned list.

```
# call the fedirt function
result <- fedirt(inputdata, model_name = "2PL")
```

At last, print the results or use the parameters for further analysis.

```
print(result$a)
print(result$b)
```

Apart from using the results for further analysis, we can also use `summary()` to generate a snapshot of the result. Here is a sample below.

```
summary(result)
```

Then, the result will be printed in the console as follows:

```
Summary of FedIRT Results:
```

```
Counts:
function gradient
      150      68
```

```
Convergence Status (convergence):
Converged
```

```
Log Likelihood (loglik):
```

```

103 [1] -957.1493
104
105 Difficulty Parameters (b):
106 [1] 0.1699265 -2.8088090 1.1167600 0.9893799 -2.5409030 -1.1789985 -
107 0.5258475 0.4560620 1.3792979
108 [10] 1.4247369
109
110 Discrimination Parameters (a):
111 [1] 0.6627037 0.2495989 2.1162137 0.8155786 0.4177167 0.4228183 0.5176585 0.6663483 0.8
112
113 School effect:
114 [1] 1.517844
115
116 Ability Estimates:
117 School 1:
118 [1] -1.187698446 -0.552434825 -0.899668043 -0.206272962 0.387992333 0.678869288 -
119 0.146111320
120 [8] 0.318665280 0.408994368 -0.139117072 0.496885125 0.562515435 0.392422237 0.5
121 In summary, we read a dataset and split it into different sites. Note that the dataset should
122 indicate different sites. Then call the function fedirt with corresponding arguments. At last,
123 print the results in need or use the part of results needed.

```

124 Sample of the personscore function

125 We provide a function personscore in the package. The detailed usage of the function is
126 shown in the user manual. We demonstrate a sample here.

```

personscoreResult = personscore(result)
summary(personscoreResult)

```

127 We called the function FedIRT::fedirt() to get the result. Then use personscore function
128 to get the person score result from result by personscore(result)

129 Summary of FedIRT Person Score Results:

```

130
131 Ability Estimates:
132 School 1:
133 [1] -1.184622030 -0.549302071 -0.902091156 -0.206457442 0.386666115 0.679138950
134 [7] -0.147088887 0.318980882 0.411915776 -0.140957745 0.496036332 0.562691489

```

135 Sample of the personfit function

136 We provide a function personfit in the package. The detailed usage of the function is shown
137 in the user manual. We demonstrate a sample here.

```

personfitResult = personfit(result)
summary(personfitResult)

```

138 We called the function FedIRT::fedirt() to get the result. Then use personfit function
139 to get the person score result from result by personfit(result)

140 Summary of FedIRT Person Fit Results:

```

141
142 Fit Estimates:
143 School 1:
144           Lz           Zh           Infit           Outfit
145 [1,] -0.925209243 -1.6737819927 0.23991579 0.24675914
146 [2,] -0.443958144 -0.7745951619 0.29205337 0.29671611

```

147 Sample of the Shiny App

148 To provide wider access for practitioners, we include the Shiny user interface in our package.
149 A detailed manual was provided in the package. Taking the 2PL as an example, we illustrate
150 how to use the Shiny app below.

151 In the first step, the server end (e.g., test administrator, school board) can be launched by running
152 the Shiny app (`runserver()`) with the interface shown below:

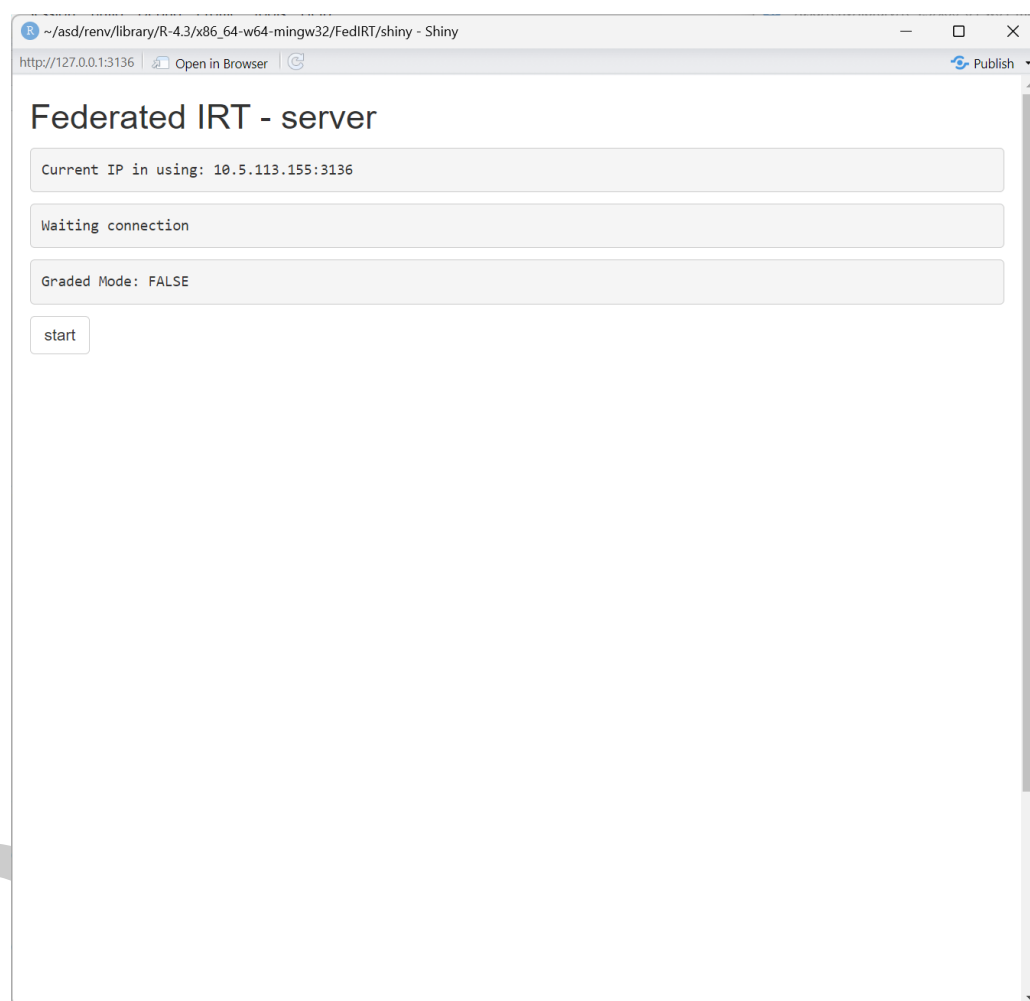


Figure 3: The initial server interface.

153 Then, the client-end Shiny app can be initialized (`runclient()`).

154 When the client first launches, it will automatically connect to the localhost port 8000 as
155 default.

156 If the server is deployed on another computer, type the server's IP address and port (which
157 will be displayed on the server's interface), then click "reconnect". The screenshots of the user
158 interface are shown below.

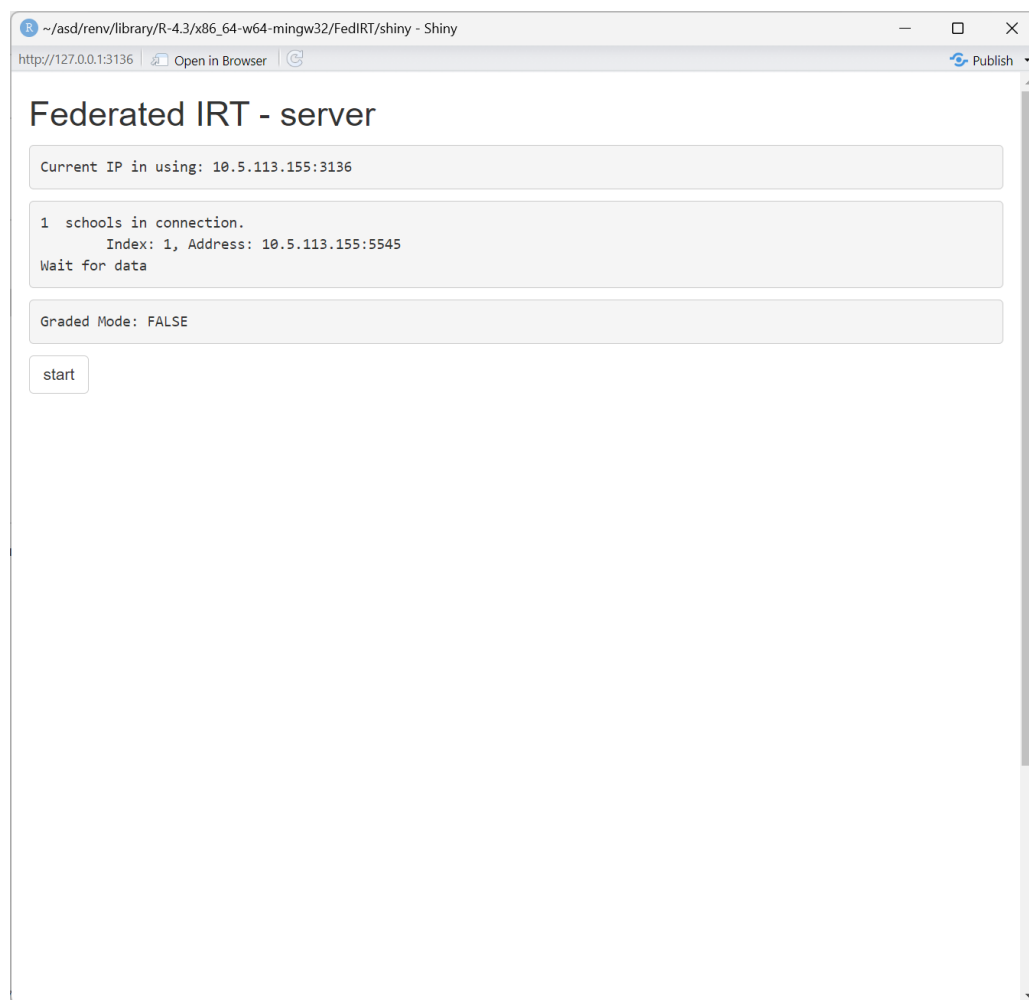


Figure 4: Server interface when one school is connected.

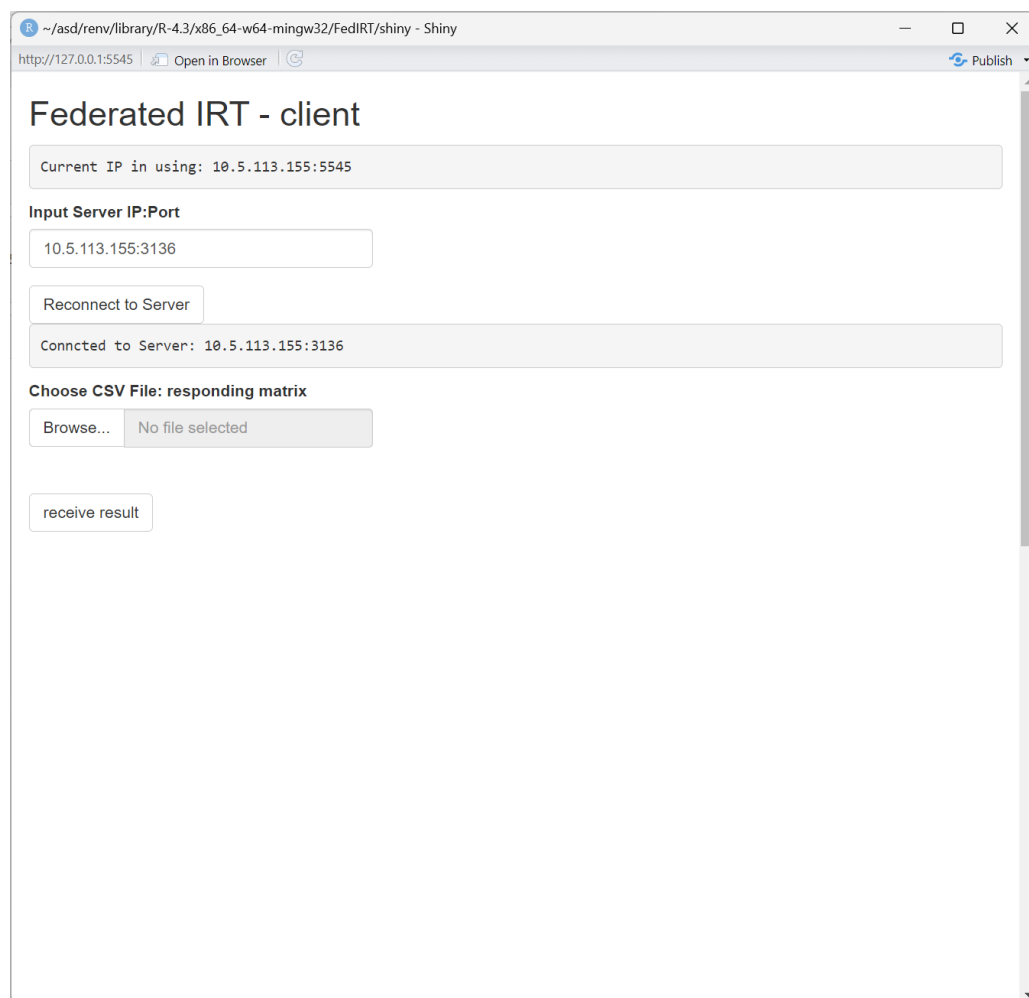


Figure 5: Client interface when connected to server.

159 Then, the client should choose a file to upload to the local Shiny app to do local calculations,
160 without sending it to the server. The file should be a csv file, with either binary or graded
161 response, and all clients should share the same number of items, and the same maximum
162 score in each item (if the answers are polytomous), otherwise, there will be an error message
163 suggesting to check the datasets of all clients.

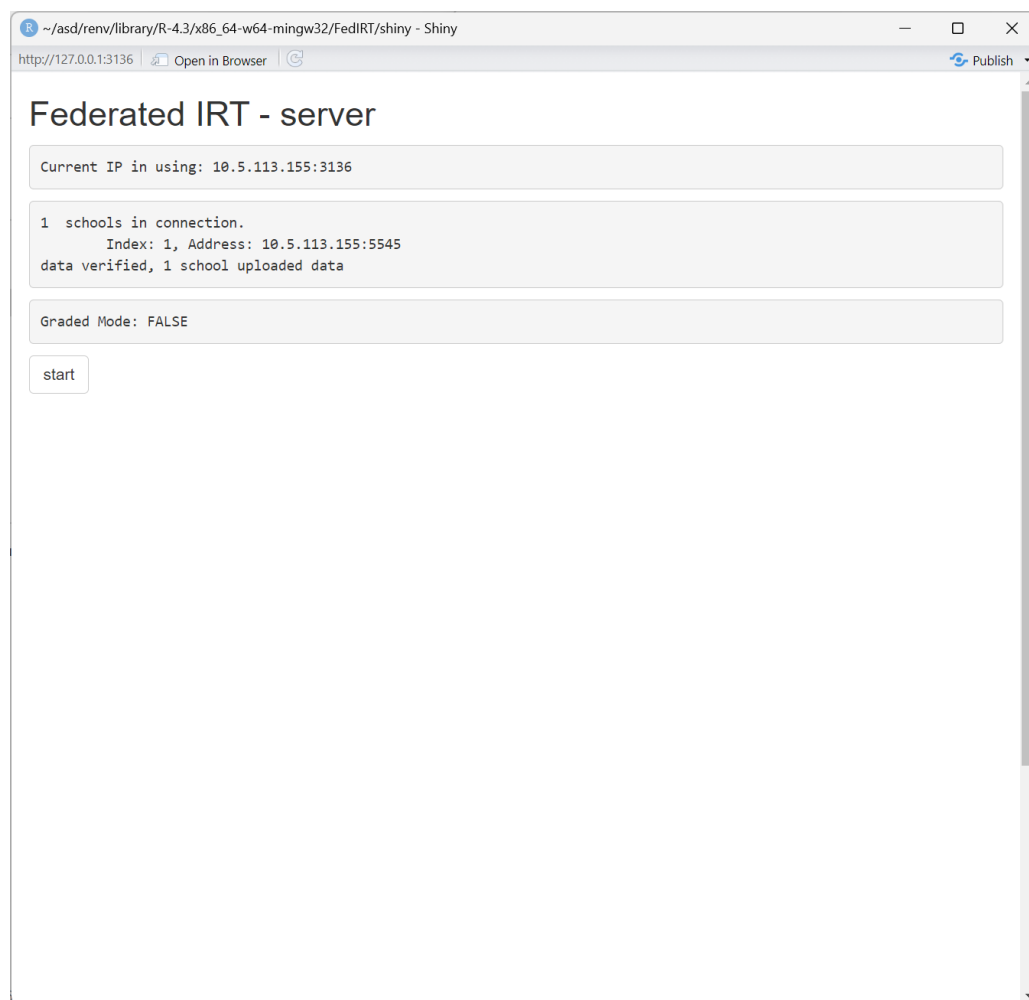


Figure 6: Server interface when one school uploaded dataset.

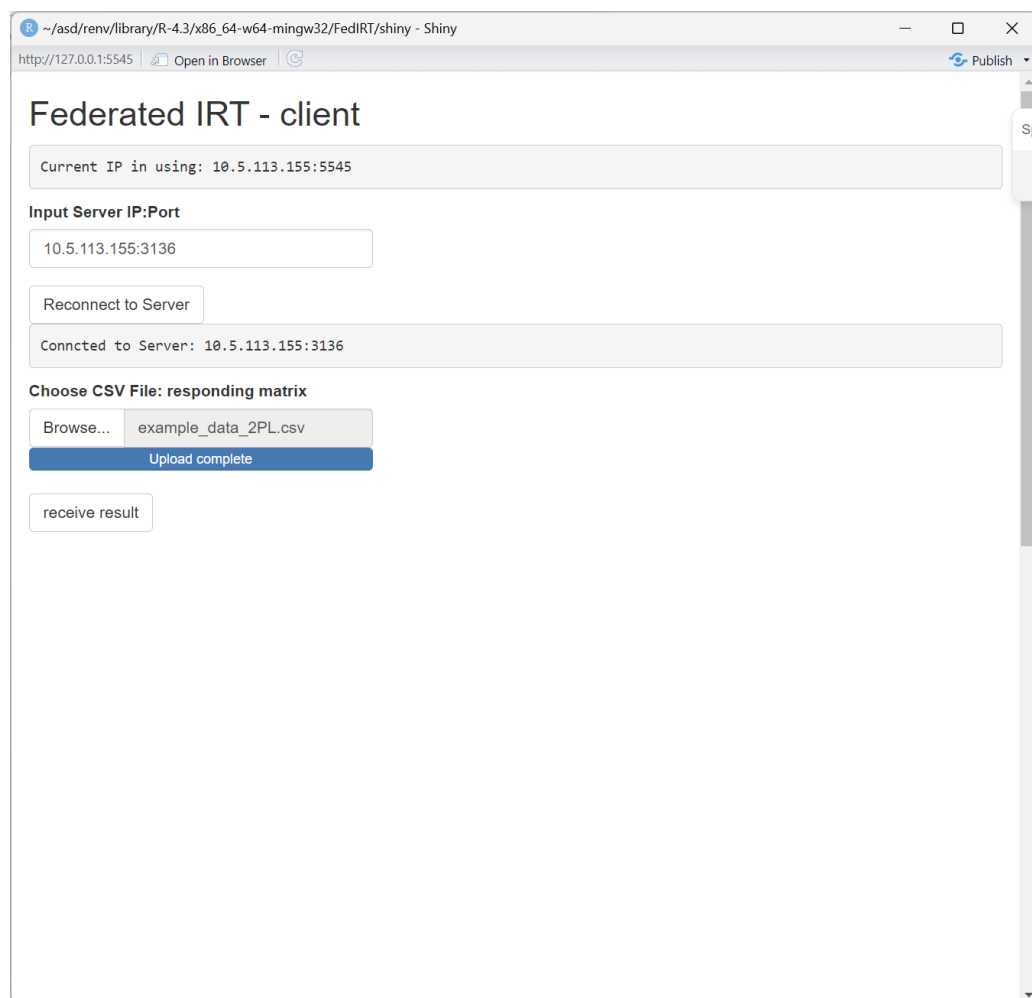


Figure 7: Client interface when a dataset is uploaded successfully.

164 After all the clients upload their data, the server should click “start” to begin the federated
165 estimates process and after the model converges, the client should click “receive result”. The
166 server will display all item parameters and the client will display all item parameters and
167 individual ability estimates.

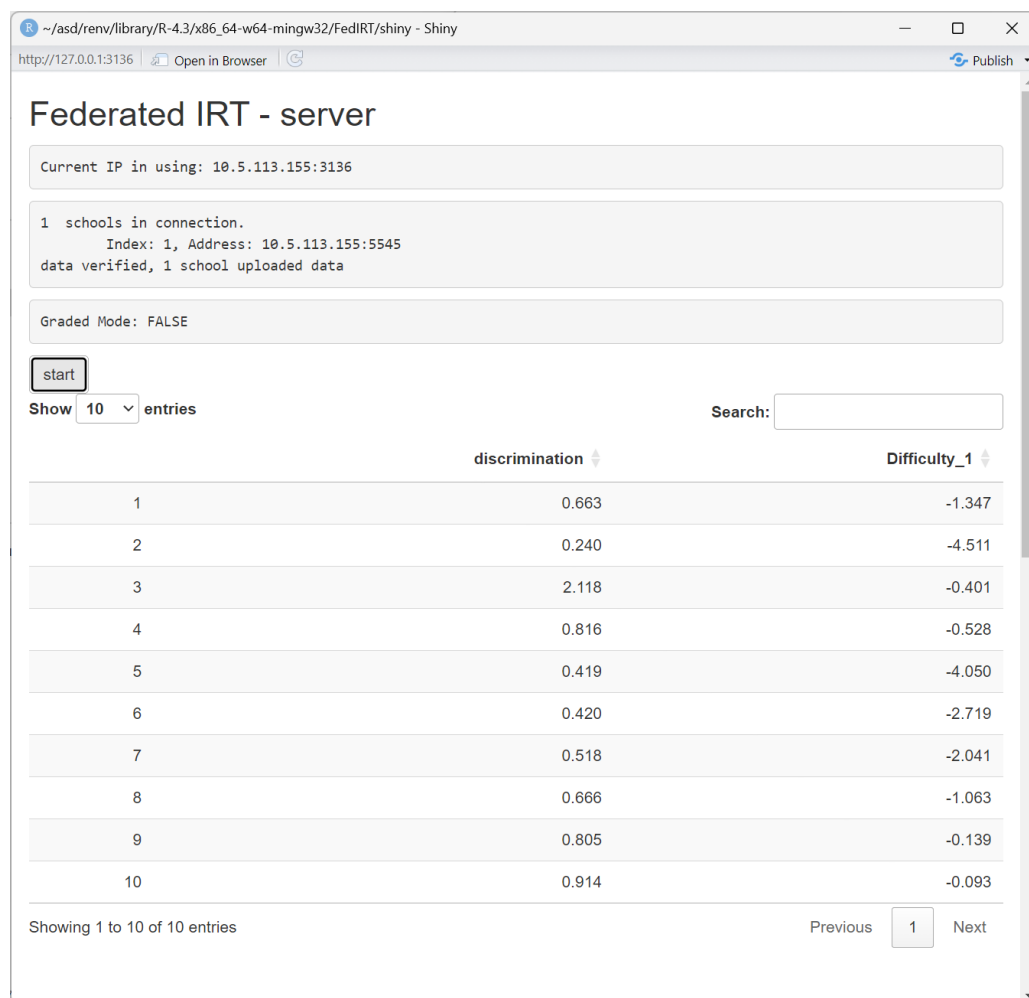


Figure 8: Server interface when estimation is completed.

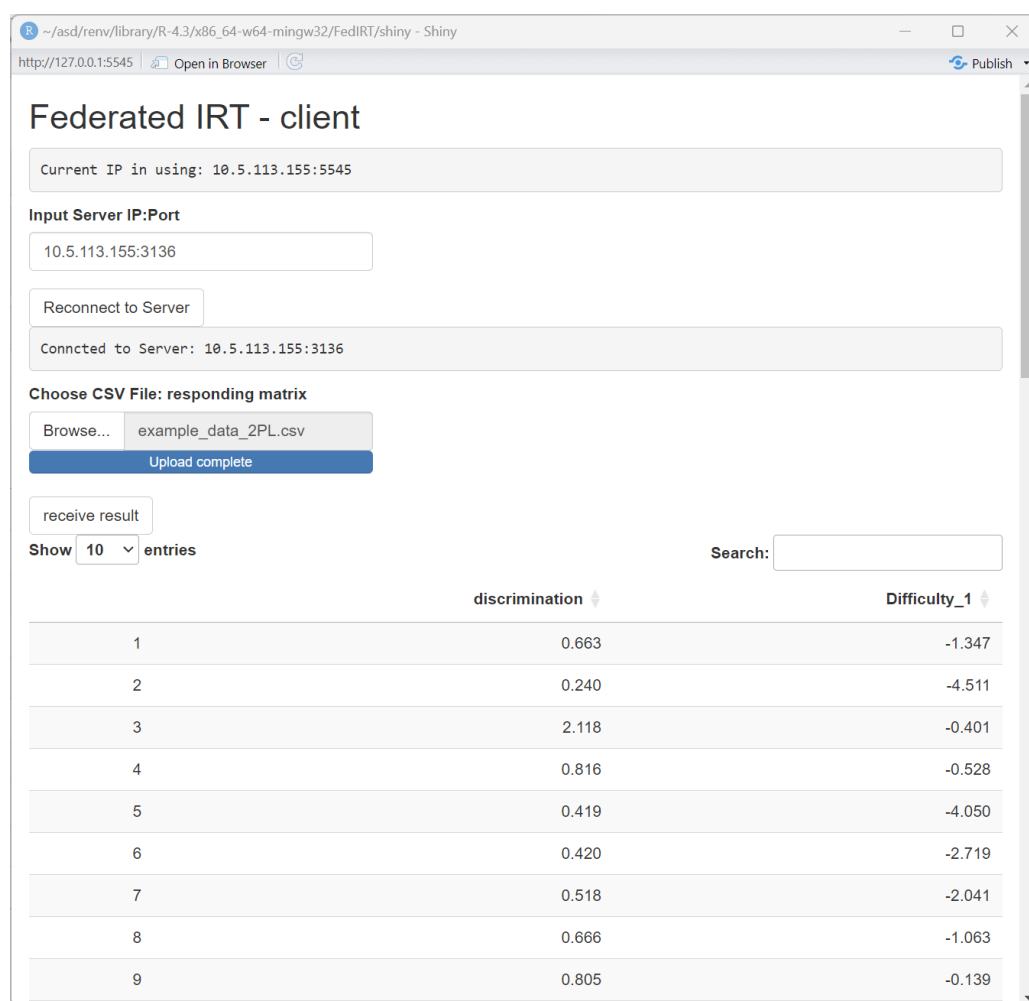


Figure 9: Client interface when the results received.

The clients will also display bar plots of the ability estimates.

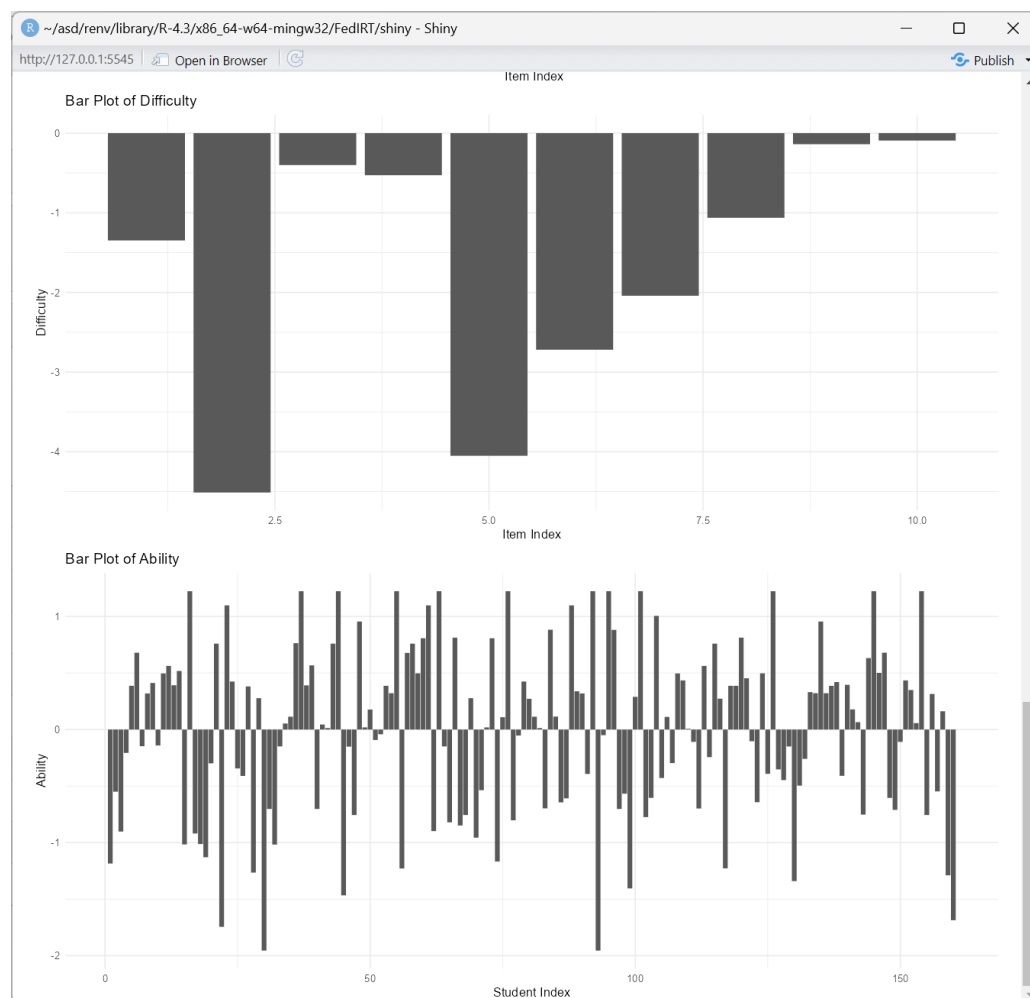


Figure 10: Client interface for displaying results.

References

- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: Application of an EM algorithm. *Psychometrika*, 46(4), 443–459. <https://doi.org/10.1007/BF02293801>
- Chalmers, R. P. (2012). Mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, 48, 1–29. <https://doi.org/10.18637/jss.v048.i06>
- Embretson, S. E., & Reise, S. P. (2013). *Item response theory*. Psychology Press. <https://doi.org/10.4324/9781410605269>
- Lemons, M. Q. (2014). *Predictive modeling of uniform differential item functioning preservation likelihoods after applying disclosure avoidance techniques to protect privacy* [PhD thesis]. Virginia Polytechnic Institute; State University.
- Liu, Y., Zhang, L., Ge, N., & Li, G. (2020). A systematic literature review on federated learning: From a model quality perspective. *arXiv Preprint arXiv:2012.01973*. <https://doi.org/10.48550/arXiv.2012.01973>
- McMahan, B., Moore, E., Ramage, D., Hampson, S., & Arcas, B. A. y. (2017). Communication-

- 185 efficient learning of deep networks from decentralized data. *Artificial Intelligence and*
186 *Statistics*, 1273–1282.
- 187 Rizopoulos, D. (2007). Ltm: An r package for latent variable modeling and item response
188 analysis. *Journal of Statistical Software*, 17, 1–25. <https://doi.org/10.18637/jss.v017.i05>
- 189 Zhou, B., & Ji, F. (2023). Federated psychometrics: A distributed, privacy-preserving, and
190 efficient IRT estimation algorithm. *APHD Research Gala*.
- 191 Zhou, B., & Ji, F. (2024). Federated item response theory: A distributed, privacy-preserving,
192 and efficient IRT estimation algorithm. *DPE MED Research Practicum Poster*.
- 193 Zhou, B., & Ji, F. (In submission). *Federated item response models*.

DRAFT