

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Published: unpublished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](#)).

Model formulation

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

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

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

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

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

Model estimation

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

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

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

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

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

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, 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 log-likelihood values and derivative values.

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

With estimates of α_j and β_j in 2PL or β_{jh} in GPCM, empirical Bayesian estimates of students' 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.

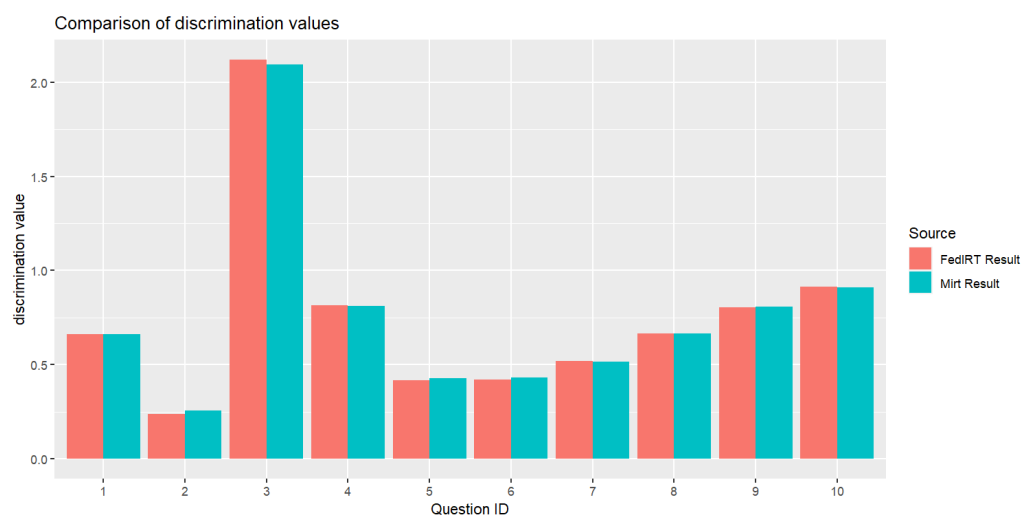


Figure 1: Discrimination parameter estimates comparison

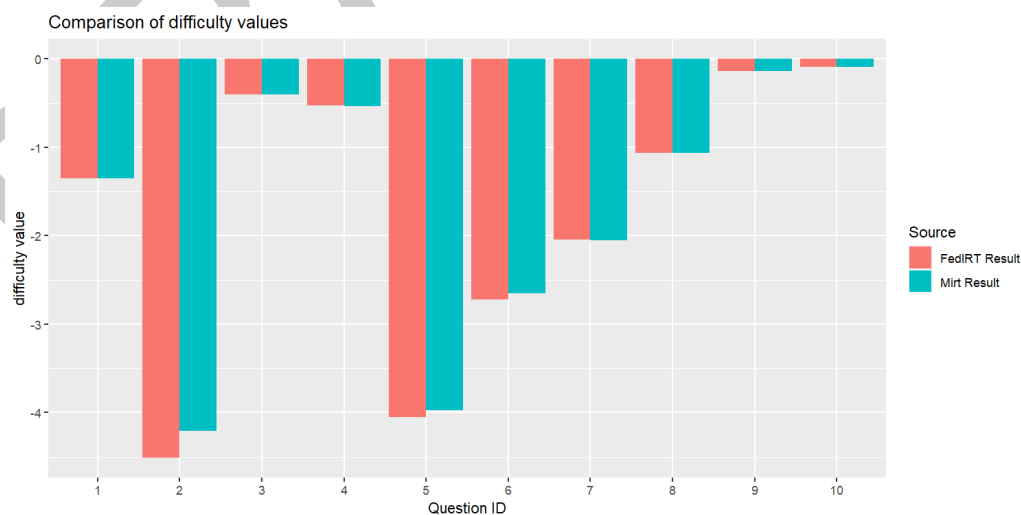


Figure 2: Difficulty parameter estimates comparison

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)
```

Example of the integrated function

We provide a function `fedirt_file()` in the package, and the detailed usage of the function is shown in the user manual. We demonstrate an example here.

Suppose we have a dataset called `dataset.csv`, and the head of this dataset is shown below. There should be one column indicating the school, for example, "site" here. Each other column indicates an item, and each row represents an answering status.

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 read the dataset.

```
# read dataset
data <- read.csv("dataset.csv", header = TRUE)
```

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

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

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

```
result$a
result$b
```

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

```
summary(result)
```

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

```
Summary of FedIRT Results:
```

```
Counts:
```

```
function gradient
      735      249
```

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

```
Converged
```

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

```
[1] -7068.258
```

```
Difficulty Parameters (b):
```

```
[1] -185.88151839    0.99524035    0.92927254    ...
```

```

107 Discrimination Parameters (a):
108 [1] 0.0028497700 0.8440140746 -0.1190176844 ...
109
110 Ability Estimates:
111 School 1:
112 [1] -1.127097195 -0.922572829 -0.993953038 ...
113 School 2:
114 [1] -1.41454573 1.78068772 1.87469389 ...
115 ...
116
117 End of Summary

```

118 Example of the personscore function

119 We provide a function personscore in the package. The detailed usage of the function is
 120 shown in the user manual. We demonstrate an example here.

```

personscoreResult = personscore(result)
summary(personscoreResult)

```

121 Summary of the person score is shown below.

122 Summary of FedIRT Person Score Results:

```

123
124 Ability Estimates:
125 School 1:
126 [1] -1.127097195 -0.922572829 -0.993953038 ...
127 School 2:
128 [1] -1.41454573 1.78068772 1.87469389 ...
129 ...
130
131 End of Summary

```

132 Example of the personfit function

133 We provide a function personfit in the package. The detailed usage of the function is shown
 134 in the user manual. We demonstrate an example here.

```

personfitResult = personfit(result)
summary(personfitResult)

```

135 After getting the result, use personfit function to get the person score result from result by
 136 personfit(result).

137 Summary of FedIRT Person Fit Results:

138 Fit Estimates:

139 School 1:

	Lz	Zh	Infit	Outfit
4	0.7584470759	0.923163304	0.002323484	0.1482672
16	-0.7562447025	-1.131668935	0.005457117	0.1799583
27	0.3417488360	0.357870094	0.005966933	0.1734402
33	-0.9244005411	-1.359789298	0.179834037	0.2266634
...				

147 School 2:

	Lz	Zh	Infit	Outfit
5	-0.90114567	-1.175767350	0.0009824580	0.1535794
8	-1.47957351	-1.888763364	0.1491518127	0.2255230

```

151 18 -0.13292541 -0.228824721 0.1104556086 0.2007658
152 19 -0.17257549 -0.277699184 0.0075031313 0.1350857
153 ...

```

154 Standard error (SE) calculation

155 We follow a typical process of calculating SE in MLE. After obtaining the MLE estimates, the
 156 Hessian matrix, which is the matrix of second-order partial derivatives of the log-likelihood
 157 function with respect to the parameters, is computed at the estimated parameters. The SEs
 158 are then derived from the square roots of the diagonal elements of the inverse Hessian matrix.

159 In our package, call the `SE()` function and input a `fedirt` object to display standard errors of
 160 item parameters.

```
SE(result)
```

161 Below is the result of SE.

```

162 $a
163 [1] 0.0041815497 0.1638884452 0.1204696925 ...
164 $b
165 [1] 272.43863961 0.20737386 1.25896302 ...

```

166 Example of the Shiny App

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

170 In the first step, the server end (e.g., test administrator, school board) can be launched by running
 171 the Shiny app `runserver()` and the client-end Shiny app can be initialized with `runclient()`
 172 with the interface shown below:

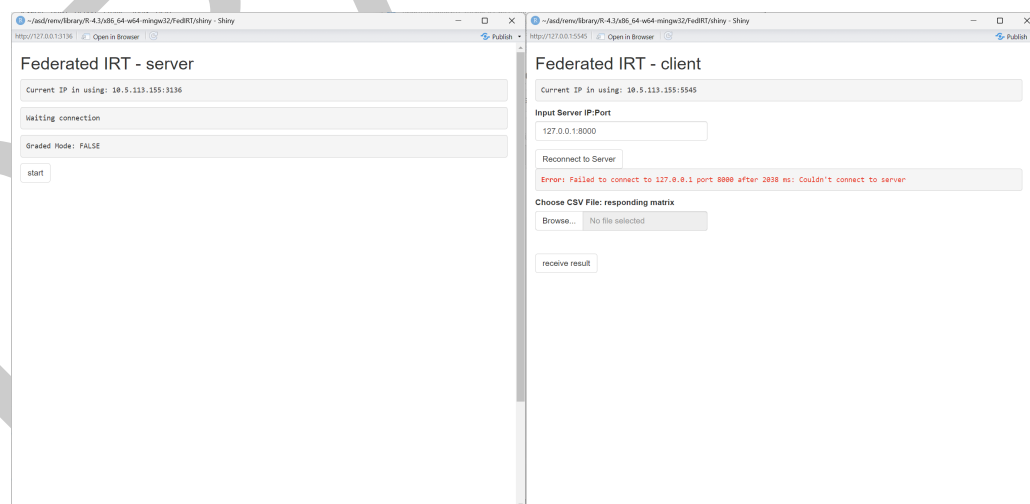


Figure 3: The initial server and client interface.

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

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

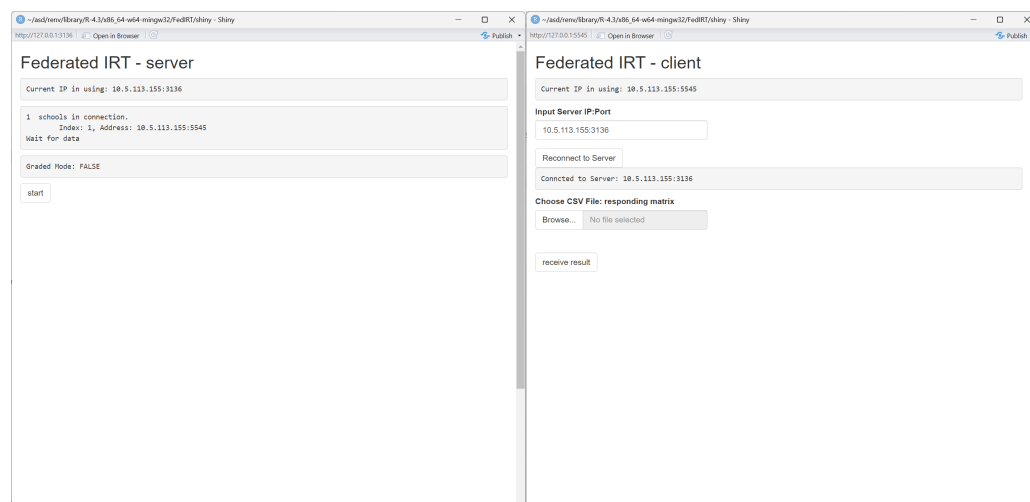


Figure 4: Server and client interface when one school is connected.

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

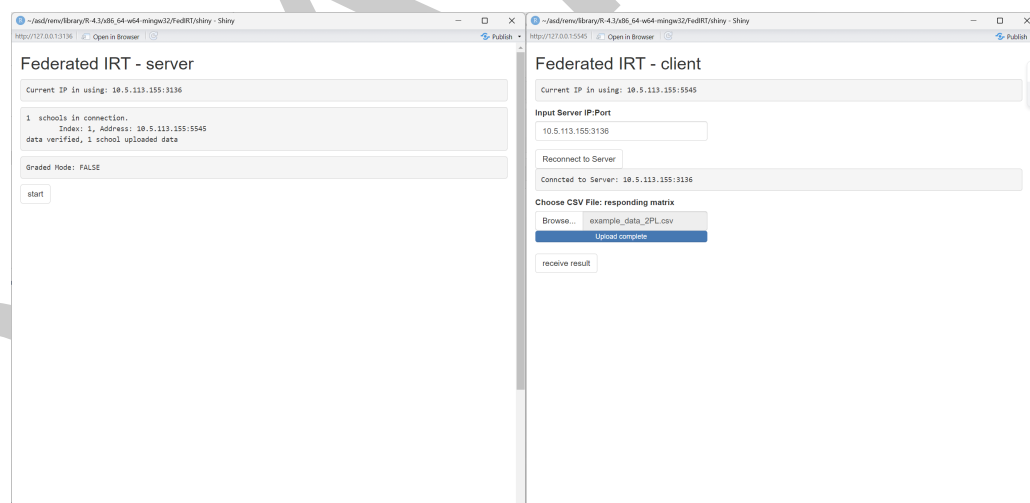


Figure 5: Server interface when one school uploaded dataset and client interface when a dataset is uploaded successfully.

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

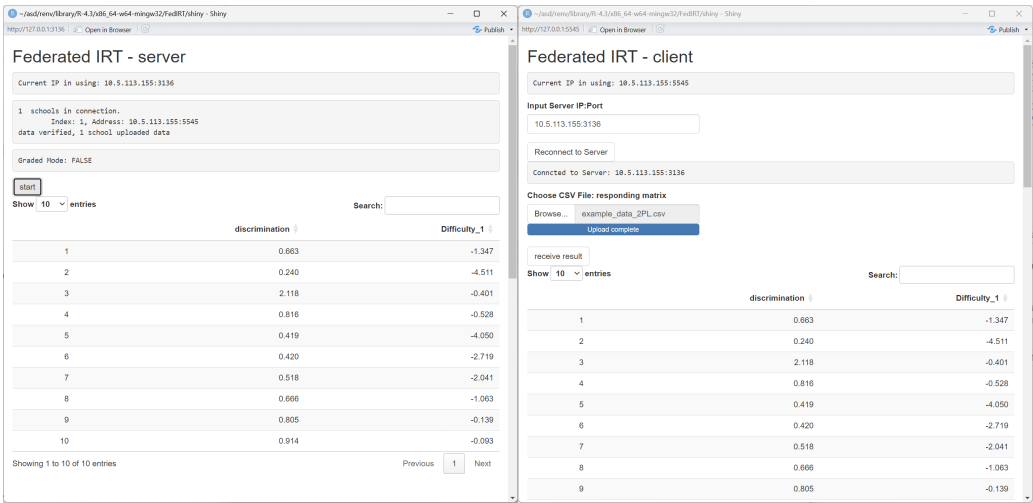


Figure 6: Server - server interface when estimation is completed and client interface when the results received.

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

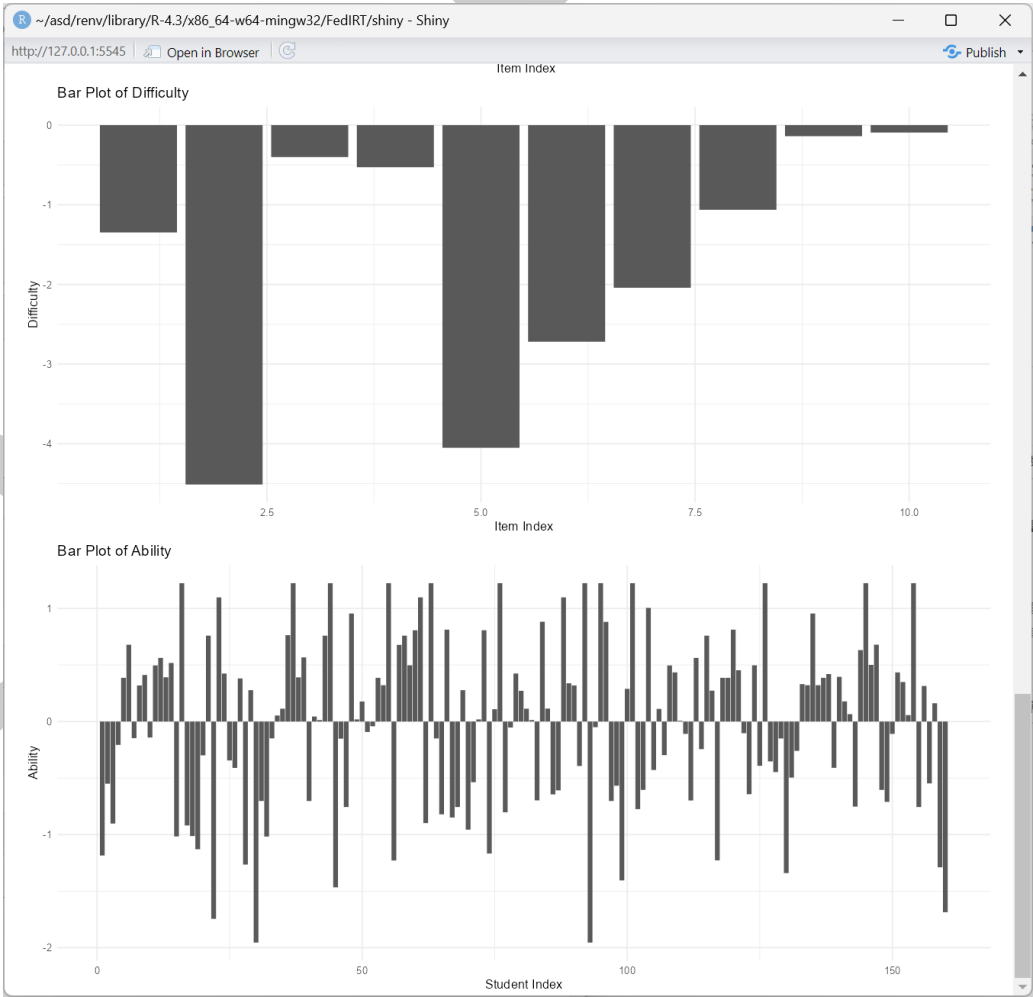


Figure 7: Client interface for displaying results.

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