

Probing Influences on Singaporean Academia

a user guide to the journey

HOME PAGE



- Overview and Objectives of Project PISA
- Variable List
Details of all the predictors available for exploration and analysis
- Key Statistics
Includes distribution of subject scores, breakdown of gender and school type

DATA EXPLORATION



- Variable Overview
For univariate and bivariate analysis. Explore distribution, proportions and key summary statistics of each variable
- Sankey Diagram
Visualise the flow from one variable to another

CONFIRMATORY ANALYSIS



- Analysis of Scores
Evaluate if a user-defined value significantly deviates from the population mean
- Analysis of Predictors
Ascertain whether there exists a statistically significant distinction among two or more categorical groups

CLUSTER ANALYSIS



- Association between Predictors
Measure dependency/association between two variables
- Latent Class Analysis
Show the distribution of categorical variables within each latent class.

REGRESSION ANALYSIS



- *Users can select target and predictor variables as well as the train-test partition split to be applied across the three available models – regression tree, random forest, and gradient boosting*
- *Within each model, the user can further calibrate the model, and select resampling options*
- *Model evaluation methods include plots of model fit as well as statistical metrics such as R-square, Root Mean Square Error, and Mean Absolute Error.*

1. Home Page

1 Welcome to Project PISA

According to the latest OECD's Programme for International Student Assessment (PISA) 2022, which measures 15-year-olds' ability to use their reading, mathematics, and science knowledge and skills to meet real-life challenges, socioeconomic status accounted for 17% of the variation in mathematics performance in Singapore (compared to 15% on average across OECD countries). Clearly, Singapore's success does not translate to success for every student. Why then do some students outperform others? And is socioeconomic status the only factor for success?

Our team believes that knowledge is power. While causality cannot and should not be easily drawn between the various forces of influence and academic performance, a more detailed and nuanced understanding of these factors would highlight potential areas to focus on when engaging parents and students as well as when developing education and socioeconomic policies for a more inclusive and equitable society.

2 Variable List

Show entries Search:

Variables	Description
1 Gender	Student (Standardized) Gender
2 SchoolType	Public or Private
3 Loneliness	I feel lonely at school.
4 ClassroomSafety	Agree/disagree: I feel safe in my classrooms at school.
5 TeacherSupport	How often: The teacher helps students with their learning.
6 Homework_Math	How much time spent on homework in Math homework
7 Homework_Reading	How much time spent on homework in Reading homework
8 Homework_Science	How much time spent on homework in Science homework
9 ParentsEducation	Highest level of education of parents (ISCED)
10 Immigration	Index on immigrant background (OECD definition)

Showing 1 to 10 of 25 entries Previous 2 3 Next

3

5158 Participating Students

164 Participating Schools

2616, 51% Male Students

2542, 49% Female Students

4804, 93% Public School

354, 7% Private School

4 Distribution of Subject Scores

Science

Reading

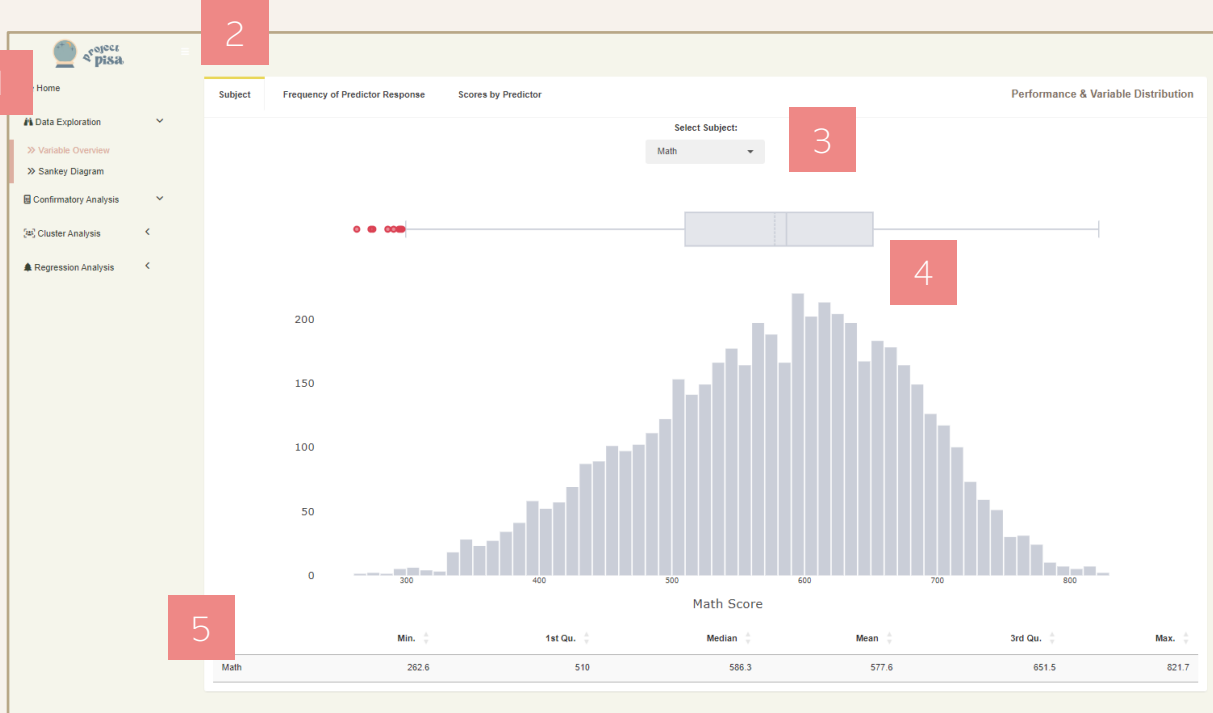
Math

200 400 600 800 Scores

On this page, the following information can be found:

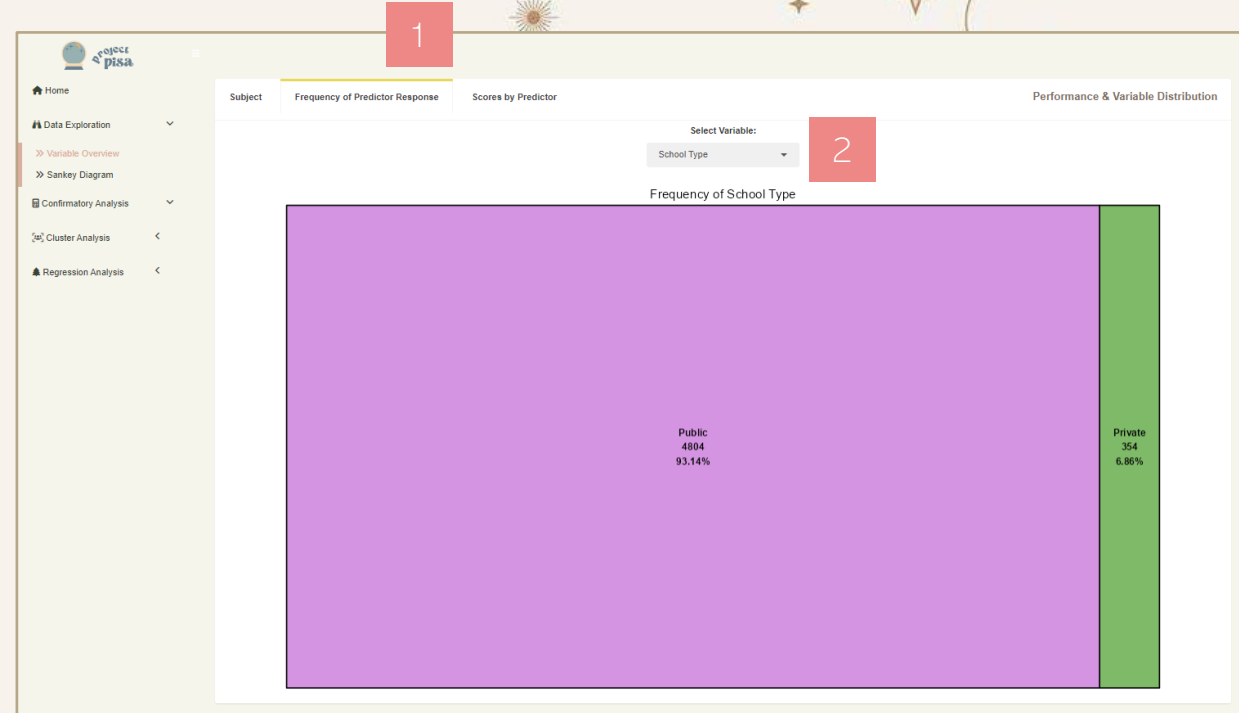
1. **Overview and objectives** of the project, accompanied with a short description of the web application.
2. **Variable List** which consists of the complete set of variables and what each of them represents.
3. **Key statistics** of the dataset which includes the total number of participants and schools involved, the split by gender and school type.
4. **Distribution of target variables** (i.e. Subject Scores).

2. Data Exploration



The first feature of the Data Exploration tab is the **Variable Overview** page.

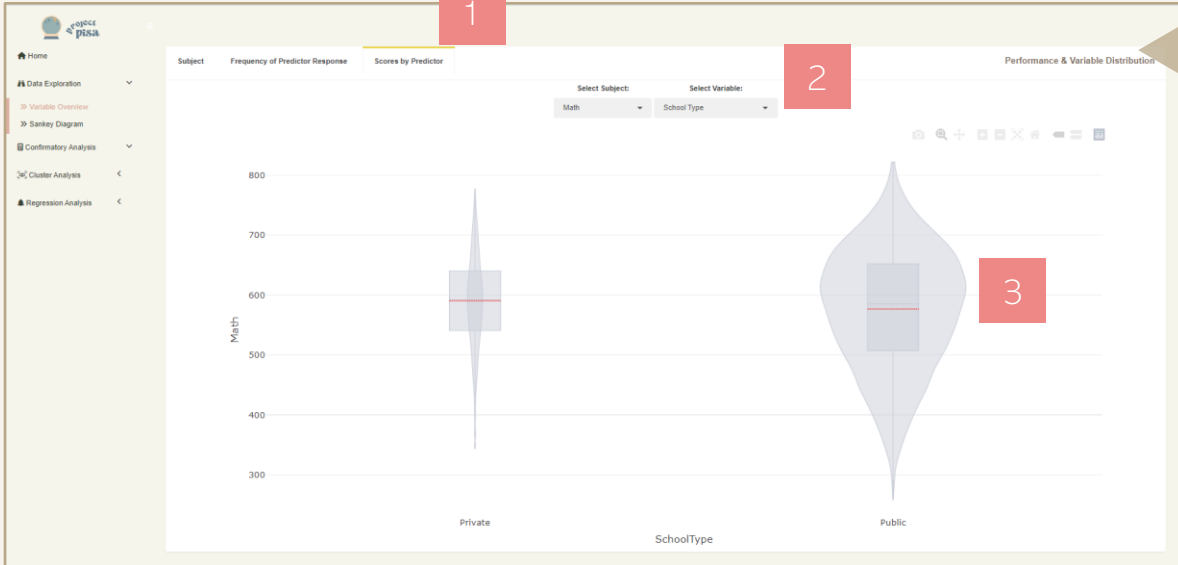
1. First, click on **∇** beside **Data Exploration** on the side menu to show the submenu items. Click on **Variable Overview**.
2. Click on the **Subject** tab on the right panel.
3. Click on the **∇** at the dropdown menu to select the target variable.
4. Once a selection is made, users can **hover over the interactive histogram** to explore the distribution of data. Hovering over the **interactive boxplot** allows users to see the key summary statistics of the data such as the mean, median, minimum and maximum values.
5. The **table** allows a quick glance of the key summary statistics as well.



The next feature of the Data Exploration tab is the **Frequency of Predictor Responses** tab. The tree-map allows users to visualize the proportion of data represented using area size.

1. Click on the **Frequency of Predictor Responses** tab at the top of the right panel.
2. Click on the **∇** at the dropdown menu to select the variable of interest
3. Once a selection is made, the tree-map would be refreshed to show users the split of responses for each variable.

2. Data Exploration

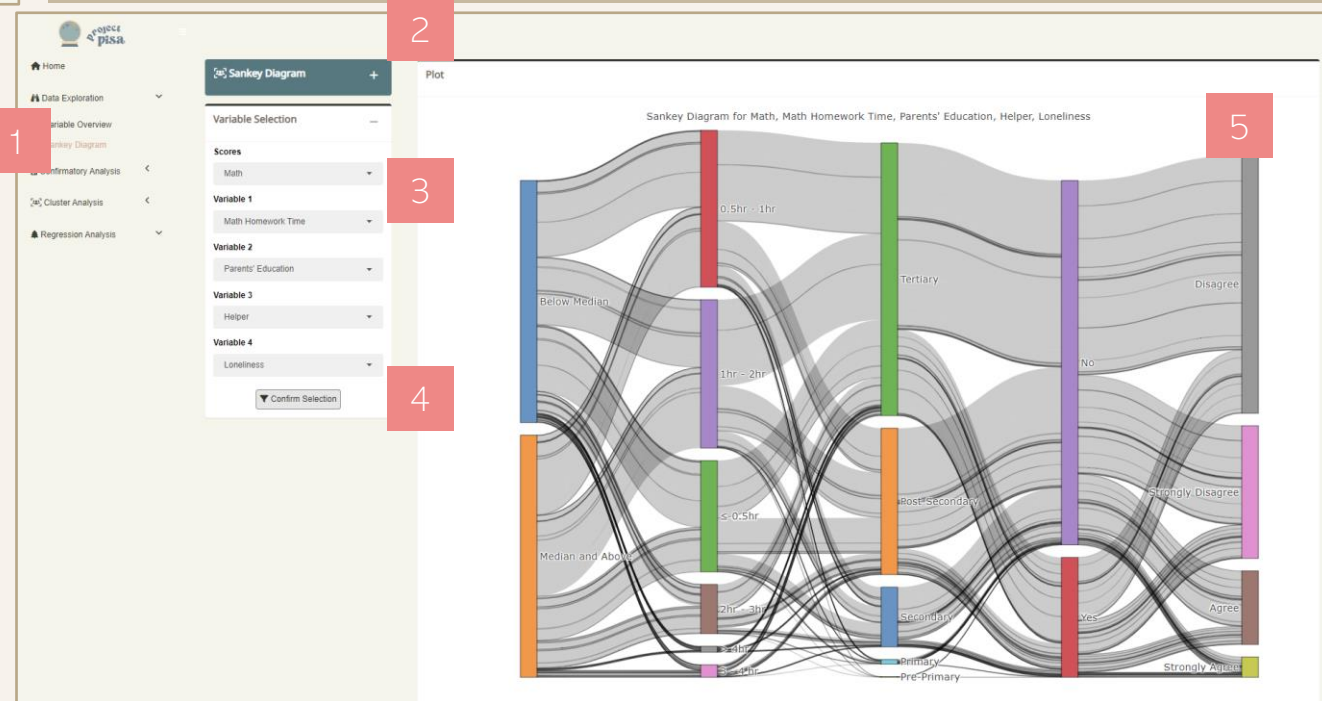


The next feature of the Data Exploration tab is the **Scores by Predictor** tab. The box-violin plot combines a boxplot and a kernel density plot. It is used to visualize the distribution of Subject Scores (i.e. target variable) between the different categorical responses for each variable.

1. Click on the **Scores by Predictor** tab at the top of the right panel.
2. Click on the ▾ of each dropdown menu to select the target variable and predictor of interest.
3. Once a selection is made, the box-violin plot would be refreshed to show the distribution of Subject Scores across the various responses for the selected variable.
4. Users can hover over the box-violin plot to access key summary statistics of the distribution such as the mean, median, minimum, and maximum values.

A **Sankey Diagram** illustrates the flow of data between different categories or clusters, with the line thickness corresponding to the frequencies of the flows.

1. Click on ▾ beside **Data Exploration** on the side menu to expose the submenu if the submenu is collapsed. Click on **Sankey Diagram**.
2. The + beside the tab header provides a short description for users to understand the purpose of the illustration of the tab.
3. Click on the ▾ of each dropdown menu to select 1 target variable interest and 4 predictor variables.
4. Click on the ▾ **Confirm Selection** button once selection is made.
5. Hover over the Sankey Diagram to explore how the responses for each variable flows to another.



3. Confirmatory Analysis

The Analysis of Scores tab allow users to evaluate if a user-defined value significantly deviates from the population mean.

1. Click on ▾ beside **Confirmatory Analysis** on the side menu to expose the submenu if the submenu is collapsed. Select **Analysis of Scores**.
2. The + beside the header provides a short description for users to understand the purpose of the tab.
3. Click on the ▾ of dropdown menu to select a **target variable interest**.
4. User can input a **test score** to perform a statistical test on. User can also hover on the box to display the up/down arrows to change the input.
5. Input a **desired confidence level**.
6. Select the **desired bin width**.
7. Select **test type**.
8. Select the **type of effect size**.
9. Click on **Run Analysis**.
10. The top-right panel displays the data distribution with the statistical results. Results in the header includes the type of test performed, the test statistics and the p-value of the test statistic. Table on the right shows type of test performed for each test type >>

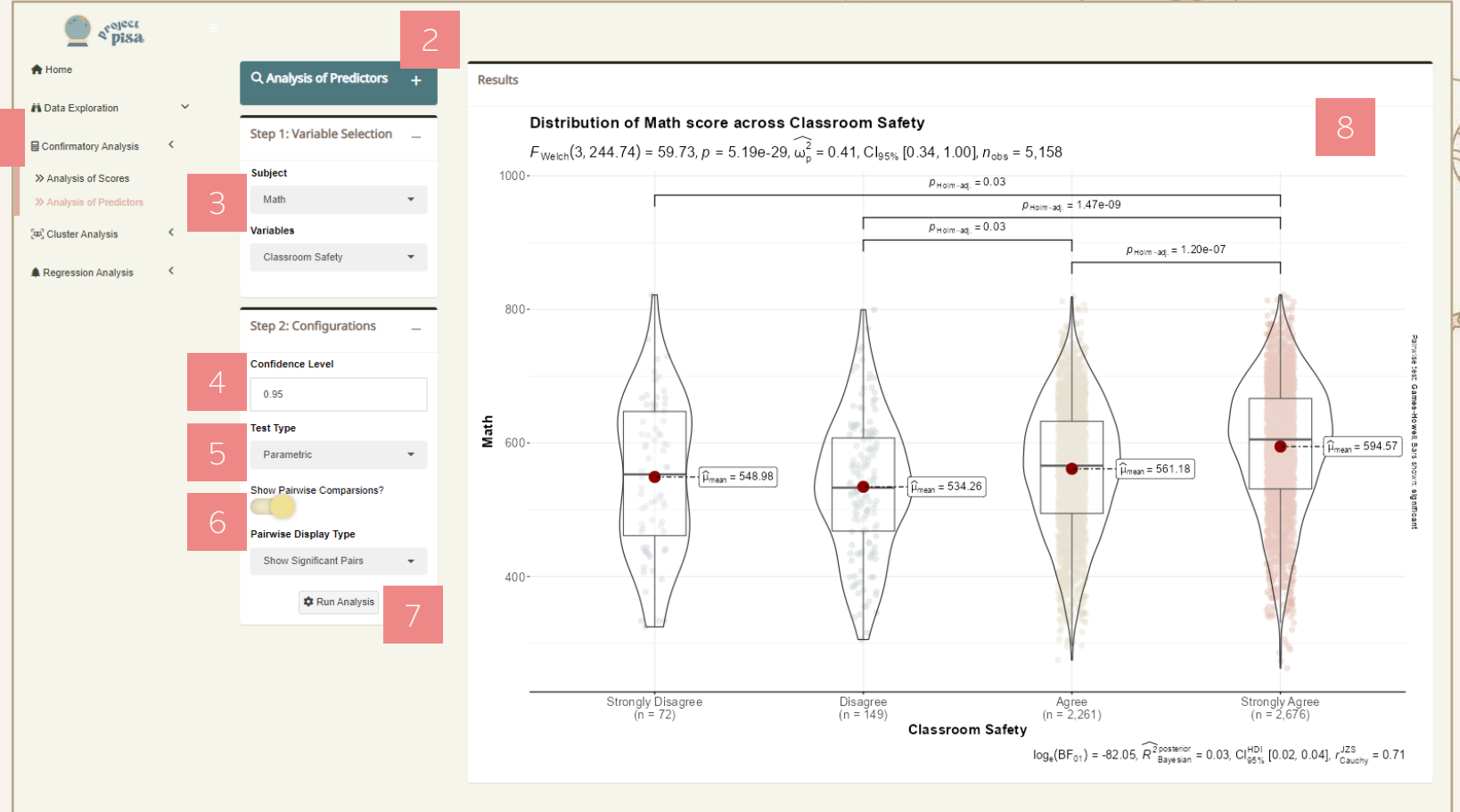
Type	Centrality Type	Test
Parametric	Mean	One-sample Student's <i>t</i> -test
Non-parametric	Median	One-sample Wilcoxon test
Robust	Trimmed Mean	Bootstrap- <i>t</i> method for one-sample test
Bayesian	MAP Estimator	One-sample Student's <i>t</i> -test

11. Points on the **QQ plot** provide an indication of normality of the dataset. If the data is normally distributed, the points will fall on the 45-degree reference line. If the data is not normally distributed, the points will deviate from the reference line.
12. The **Anderson-Darling test** is a statistical test that determines if a data set follows a **normal distribution**, and then determining the p-value for the statistic. P-value is used to determine whether a null hypothesis should be accepted or rejected.

3. Confirmatory Analysis

An ANOVA test is a statistical method utilized to ascertain whether there exists a statistically significant distinction among two or more categorical groups by examining differences in means using variance.

1. Click on ∇ beside **Confirmatory Analysis** on the side menu to expose the submenu if the submenu is collapsed. Click on **Analysis of Predictors**.
2. The **+** beside the header provides a short description for users to understand the purpose of the tab.
3. Click on the ∇ of dropdown menu to select a **target variable** and a **predictor variable**.
4. Input a **desired confidence level**.
5. Select **test type**.
6. User should toggle the switch to the right should they wish to display the **Pairwise Comparisons** in the plot. If Pairwise Comparisons is turned on, the option for Pairwise Display Type would be displayed.
7. Click on **Run Analysis** once options are selected.



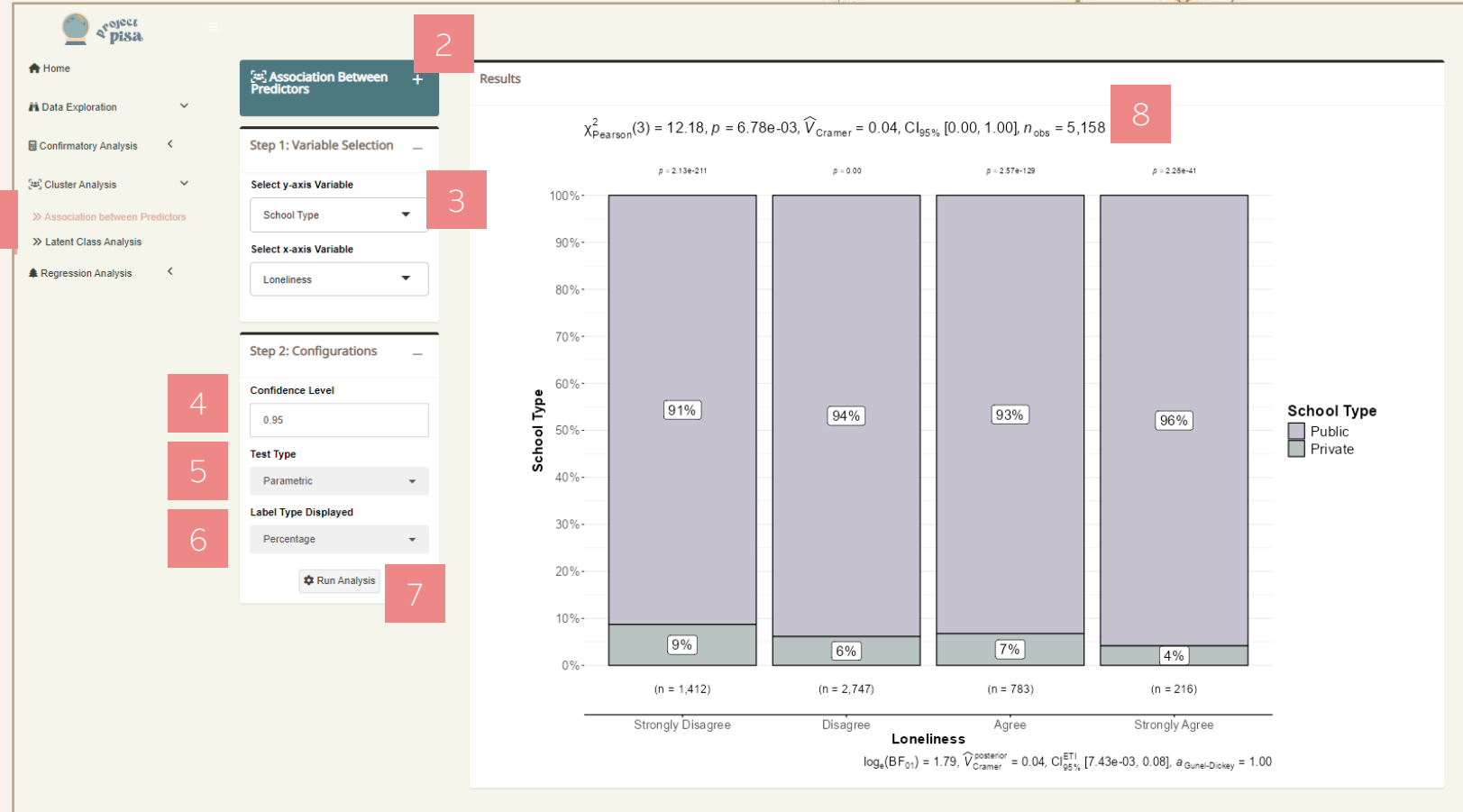
8. The right panel displays the data distribution of each response of the predictor variable with the statistical results. Results in the header includes the type of test performed, the test statistics and the p-value of the test statistic. See table on the right for type of test performed for each test type selection >>

Type	Centrality Type	Test
Parametric	Mean	Student's or Welch's t-test
Non-parametric	Median	Mann-Whitney U test
Robust	Trimmed Mean	Yuen's test for trimmed means
Bayesian	MAP	Student's t-test

4. Cluster Analysis

A test of association is a hypothesis test designed to establish and quantify the relationship between two distinct factors. When dealing with categorical variables, multicollinearity can be identified using the chi-square test. Multicollinearity can result in less reliable statistical inferences and may lead to skewed or misleading results.

1. Click on ∇ beside **Cluster Analysis** on the side menu to expose the submenu if the submenu is collapsed. Click on **Association between Predictor**.
2. The **+** beside the header provides a short description for users to understand the purpose of the tab.
3. Click on the ∇ of dropdown menu to select a variable to plot on the y-axis and x-axis respectively.
4. Input a desired confidence level.
5. Select **test type**.
6. Users have the option to select the type of labels to be displayed on the plot. This can be in count or percentage or both.
7. Click on **Run Analysis** once options are selected.



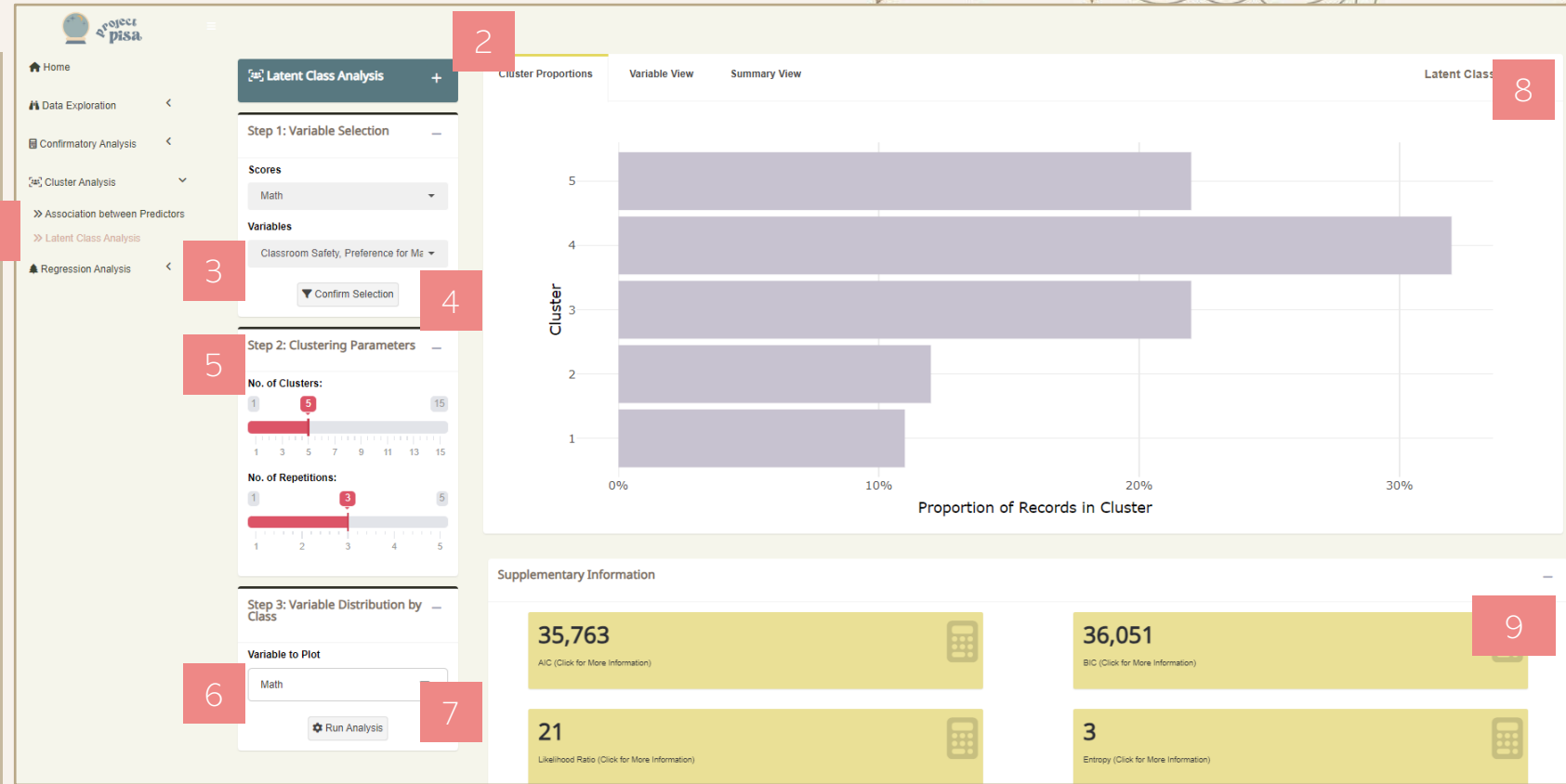
8. The right panel displays the data distribution of each response of the predictor variable with the statistical results. Results in the header includes the type of test performed, the test statistics and the p-value of the test statistic. See table on the right for type of test performed for each test type selection >>

Type	Test
Parametric	Pearson's chi-squared test
Non-parametric	Bayesian Pearson's chi-squared test
Robust	McNemar's chi-squared test
Bayesian	-

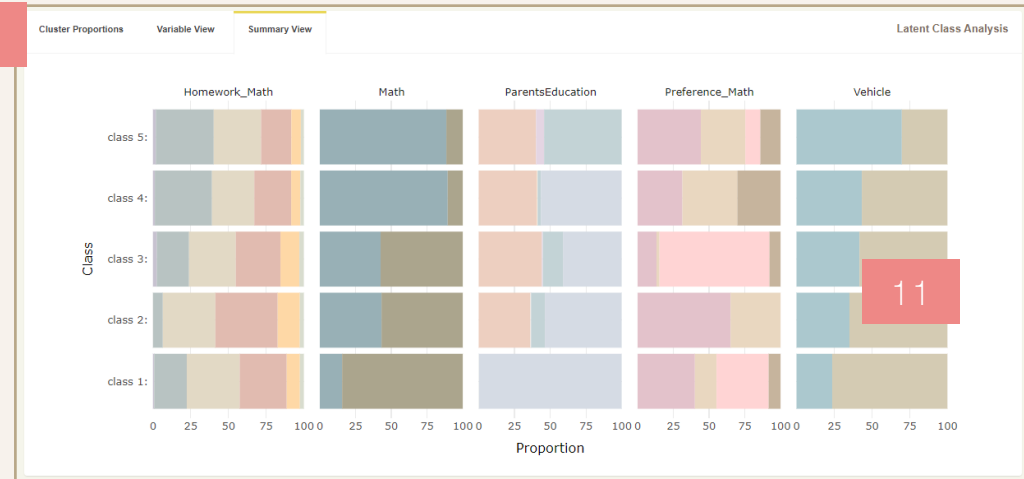
4. Cluster Analysis

A Latent Class Analysis (LCA) Bar Plot is used to show the distribution of categorical variables within each latent class.

1. Click on ▾ beside **Cluster Analysis** on the side menu to expose the submenu if the submenu is collapsed. Click on **Latent Class Analysis**.
2. The + beside the header provides a short description for users to understand the purpose of the tab.
3. Click on the ▾ of dropdown menu to select target variable(s) and predictor variable(s).
4. Click on **Confirm Selection** once variables have been selected.
5. Select the desired number of clusters and number of repetitions.
6. The selection for **Variable to Plot** is dependent on the variables selected under **Variable Selection**. This will affect the Variable View tab.
7. Click on **Run Analysis** once options are selected.
8. The **Cluster Proportions** tab shows the proportion of data that falls under each cluster. Hover over the plot to show the exact proportion amount.
9. **Supplementary Information** box can be expanded to show the relevant model diagnostic statistic.



10. The **Variable View** and **Summary View** tabs show the stacked bars of the LCA by class and individual variable and by class and all variables, respectively.
11. Hover over the stacked bars to explore the attributes and proportions of each variable for each class.



5. Regression Analysis

A Regression Tree is a machine learning algorithm that partitions the data into subsets. The goal of a regression tree is to encapsulate the training data in the smallest possible tree, i.e. simplest possible explanation for the variation in scores.

1. Click on ▾ beside **Regression Analysis** on the side menu to expose the submenu if the submenu is collapsed. Click on **Regression Tree**.
2. Click on the ▾ of dropdown menu to select the target variable, predictor variable(s), and train-test partition ratio of choice for regression analysis.
3. The + beside the header provides a short description for users to understand the purpose of the tab.
4. Click on **Run Analysis** to initiate the model building process.

The screenshot displays the Project DISA Regression Tree interface. It is divided into several sections:

- Step 1: Construct Model:** Includes a sidebar with a dropdown menu for 'Subject' (set to 'Math') and 'Variables' (set to 'School Type, Lonelin'). A 'Train-Test Partition' slider is set to 0.8. A 'Continue to Step 2 >>' button is visible.
- Step 2: Model Initiation:** Features a 'Run Analysis' button.
- Step 3: Tuning Parameters:** Includes sliders for 'Minimum Split' (set to 5) and 'Maximum Depth' (set to 10). A 'Tune Model' button is present.
- Fit Assessment:** Contains two scatter plots: 'Predicted vs Actual' and 'Residuals vs Actual'. The 'Predicted vs Actual' plot shows a positive correlation with a regression line. The 'Residuals vs Actual' plot shows a negative correlation with a regression line.
- Complexity Parameter:** A plot showing 'X-val Relative Error' vs 'number of splits'. The error decreases as the number of splits increases, reaching a minimum at 131 splits. A yellow box highlights the 'Best Complexity Parameter' as 0.00267.
- Evaluation Metrics:** Three yellow boxes display: R-Square (0.22), RMSE (89.645), and MAE (71.392).
- Regression Tree and Variable Importance:** A tree diagram is shown on the right side of the interface.

5. The **Tuning Parameters** panel is only displayed after the initial model has been triggered. The model can subsequently be calibrated and refined by selecting the minimum split and maximum depth.

6. If the “**Model with the best complexity parameter**” is switched on, no further action is required. Else, input a complexity parameter value.
7. Click on **Tune Model**. Re-trigger the button to apply any further changes.

8. Users explore the model fit plots, model evaluation statistics, interactive regression tree, variable importance plot. **Any changes to the variables or parameters requires retriggering of 'Run Analysis' before 'Tune Model' can be used.**

5. Regression Analysis

Random Forest combines the opinions of many “trees” (individual models) to make better predictions, creating a more robust and accurate overall model.

1. Click on ▾ beside **Regression Analysis** on the side menu to expose the submenu if the submenu is collapsed. Click on **Random Forest**.
2. The selection for **Step 1: Construct Model** follows the selection made previously. If not previously selected, click on the dropdown menus to make the selection.
3. The + beside the header provides a short description for users to understand the purpose of the tab.
4. To further calibrate the model, users have the option to **change the number of trees and minimum node size** by toggling the sliding bar. **Variable importance measure and splitting rule** can be changed by selecting the **desired selection** from the dropdown menu.

The screenshot displays the Project PISA interface for a Random Forest regression analysis. The interface is divided into several sections:

- Step 1: Construct Model:** Includes dropdowns for Subject (Math) and Variables (School Type, Lonelin). A Train-Test Partition slider is set to 0.8.
- Step 2: Tuning Parameters:** Includes a No. of Trees slider (set to 50), a Variable Importance Measure dropdown (Gini Importance), a Split Rule dropdown (Variance), and a Minimum Node Size slider (set to 5).
- Step 3: Resampling Options:** Includes a Resampling Method dropdown (Cross Validation) and a K-fold dropdown (10).
- Summary Metrics:** R-Square (0.346), RMSE (80.328), and MAE (65.216).
- Fit Assessment:** A scatter plot of Predicted vs Actual values with a regression line and shaded confidence intervals.
- Residuals vs Actual:** A scatter plot of Residuals vs Actual values with a regression line and shaded confidence intervals.
- Variable Importance:** A horizontal bar chart showing the importance of various variables, with ParentsEducationTertiary being the most important.

Numbered callouts (1-7) highlight specific UI elements: 1 points to the Regression Analysis menu, 2 to the Step 1 header, 3 to the Random Forest tab header, 4 to the No. of Trees slider, 5 to the Resampling Method dropdown, 6 to the Run Analysis button, and 7 to the Fit Assessment and Variable Importance sections.

5. Resampling options can also be amended using the dropdown menu. K-Fold and Repeat Count would only be displayed if Repeated Cross Validation is selected. K-Fold option would be displayed when Cross-Validation is selected. No further options if a user chooses to use Bootstrap.
6. Trigger the **Run Analysis** button once selections are made.
7. User can subsequently explore the statistical results, variable importance, and model fit plots. Any further changes to the options will require retriggering of the **Run Analysis** button to refresh the results.

5. Regression Analysis

Gradient Boosting is an ensemble machine learning technique that combines the predictions from several models to improve the overall predictive accuracy.

1. Click on ∇ beside **Regression Analysis** on the side menu to expose the submenu if the submenu is collapsed. Click on **Gradient Boosting**.
2. The selection for **Step 1: Construct Model** follows the selection made previously. If not previously selected, click on the dropdown menus to make the selection.
3. The **+** beside the header provides a short description for users to understand the purpose of the tab.

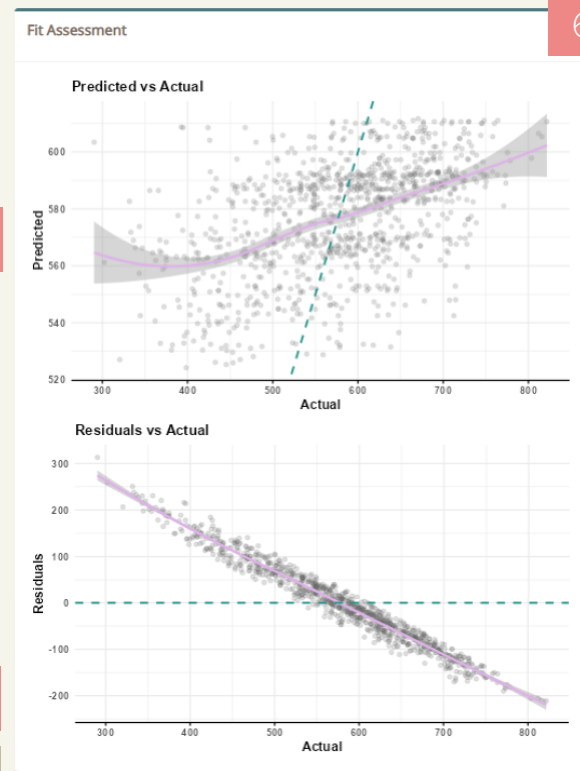
The screenshot shows the 'Gradient Boosting' configuration interface. It is divided into three main steps:

- Step 1: Construct Model:** Shows 'Subject' set to 'Math' and 'Variables' set to 'School Type, Lonelin'. The 'Train-Test Partition' is set to 0.8.
- Step 2: Resampling Options:** Shows 'Resampling Method' set to 'Cross Validation' and 'K-fold' set to 10.
- Step 3: Tuning Parameters:** Shows sliders for 'Min. Node Size' (set to 5), 'Max. Tree Depth (select range)' (set to 2-3), 'Boosting Iterations (choose 2 for comparison)' (set to 10-50), and 'Learning Rate' (set to 0.01).

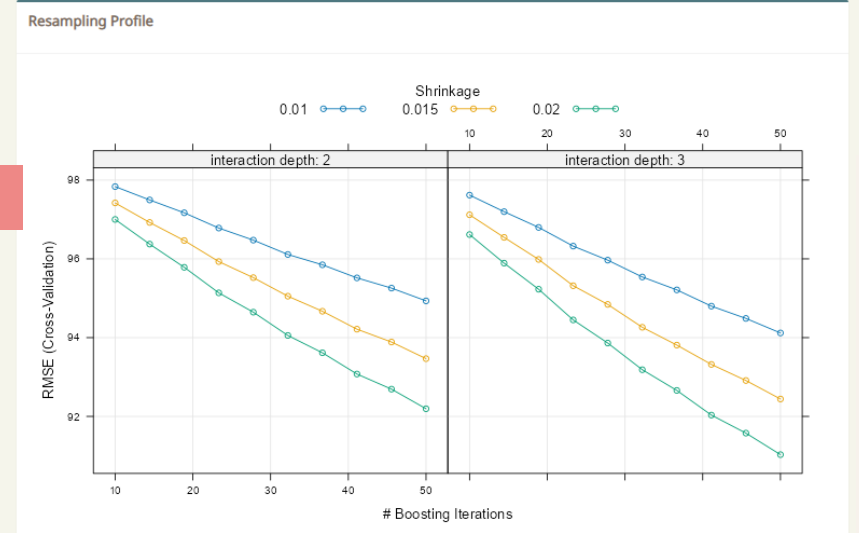
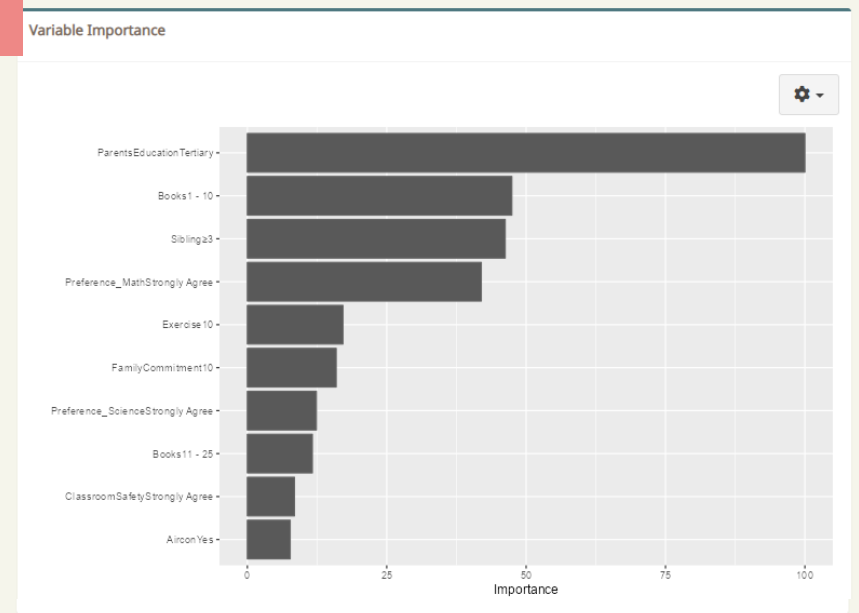
4. To further calibrate the model, users have the option to change the min. node size, max. tree depth range, and boosting iteration selections by toggling the respective sliding bars. Learning rate can be changed by inputting or using the up/down arrows.

Summary statistics for the regression analysis:

- 0.241** R-Square
- 89.172** RMSE
- 72.674** MAE



Number of Trees	Max. Tree Depth	Learning Rate	Min. Node Size	
60	50	3	0.02	5



5. Regression Analysis

5. Trigger the **Run Analysis** button once selections are made.
6. User can subsequently explore the statistical results, variable importance, and model fit plots.
7. **Best Tune** and **Resampling Profile** panels are displayed only when **Cross Validation** or **Repeated Cross Validation** methods are selected under **Step 2: Resampling Options**.

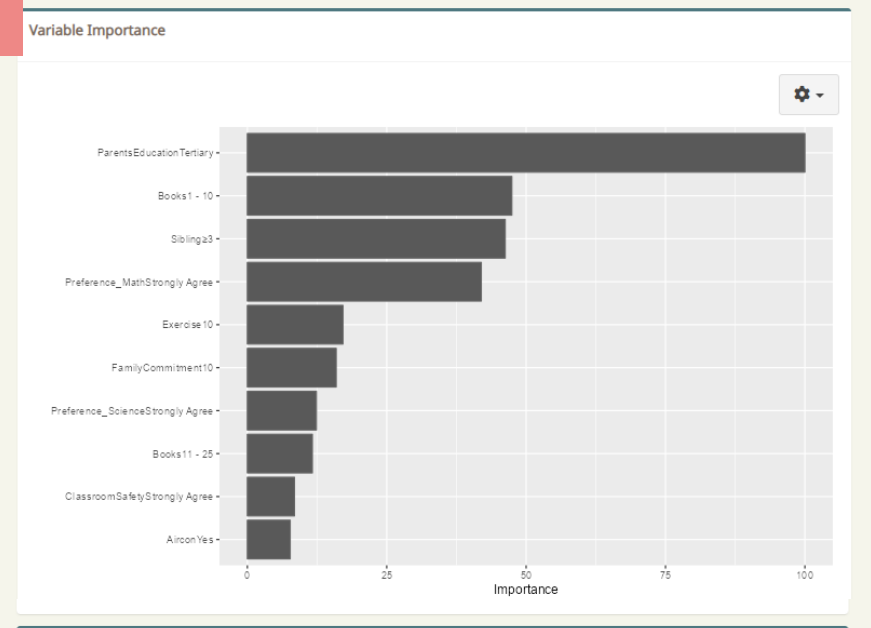
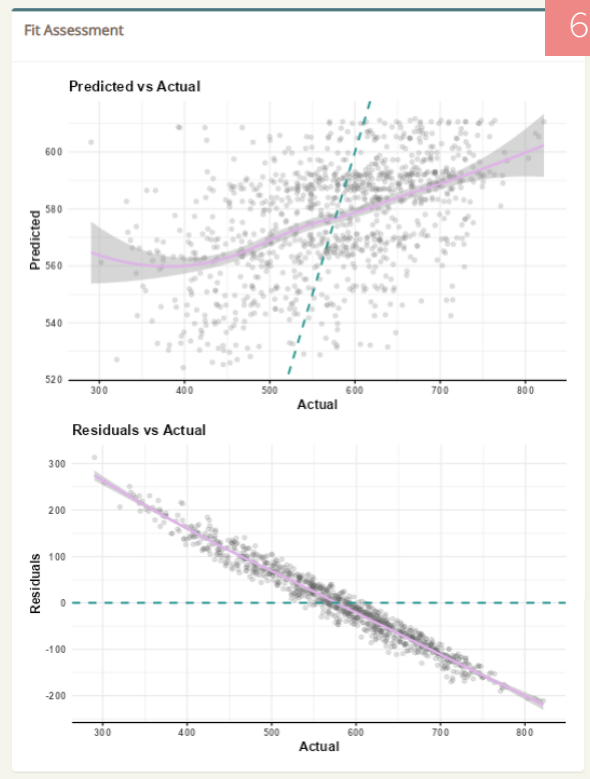
The interface shows the following configuration steps:

- Step 1: Construct Model:** Subject: Math; Variables: School Type, Lonelin; Train-Test Partition: 0.8.
- Step 2: Resampling Options:** Resampling Method: Cross Validation; K-fold: 10.
- Step 3: Tuning Parameters:** Min. Node Size: 5; Max. Tree Depth: 10; Boosting Iterations: 50; Learning Rate: 0.01.

A **Run Analysis** button is visible at the bottom of the tuning parameters section.

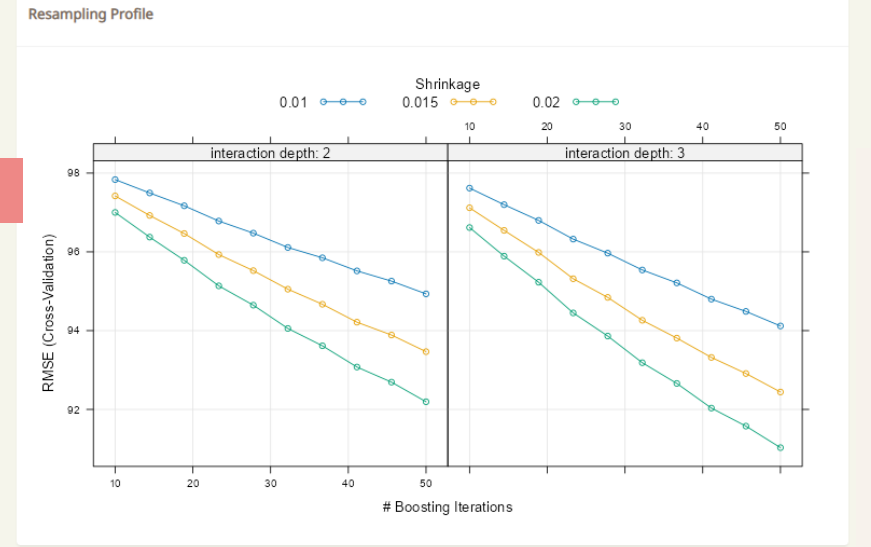
Summary statistics for the model fit:

- 0.241** R-Square
- 89.172** RMSE
- 72.674** MAE



Best Tune

Number of Trees	Max. Tree Depth	Learning Rate	Min. Node Size
60	50	3	0.02
			5



Best Tune and Resampling Profile panels show the best tuning parameters amongst the range selected by the user. This allow users to further finetune the model to obtain the best results.