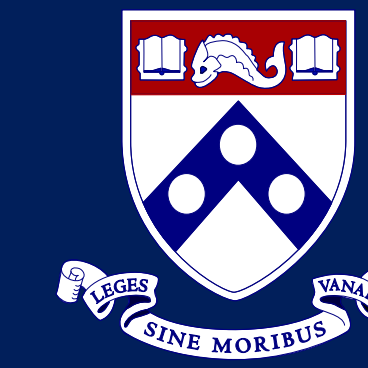# Automation of Statistics Summary and Analysis

Lihai Song, MS[1], Zi Wang, MS[1], Jing Huang, PHD,[1,2]
1.Department of Biomedical and informatics, The children's Hospital of Philadelphia
2.Department of Biostatistics, Epidemiology and Informatics
University of Pennsylvania Perelman School of Medicine

## BACKGROUND

**Motivation**

Basic statistical analyses, e.g. summary statistics, regression analysis and testing of association, are used in almost all the medical research. The analyses can be repetitive and require tedious work of logging to ensure reproducibility of the study. To make statistics easier for medical searchers, we developed a User Interface (UI) of statistical analysis which automates basic statistical analysis, produces publishable tables and figures with just simple clicks. Our UI will provide a convenient tool for medical researchers to explore data, conduct basic analysis and increase reproducibility.
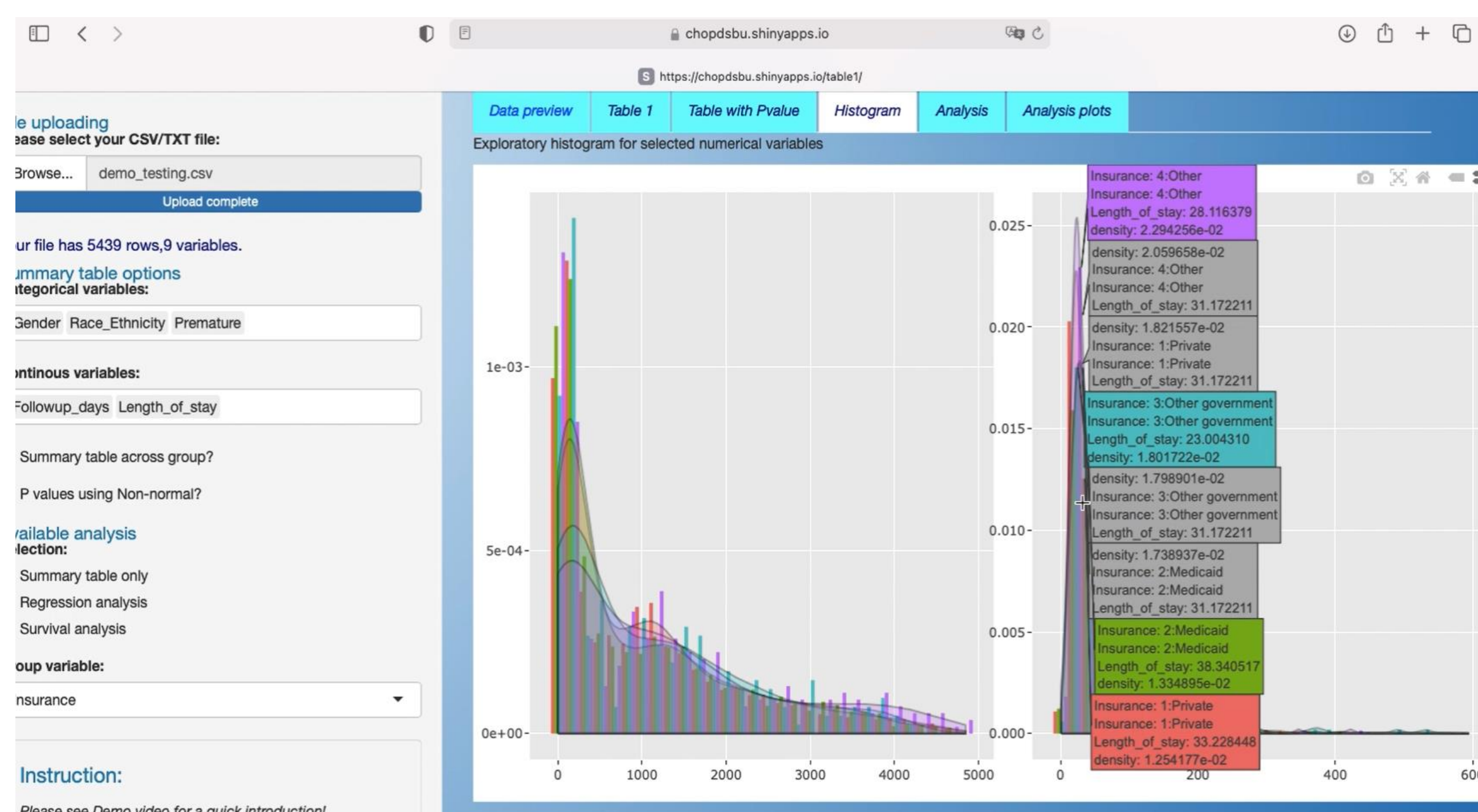
**The developed user interface**

●Easy to use for researchers with less statistical programming experience.
●Fast response substantially improve research efficiency.
●Reduce errors by hand coding and produce publishable tables and figures.

## OBJECTIVE

●To reduce research burden on researchers by providing automated statistical analysis tools.
●To provide various analysis with simple clicks.
●To provide downloadable and publishable tables and figures.
●To improve research reproducibility.

## FEATURES

●**User friendly interface:**
●This UI provides typing-free click & select interface including a built-in tutorial video and an analysis results with download buttons.



●**Instantaneous output**
●Demographics and baseline characteristics summary statistics and interactive graph functionality of continuous variables.
● Regression analysis, e.g. logistic regression, linear regression, and survival analysis can be performed. Publication-quality tables and figures can be downloaded directly from the UI.
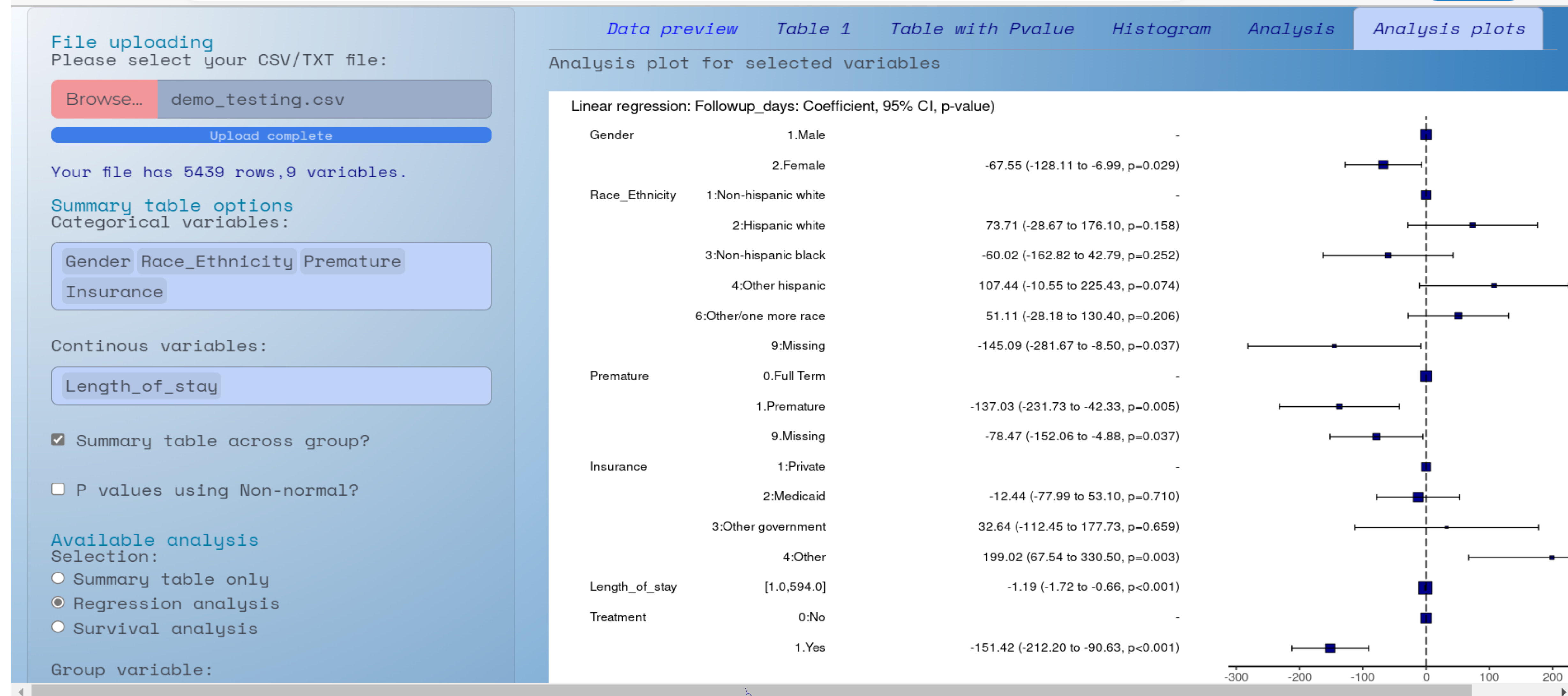
## RESULTS

### Figure 1. Main page of the UI



### Figure 2. Summary Statistics(Table 1)



### Figure 3. Survival Analysis graph



### Figure 4: Univariate and Multivariate Logistic Regression Analysis Results



## LIMITATIONS

● Subgroup is not supported in summary table.
● Interaction terms are not included in all analysis reports.
●The accepted CSV file has size limit.
●This UI can NOT do data cleaning and it requires high quality of analytic data set.

## CONCLUSIONS

●The developed UI automates basic statistical analysis, substantially reduces errors of hand coding, and improves the productivity of research.
●The developed UI makes data exploration and analysis easy for medical researchers, it increases the efficiency and productivity of medical research.
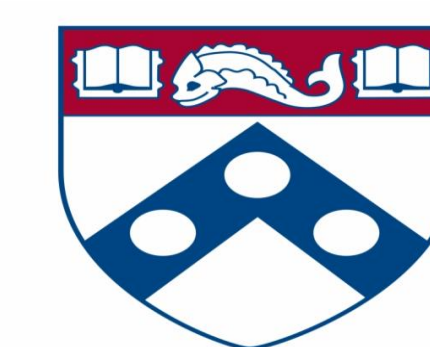●Analysis using advanced statistical methods are under-development.

## FUTURE WORK

●**More to come:**
●We are working on other commonly used basic statistical analysis.
●Advanced statistical analysis methods will be included in the next version.
●More file formats will be supported in the next version.
● Summary table will allow subgroup statistics.

Children's Hospital of Philadelphia RESEARCH INSTITUTE

Perelman School of Medicine UNIVERSITY OF PENNSYLVANIA