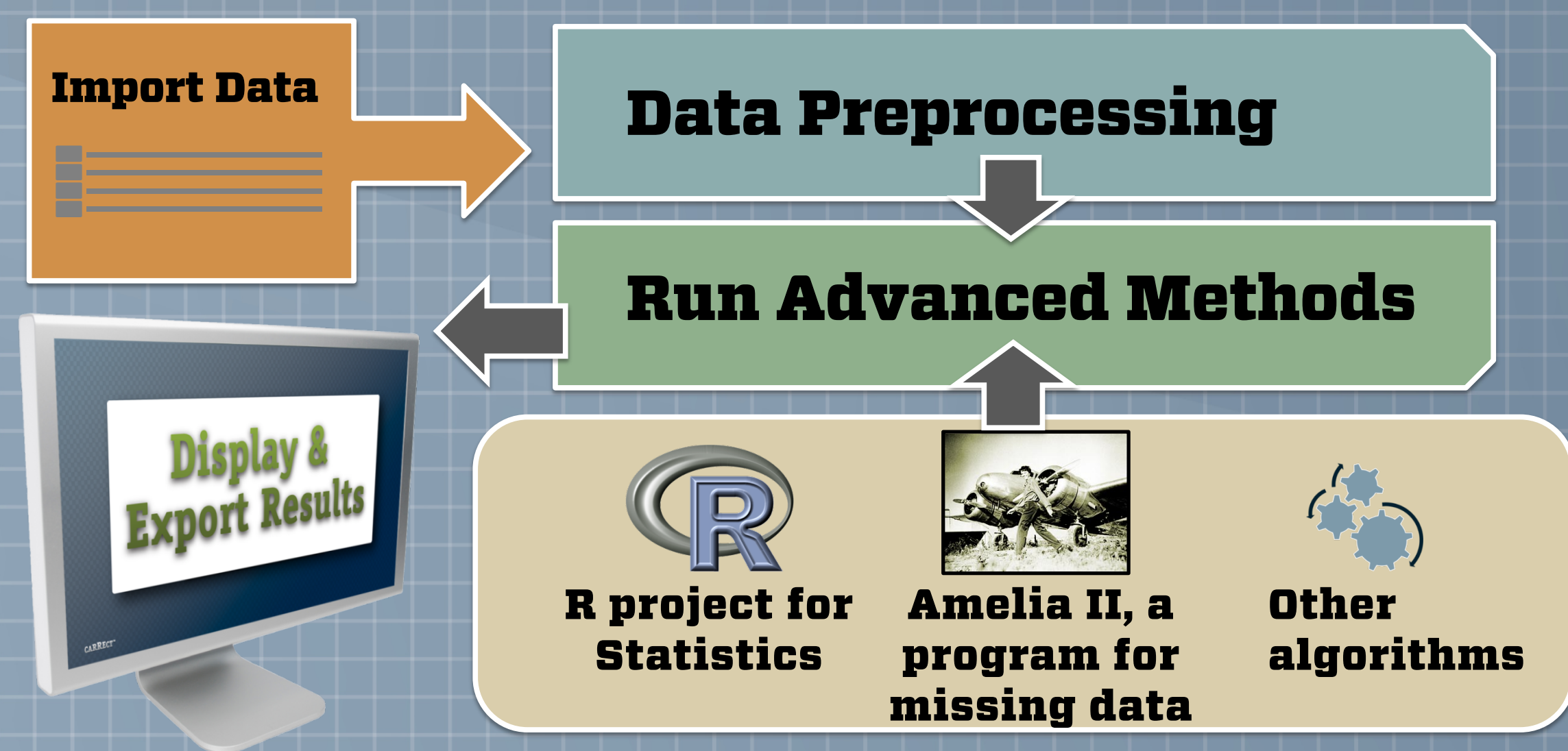# Bias Reduction in Statistical Analyses of Incomplete Data Sets

Jonathan Lansey, Paul Picciano, Ph.D., Ian Yohai, Ph.D., Robert McCormack, Ph.D. - Aptima, Inc.
Fred Grant, Ph.D., Robert Gern, Ph.D - Northrop Grumman Corp.

## Introduction

Analyses produced by epidemiologists and public health practitioners are susceptible to bias from a number of sources. It often requires a great deal of expertise to understand and apply the multitude of statistical methods that are avaialble to diagnose and correct for these biases. To address this challenge, Aptima began development of CARRECT, the Collaborative Automation Reliably Remediating Erroneous Conclusion Threats system. When complete, CARRECT will support statistical bias reduction and improved analyses and decision making by engaging the user in a collaborative process in which the technology is transparent to the analyst.

### Overview of CARRECT system architecture



## Methods

Older approaches to imputing missing data, including mean imputation and single imputation regression methods, have steadily given way to a class of methods known as multiple imputation (MI).

Under MI, the missing values are drawn multiple times, generating m complete datasets along with the estimated parameters of the model (Rubin 1987; King et al. 2001). As a result, MI will lead to valid standard errors and confidence intervals along with unbiased point estimates. CARRECT uses a bootstrapping-based MI algorithm (Honaker and King, 2010) that gives essentially the same answers as the standard Bayesian Markov Chain Monte Carlo (MCMC) or Expectation Maximization (EM) approaches, is significantly faster, and can handle larger datasets.

## Conclusion

Our approach and program were designed to make bias mitigation much more accessible to much more than only the statistical elite. We hope that it will have a wide impact on reducing bias in epidemiological studies and provide more accurate information to policymakers.

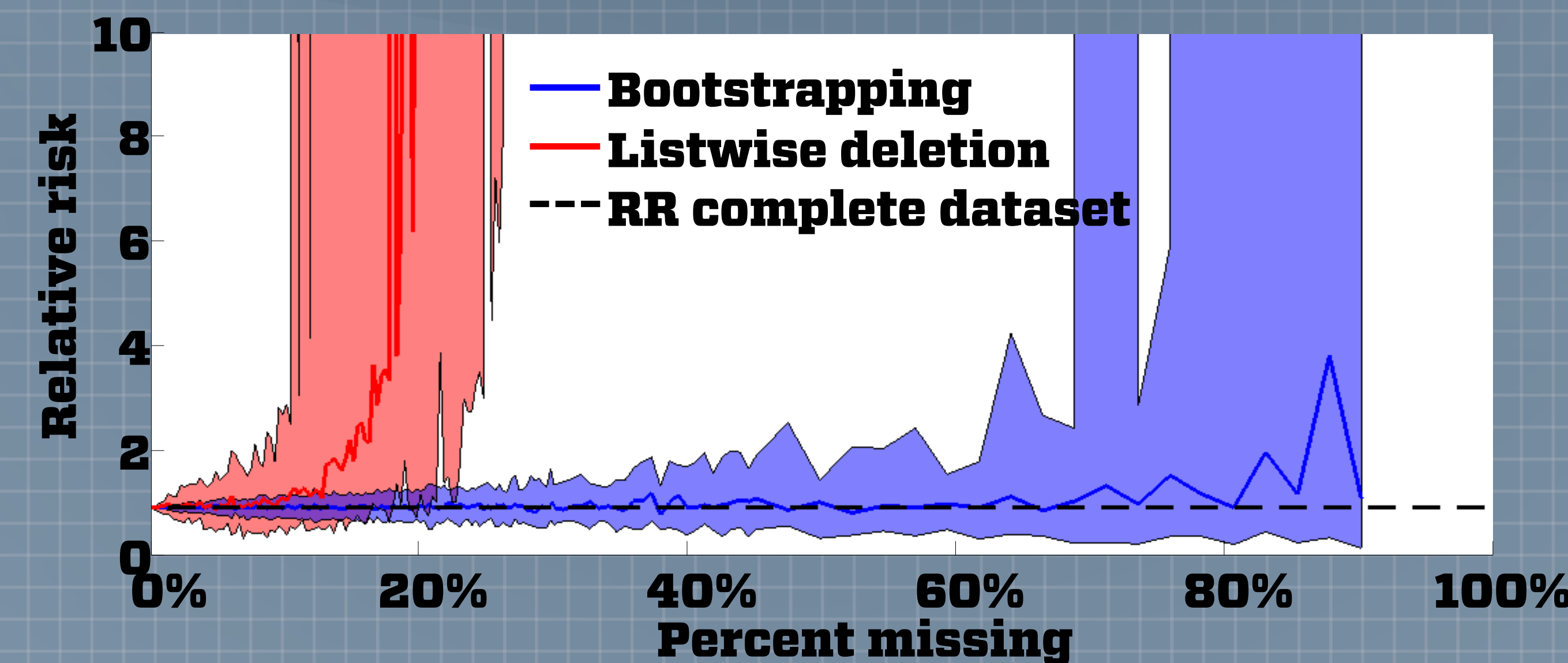**Contacts:** jlansey@aptima.com, ppicciano@aptima.com, iyohai@aptima.

## Results

Tests of the bootstrapping MI algorithm were conducted with an epidemiological dataset from the Integrated Health Interview Series (IHIS) producing verifiably unbiased results despite high missingness rates.

### Procedure:

- Created "complete" dataset. With 13 variables (selection shown at right).
- Removed increasing amounts of data completely at random
- Performed Listwise Deletion and Bootstrapping Multiple Imputation (MI)
- Calculated Relative Risk (RR) for coronary heart disease, given that the subject drinks.
- Repeated the imputation and RR calculation 50 times for each level of missingness from 1-90%

| Name | Description |
|---|---|
| AGE | Age |
| WEIGHT | Weight in pounds without clothes or shoes |
| BMI | Body mass index |
| HEIGHT | Height |
| ALC1YR | Ever had 12+ drinks in any one year |
| WALKFUNX | Times walked for leisure during past 7 days |
| ...... | ...... |
| CHEARTDIEV | Ever told had coronary heart disease |



The solid lines show the median RR over the 50 iterations for each missingness rate and the shaded regions show the 95% bounds for those 50 iterations. Bootstrapping MI provides results with greatly reduced bias, and has fewer computational demands than other methods.
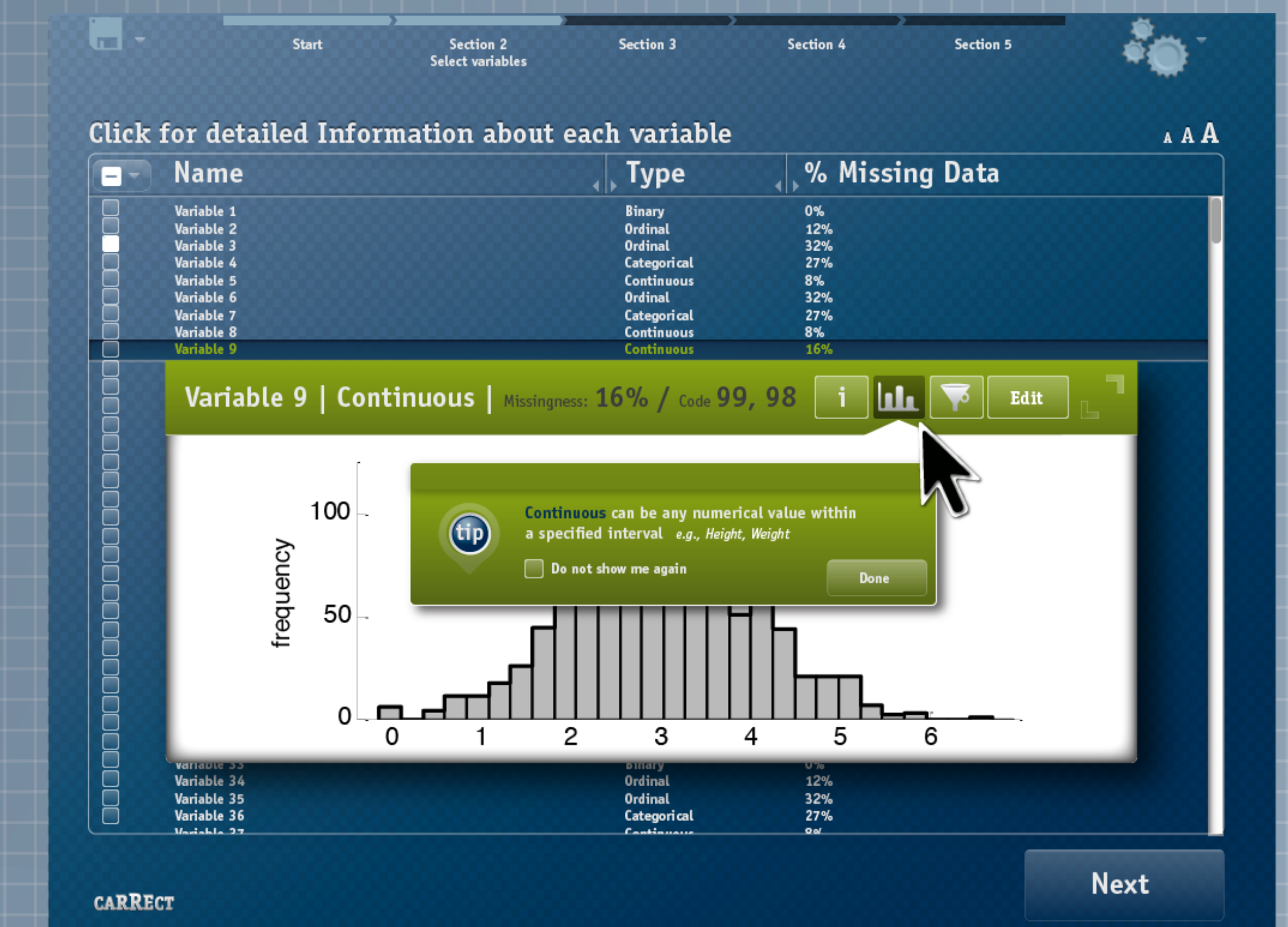
### Acknowledgements:

## User Interface Mockups

An intuitive data wizard guides the user through the analysis processes by analyzing key features of a given dataset. The wizard shows prompts the user to provide additional substantive knowledge to improve the handling of imperfect datasets.

The figure below shows how a user could select a variable type.



The figure below shows how a user is helped along with the selection of the most appropriate algorithms and models for the data