

report1

Nathaniel Brown, Huijia Yu, Angie Shen

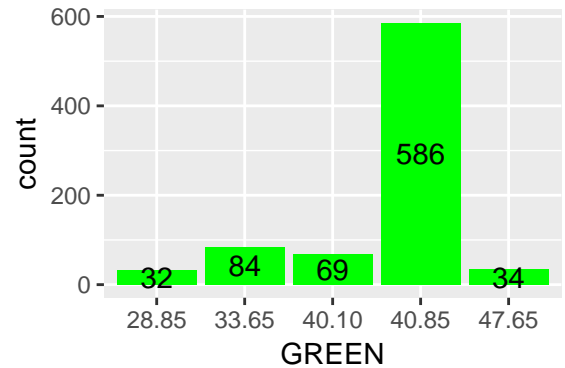
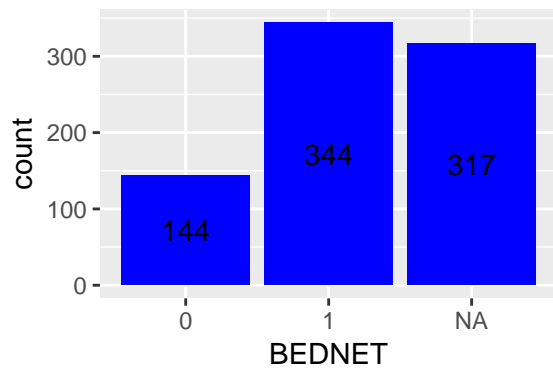
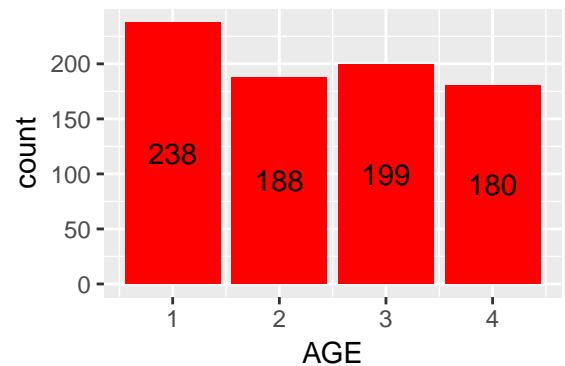
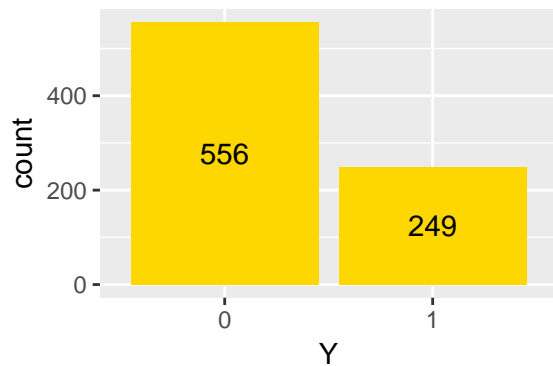
November 6, 2017

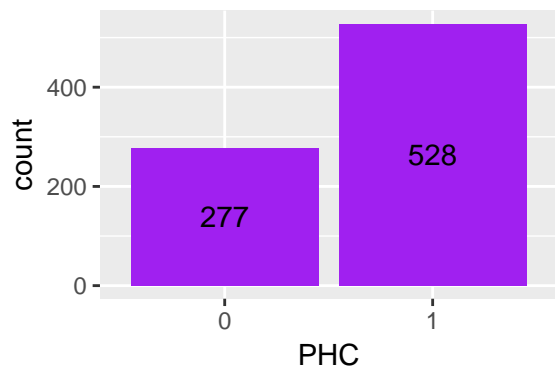
Introduction

We are interested in predicting whether malaria parasites will be found in a child's blood based on his/her age, bednet use, the amount of greenery in his/her village, and the presence of a health clinic. However, almost 40% of the observations do not have bednet use reported, so removing those observations would cause us to lose too much information. Instead, we will impute bednet values to predict the outcome.

Exploratory Data Analysis

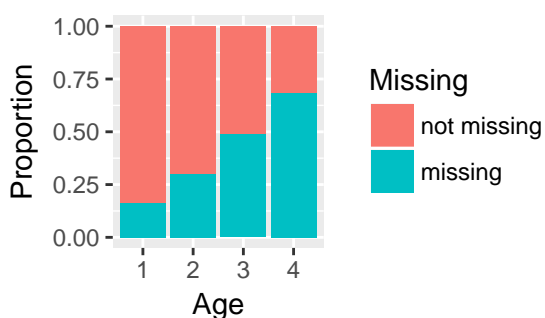
In the graphs below, we illustrate the frequency of each category within each the dataset:



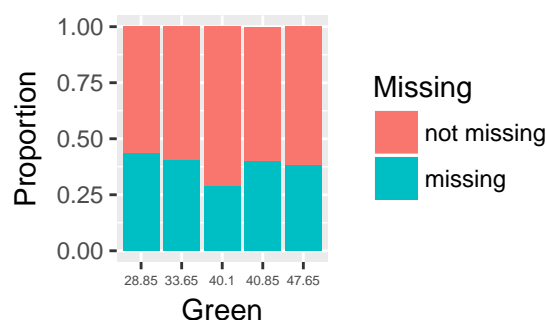


Below, we investigate the proportions of missing bednet responses by each level of the predictors.

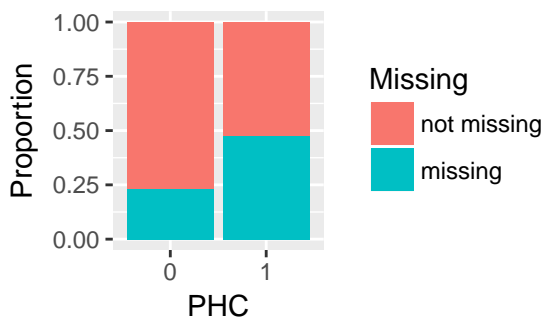
Age by Bednet



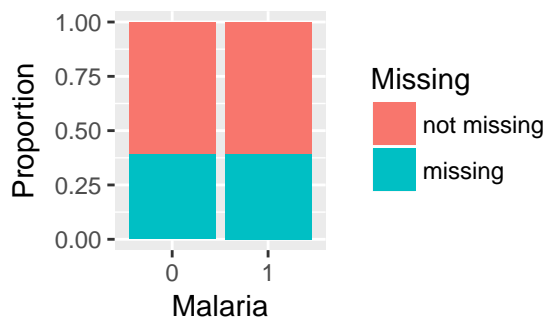
Greenery by Bednet



Health Clinic by Bednet



Parasites by Bednet



There appears to be a relationship between missingness in the bednet variable and age. As Age increases, the proportion of missing bednet responses increases. The other predictor variables and the response do not have an obvious visible relationship with missing bednet. We will use logistic regression to formally test the null hypothesis that the bednet data is missing completely at random (MCAR) versus the alternative that it is missing at random (MAR).

Test for Missingness Mechanism

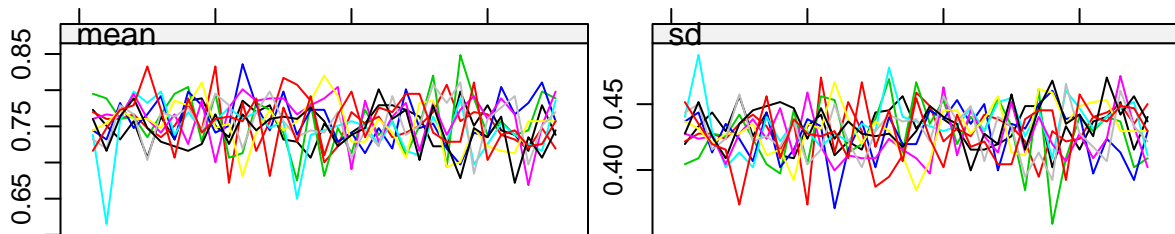
description of the three mechanisms

To test the assumption of the MCAR mechanism versus MAR, we build a logistic regression model to determine if our observed predictors (age, green, and public health clinic) are predictive of the missing bednet values. If this “full” model is not a better predictor than a “null” model with only an intercept, then we do not reject the null hypothesis that the missingness mechanism is MCAR.

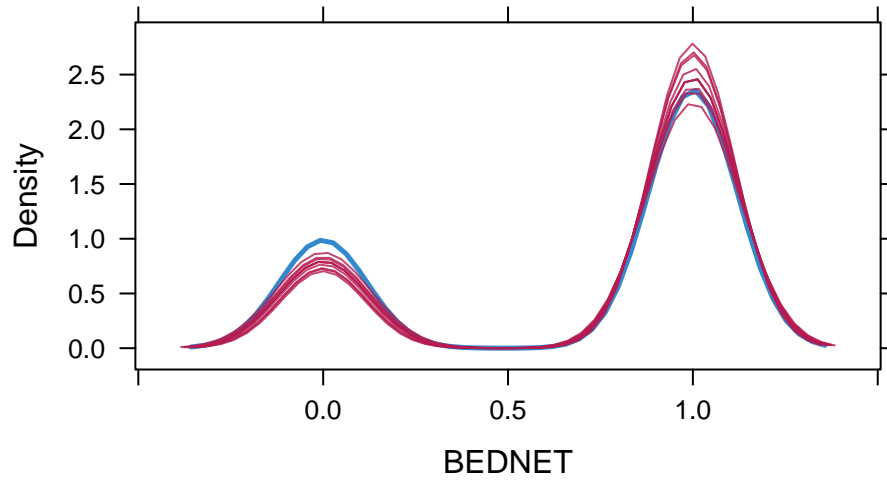
Based on a Chi-Squared test comparing the deviance of the full model to the deviance of the null model, there is extremely strong evidence that the missingness mechanism is MAR instead of MCAR ($p=0$). There is no way of knowing whether it is MAR or NMAR, since NMAR implies that the missing mechanism depends on unobserved factors.

Discussion of Approaches

Since the data is not MCAR, we cannot use bootstrapping within the bednet variable. We will instead use chained regression using the MICE (Multivariate Imputation by Chained Equations) package to impute the missing values and calculate malaria prevalence. **explain what chained regression is**



Iteration



The trace lines appear to be stationary and free of trends, indicating convergence. We can see that the density distribution of each of the imputed datasets (in red) is congruent with the original one (in blue).

	est	se	t	df	Pr(> t)	lo 95	hi 95	nmis
(Intercept)	-0.4715822	1.0021039	-0.4705921	651.52255	0.6380895	-2.4393251	1.4961607	NA
AGE	0.2621580	0.0691077	3.7934696	797.59464	0.0001598	0.1265035	0.3978125	0
GREEN	-0.0178805	0.0233223	-0.7666692	775.20987	0.4435117	-0.0636630	0.0279019	0
PHC	-0.4946417	0.1751651	-2.8238598	572.52042	0.0049101	-0.8386863	-0.1505971	0
BEDNET	0.0713020	0.2553208	0.2792643	34.18665	0.7817246	-0.4474679	0.5900719	317

After fitting a logistic regression to each of the 10 generated datasets, we can see that the bednet variable is actually not significant at the $\alpha=0.05$ confidence level.

Contributions

Nathaniel made the Exploratory Data Analysis plots on the missingness of the bednet variable, and wrote the observations. Huijia worked on the Discussion of Approaches section.