

HW5_YQ

Youqing Xiang

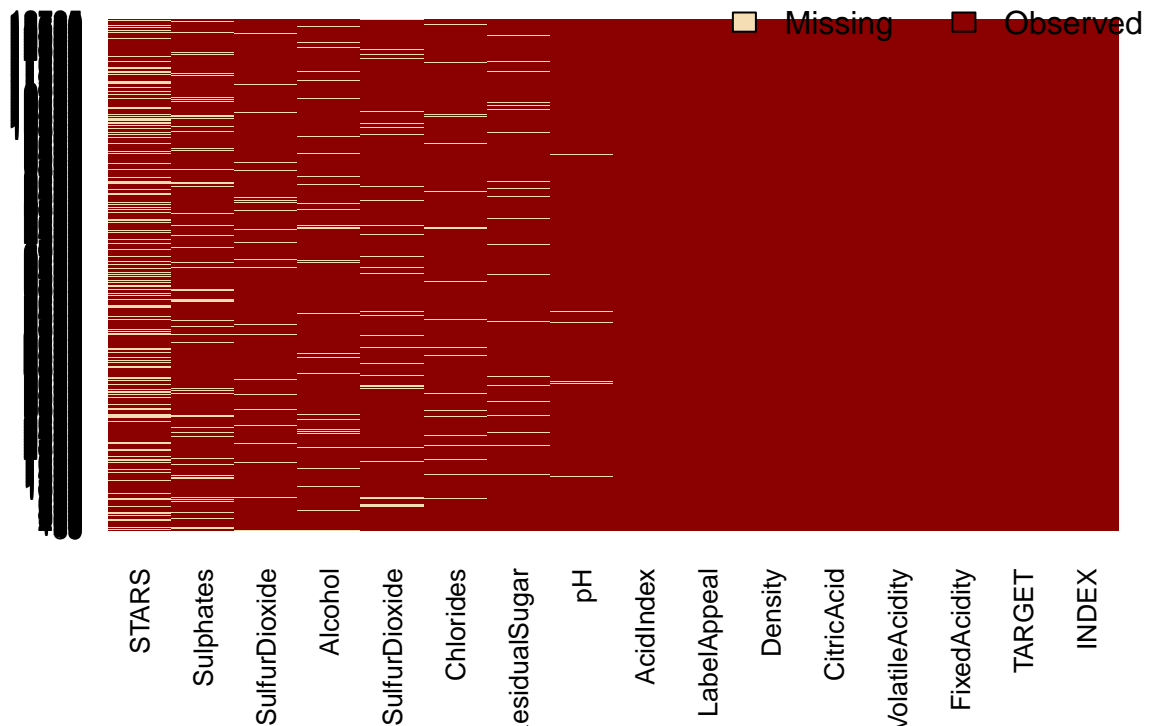
July 16, 2016

Data Exploration

Variable	NAs	Percent_NAs
TARGET	0	0.000
FixedAcidity	0	0.000
VolatileAcidity	0	0.000
CitricAcid	0	0.000
ResidualSugar	616	0.048
Chlorides	638	0.050
FreeSulfurDioxide	647	0.051
TotalSulfurDioxide	682	0.053
Density	0	0.000
pH	395	0.031
Sulphates	1210	0.095
Alcohol	653	0.051
LabelAppeal	0	0.000
AcidIndex	0	0.000
STARS	3359	0.263

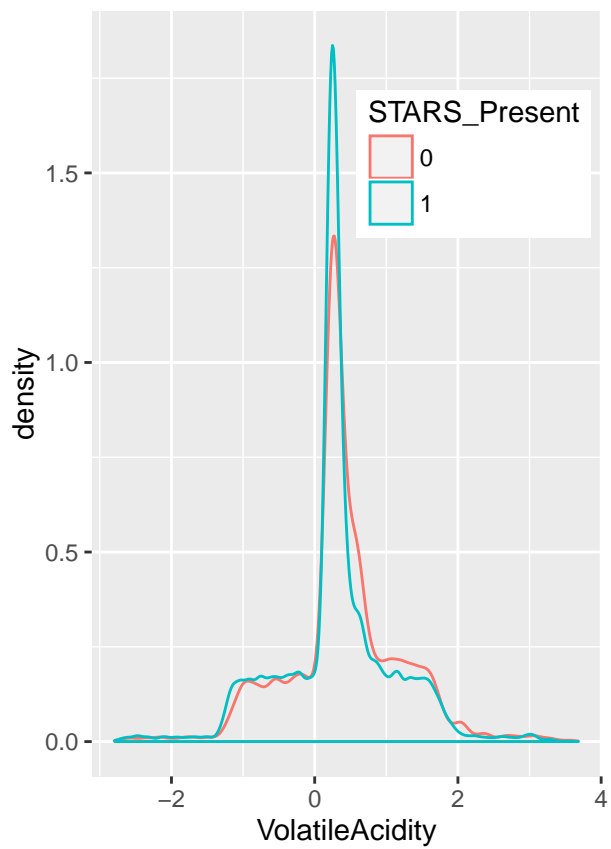
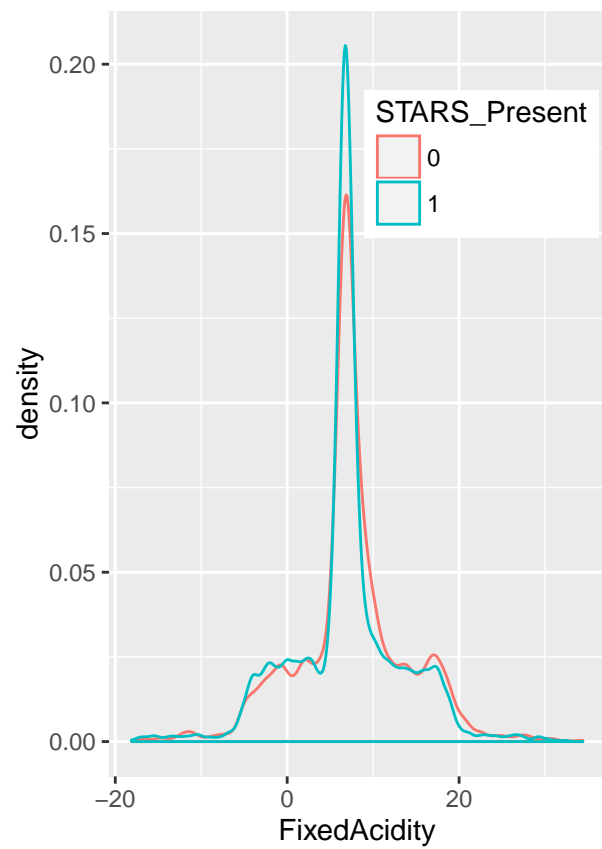
This dataset includes 12795 observations with 15 variables (**Index** columns excluded) in total. As the table shows above, there are a fair number of NAs and **STARS** variable has the most NAs.

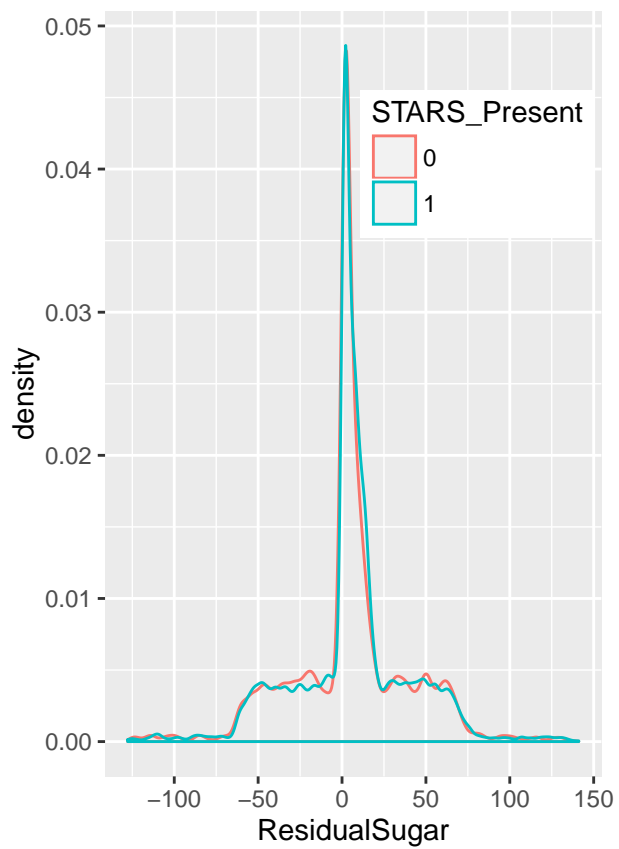
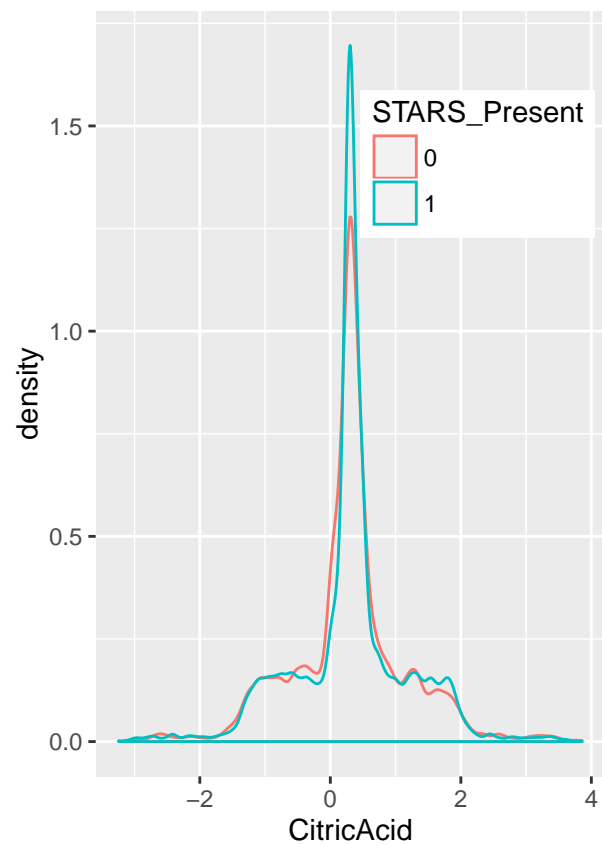
Missingness Map

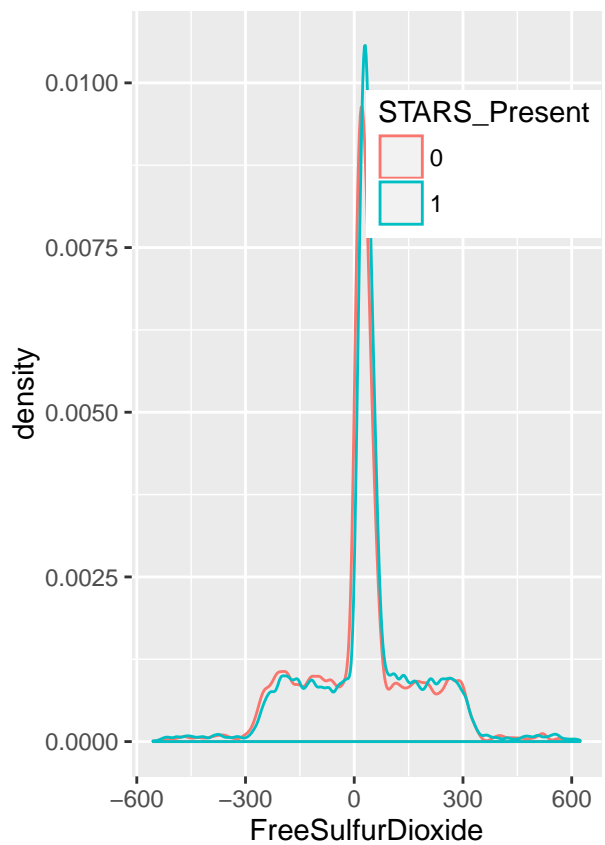


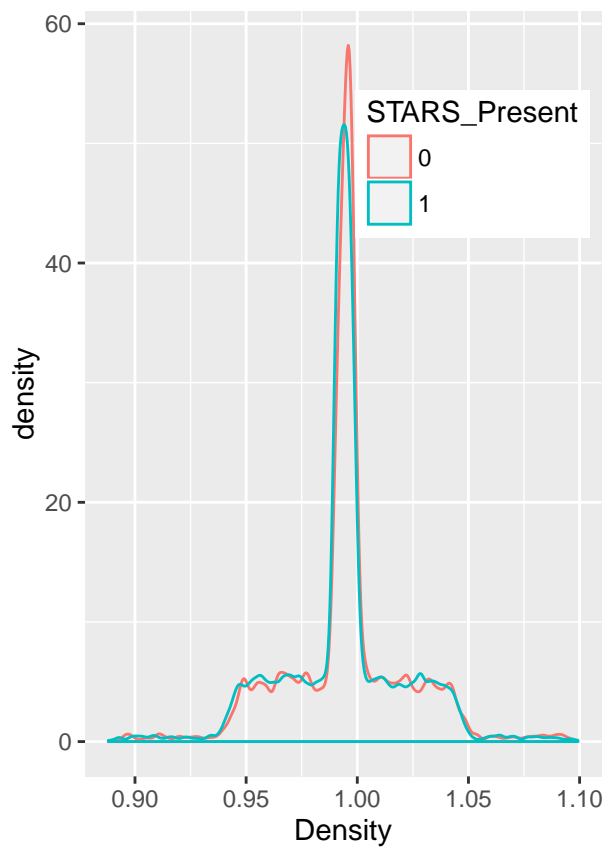
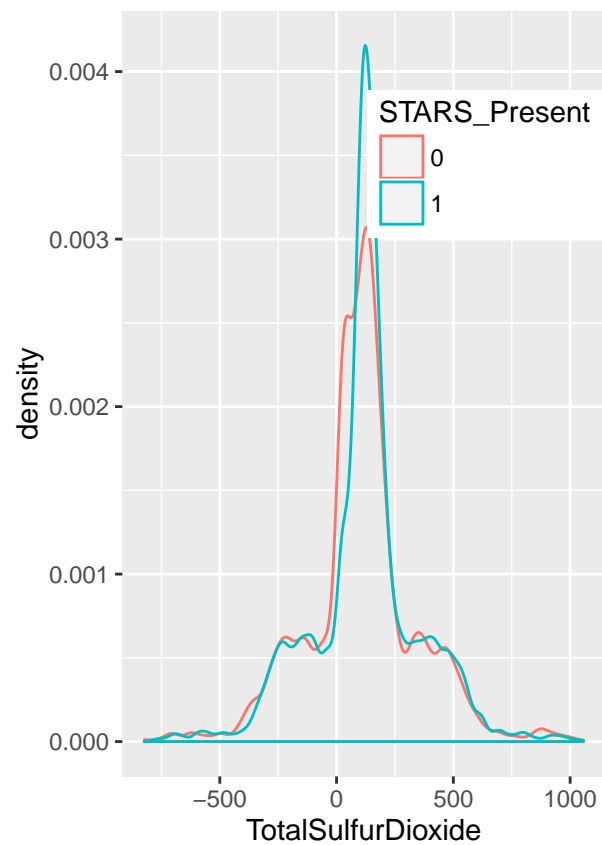
The above matrix of the NAs doesn't show any pattern. So, we could consider to replace the NA values with certain method to avoid potentially losing a large amount of data if we just simply drop off the NA values during data preparation process.

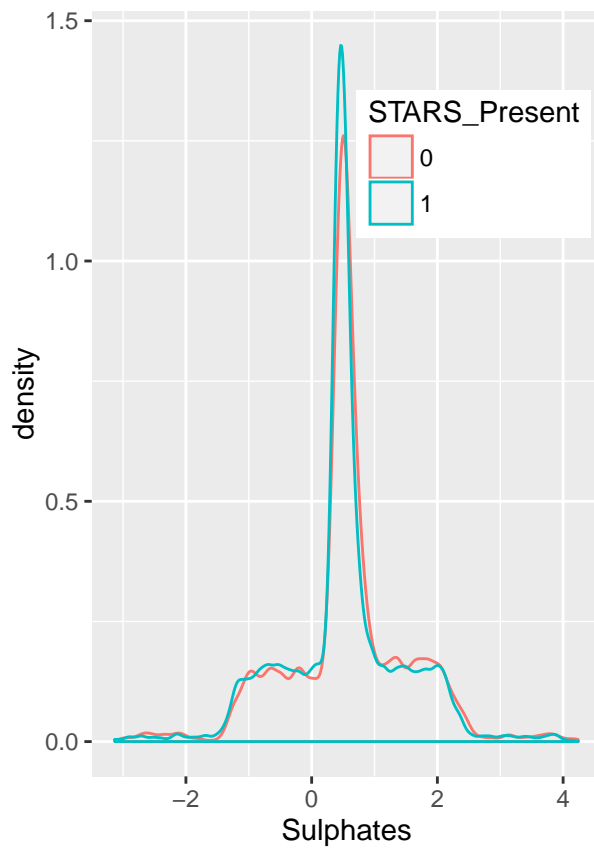
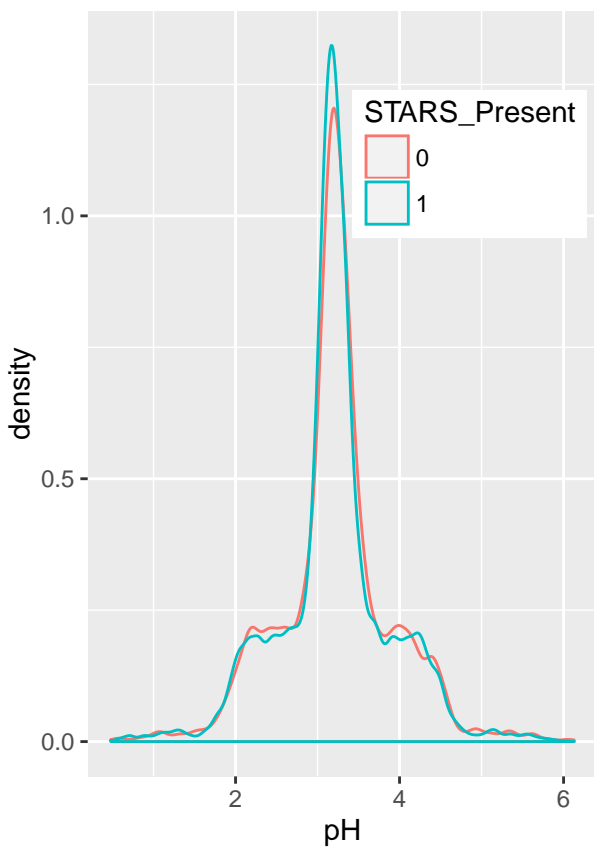
Since **STARS** variable has the most NAs, which accounts for 26.3% of data points, it is worthwhile to check how **STARS** variable NAs affect other variables, especially the **TARGET** variable. Here I created a new variable: **STARS_Present**, which is a categorical variable, equals to 0 when **STARS** value is missing and 1 when **STARS** value is present. And then I show density plots for each variable grouped by **STARS_Present**.

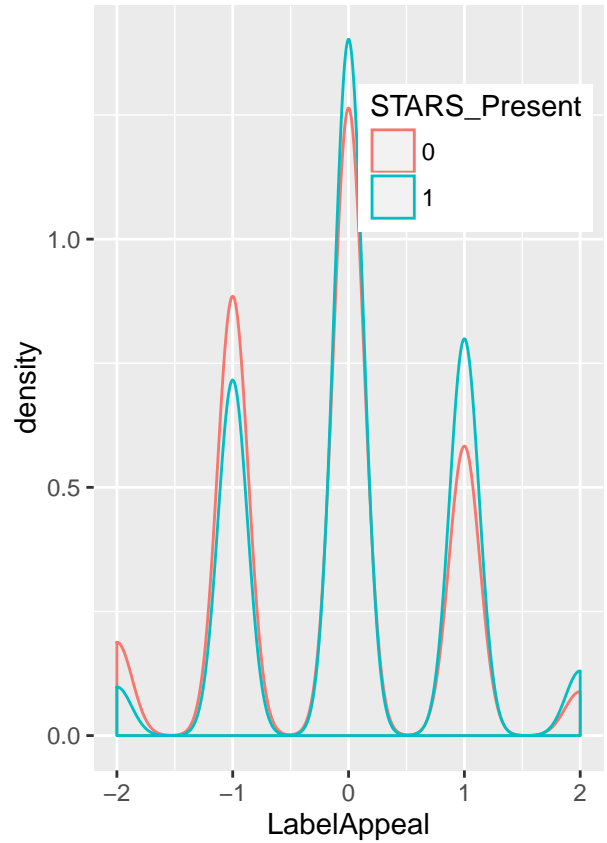
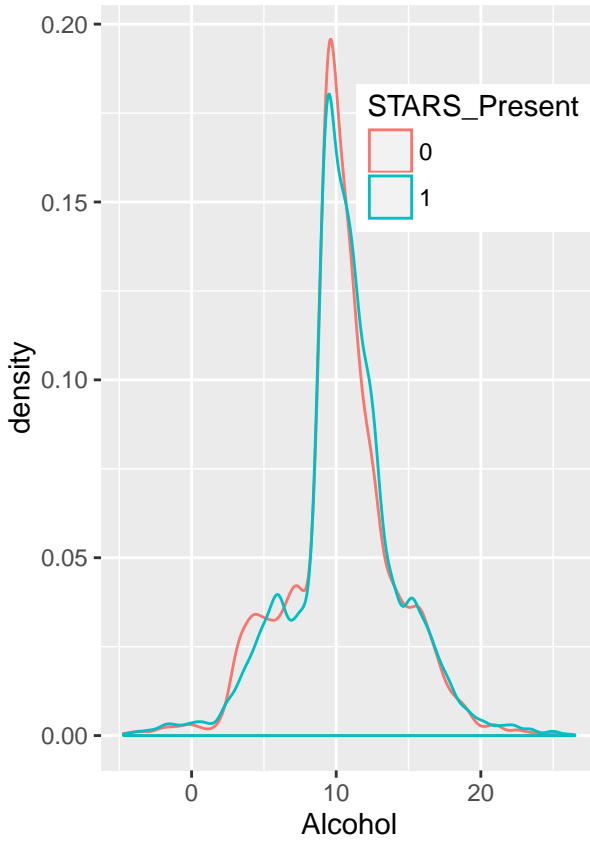


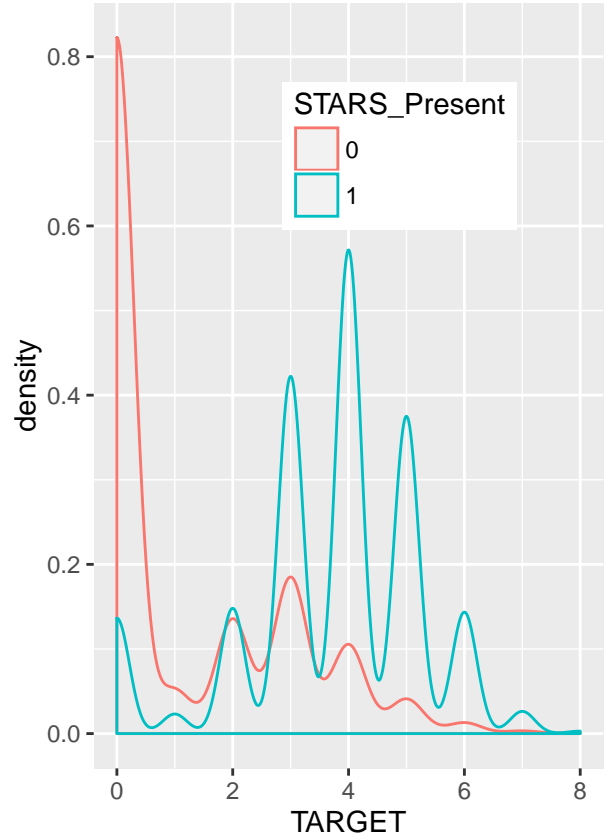
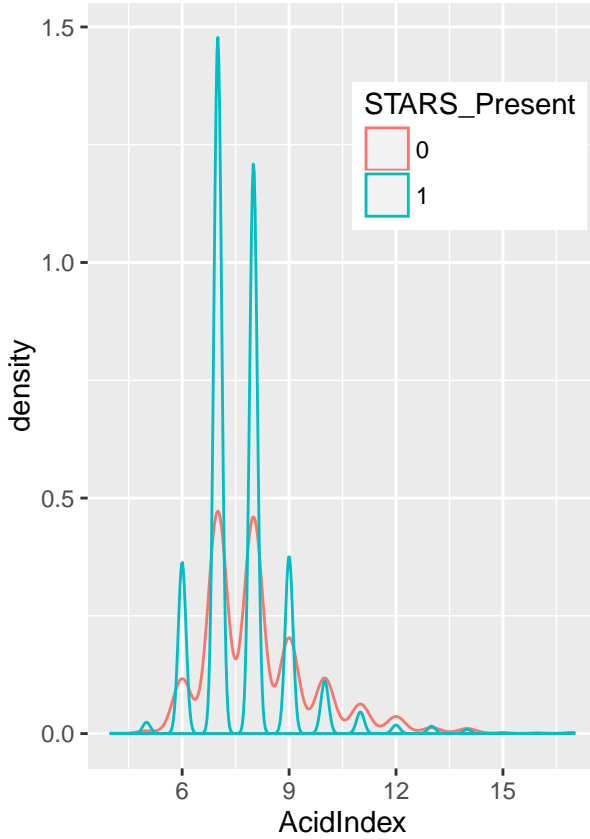












From the above density plots, we can see that **STARS** variable NAs have big effect on **TARGET** distribution but no obvious effects on other variables. And we can also conclude that **STARS** variable NAs is actually predictive of the target. Overall, this part of analysis suggests us that we should treat **STARS** variable NAs as category variable.