

Social Effects on Drug Use

Jesse Loi

Introduction



- The National Survey on Drug Use and Health (NSDUH) contains data relating general features like education to drugs use.
- We will investigate how social factors may impact marijuana use in youth
 - Social factors: opinions of peers, parental support

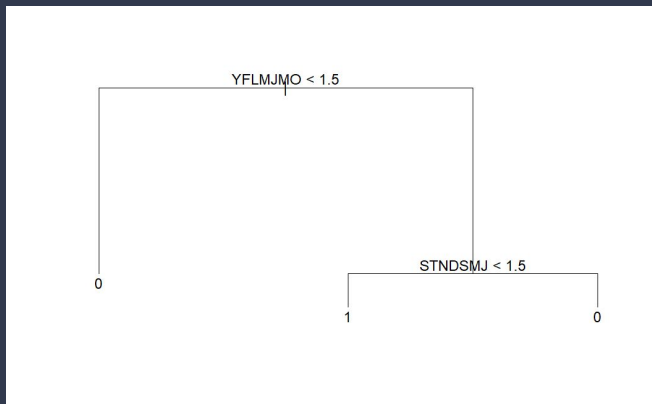
Models:

- Simple classification and regression decision trees
- Bootstrapped Aggregated (Bagging) Decision Trees
- Random Forest Decision Trees
- Boosted Decision Trees

Goal:

Minimize Mean-Squared Error and Misclassification Rate

Theoretical Background



Simple Decision Tree

Bootstrap Aggregating:

Re-sample data many times and build many trees

Flaw: one impactful variable may dominate (try Random Forest)

Random Forest:

Avoids dominance of a single feature by growing many trees with different sets of features

Boosting:

Train our data on errors of previous data

Flaw: tends to overfit (cross validate)

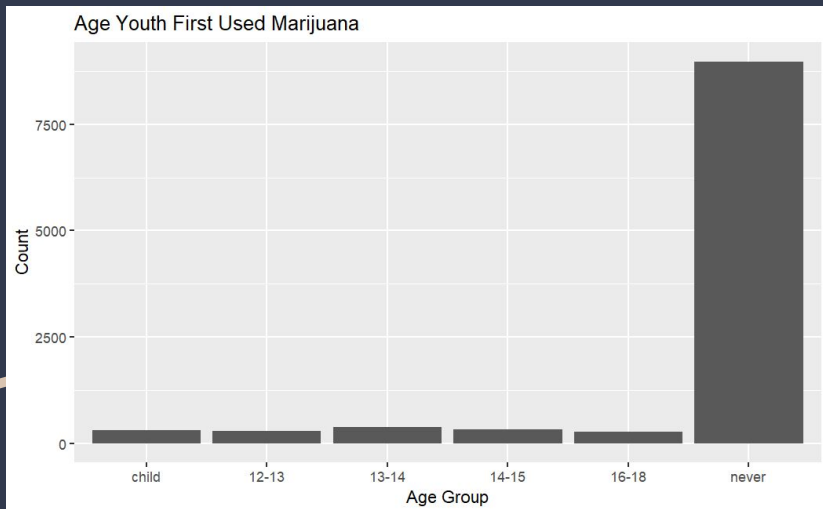
Methodology: Cleaning

Len : 3 ALCOHOL FREQUENCY PAST YEAR - IMPUTATION REVISED

RANGE = 1 - 365

991 = NEVER USED ALCOHOL

993 = DID NOT USE ALCOHOL PAST YEAR



Data Cleaning

- Convert variables to categorical
- Change Extreme Values to 0

Problem: Unbalanced Data

Solution: omit "never" entries and investigate only those who have used marijuana

Methodology: Questions for Each Model Type

Binary Classification:

How can we classify youth into having used marijuana and never having taken marijuana.

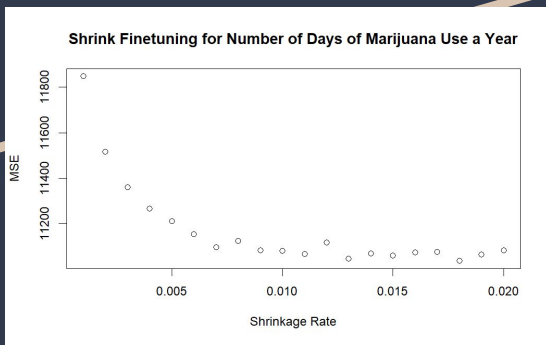
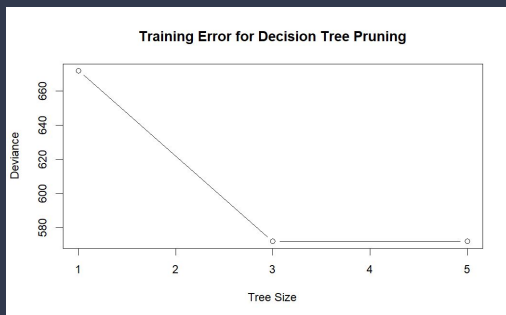
Multiclass Classification:

How can we determine when a marijuana user first used marijuana?

Regression:

Can we predict how many days a year a youth will have used marijuana?

Methodology: Finetuning (Simple, Boost)



For simple trees: try out different tree sizes to avoid overfitting

- A tree size of 3 was the best

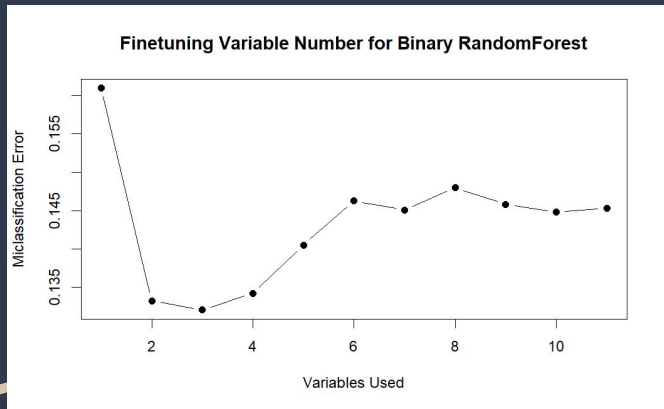
For Boosting: we change the shrinkage rate, which is the rate at which our model updates to accommodate errors.

- Increasing the shrinking parameter did not decrease test MSE by much, but we chose an elbow point of 0.011

Methodology: Finetuning (Bagging, Random Forest)

For Random Forest:

We finetune the number of variables used, ending with 3, which is very close to the size of the tree for our simple tree used before.



Results

Binary Classification:

- Predicts marijuana use with an error rate of 12.4%
- Important Variables: opinion of peers' marijuana use, number of peers using marijuana

Multiclass Classification:

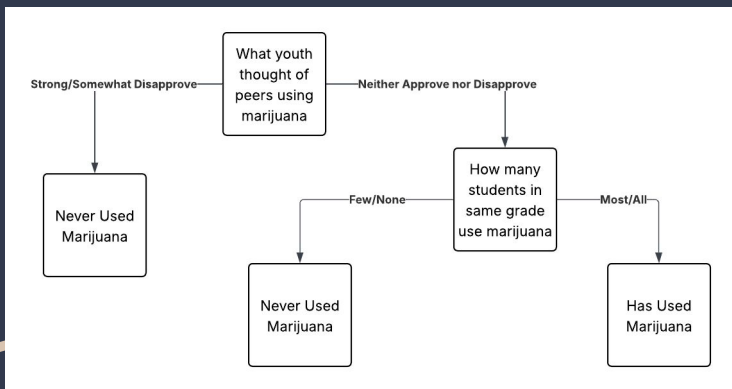
- Predicts age of first marijuana use with an error rate of 70%
- Important Variables: age of alcohol first use

Regression:

- Predicts with a mean squared error of 10000 days.
- Important Variables: days per year using alcohol, whether parents say they are proud, and monthly cigarette usage

Discussion: Social Effects Based on Variable Type

A simple decision tree



Binary Classification:

- Tree size of 3 is enough (3 terminal nodes)
- Social factors related to peers are the most impactful

Multiclass Classification:

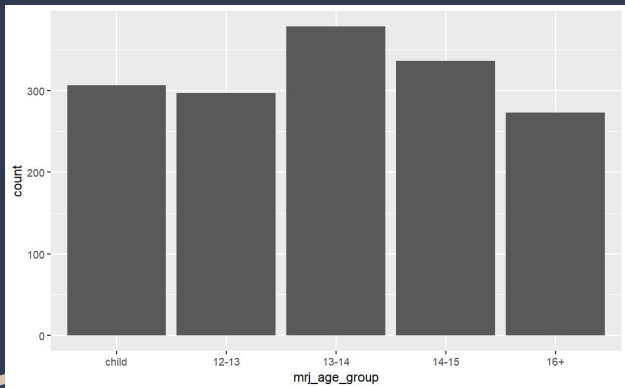
- Very difficult to determine when youth started smoking by age group.
- Most impactful factors relied on collinearity with other drug use variables (generally unreliable).

Regression:

- Frequency of drug use per year all can predict each other
- Parental intervention surprisingly helpful

Discussion: Binary, Ordinal, and Numerical Variables

Ideal for treating as ordinal



Example: Days used marijuana during the past year

Binary: categorize days into those above 182 and under or equal to 182, entailing “high” or “low” usage

Ordinal: categorize days into buckets

Numerical: have variables from 0 to 365 representing

When to use each:

Binary: when you have little or unbalanced data

Ordinal: when there aren't enough data for regression

Numerical: when there is plenty of data to work with

Conclusion

Peer programs may be far more effective at preventing marijuana use in youth

Hard to identify risk factors for when a teen will start to use marijuana (could be many ages)

- More work can be done to collect data to predict the age.

High usage of one drug tends to imply high usage of another drug.

- Can make anti-drug campaigns more transferable
- Parents cannot stop children from first using drugs, but can lower the frequency

References

1. Ripley B (2024). `_tree`: Classification and Regression Trees_. R package version 1.0-44,<<https://CRAN.R-project.org/package=tr>>.
2. Liaw A, Wiener M (2002). "Classification and Regression by randomForest." *_R News_*, *2*(3),18-22.<<https://CRAN.R-project.org/doc/Rnews/>>.
3. Ridgeway G, Developers G (2024). `_gbm`: Generalized Boosted Regression Models_. R package version 2.2.2, <<https://CRAN.R-project.org/package=gbm>>.
4. H. Wickham. `ggplot2`: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016.