

Module 1
Lesson 1 - What Is Machine Learning?

Machine Learning

Describe Associations

Search For Patterns

Make Predictions



Accuracy = Test Error Rate

Goal: Find a model that minimizes test error rate

Module 1

Lesson 2 - Machine Learning and the Bias Variance Trade-Off

Linear Regression

Accuracy = mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

Variance = change in parameter estimates across different data sets

Bias = how far off model estimated values are from true values

Low Variance ✓

Low Bias ✓

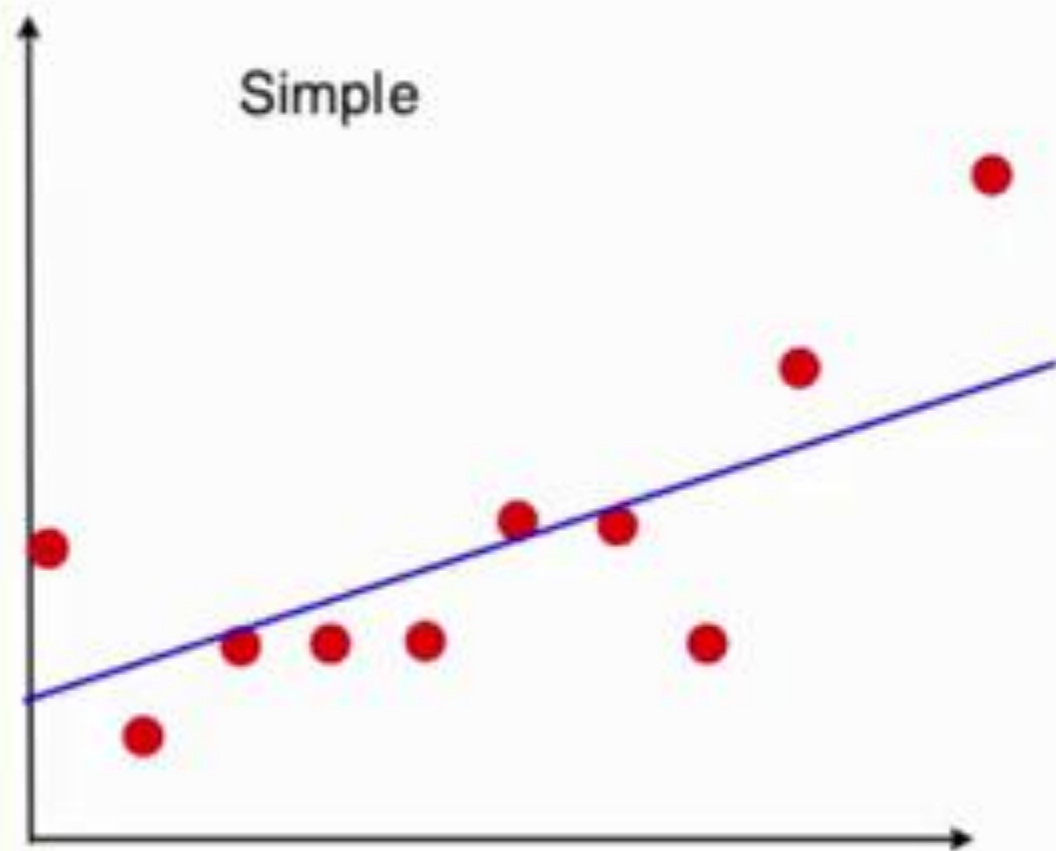


Complexity = High Variance/Low Bias



Complexity = Low Variance/High Bias

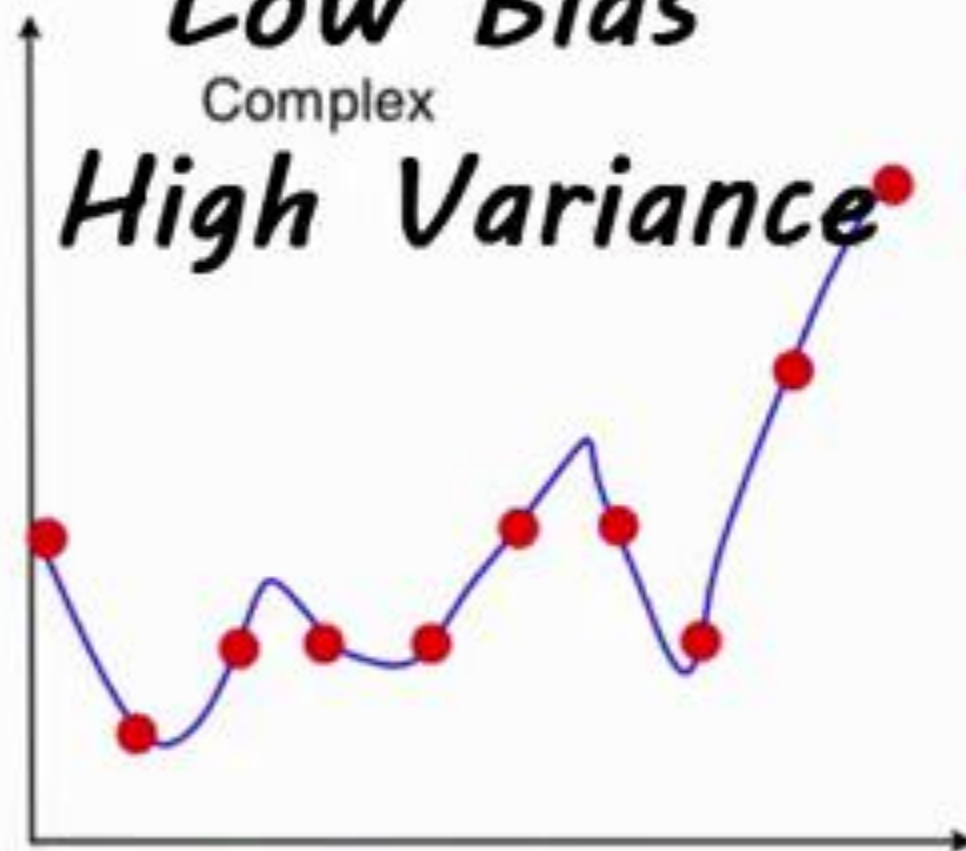
Simple



Low Bias

Complex

High Variance

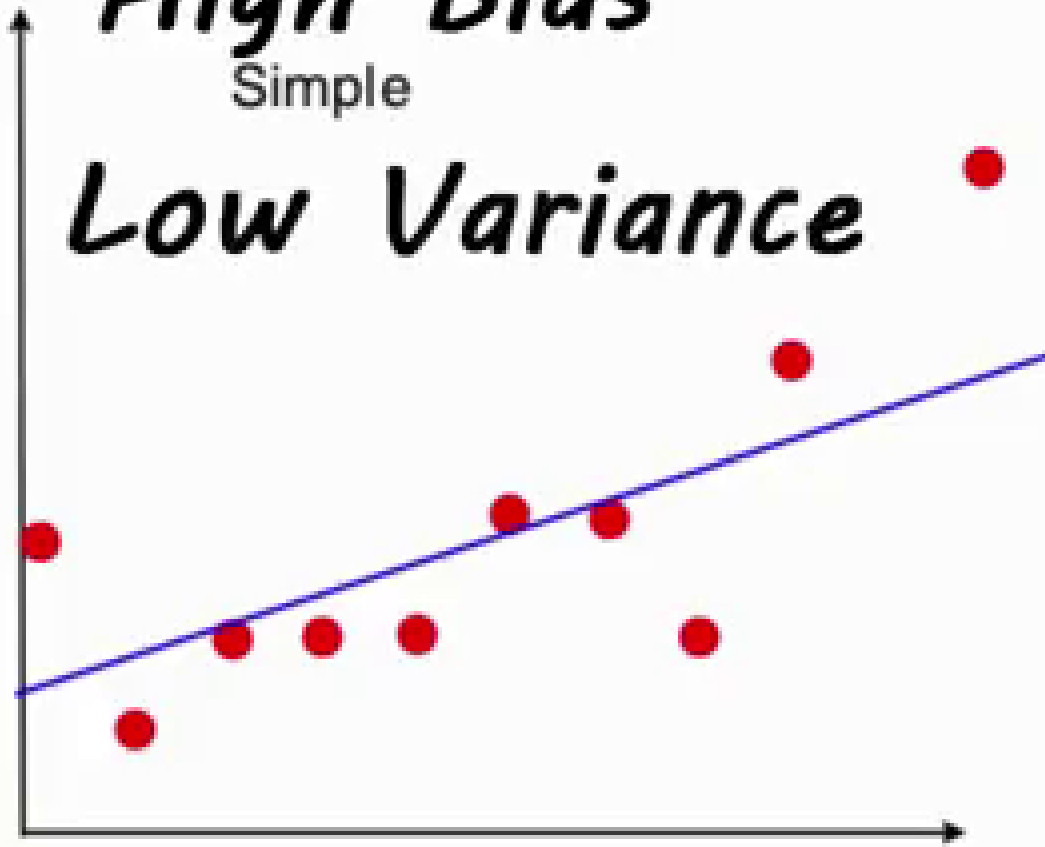


Overfitted

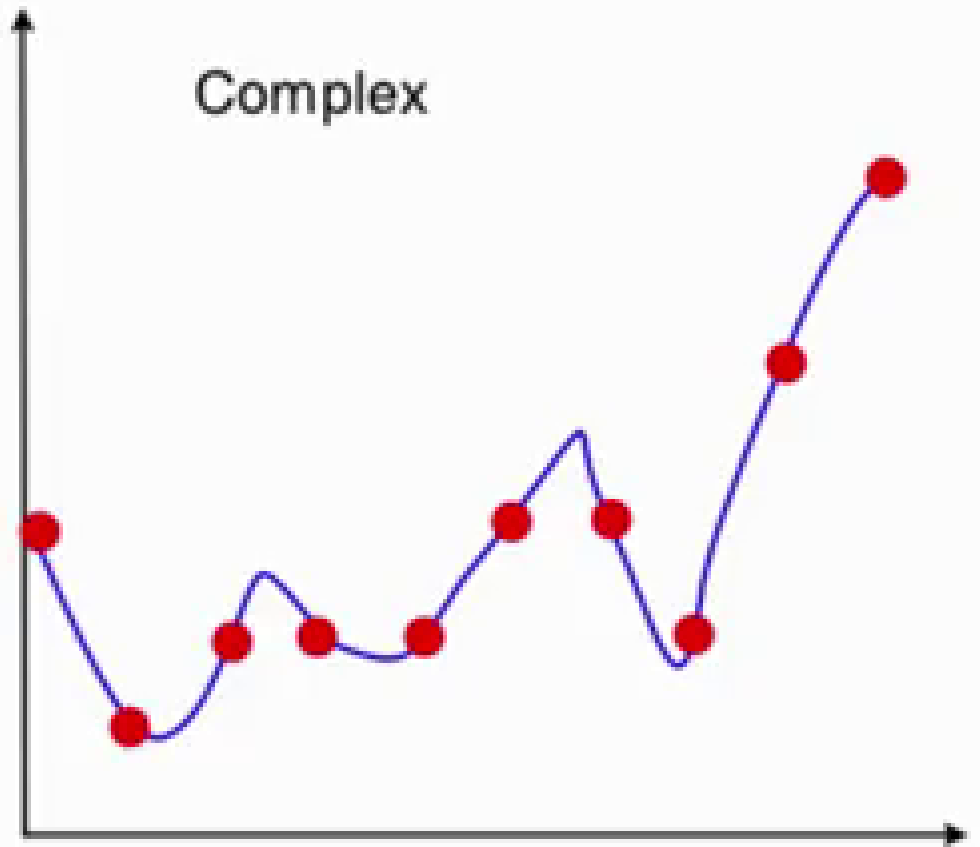
High Bias

Simple

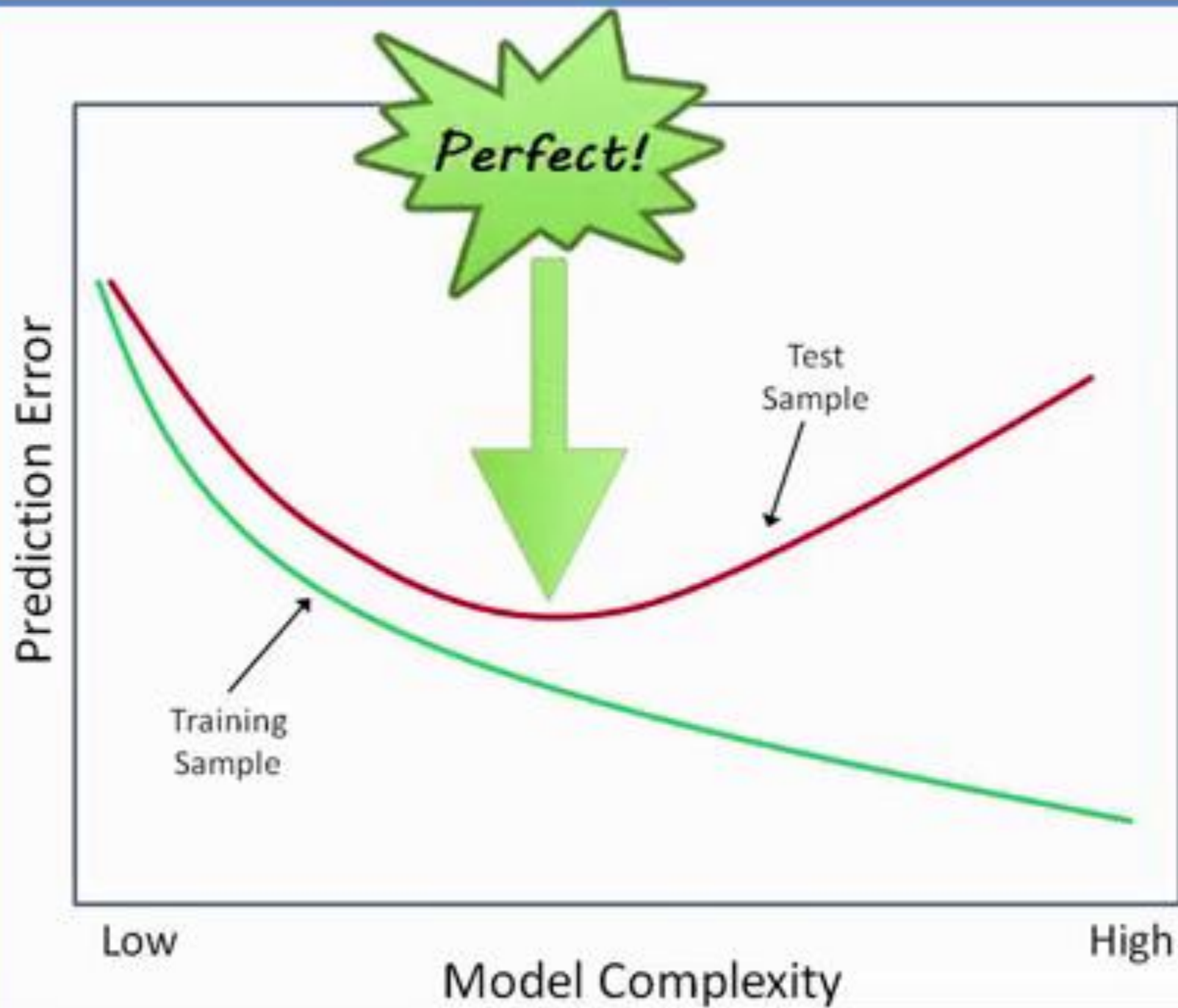
Low Variance



Complex



Underfitted



Logistic Regression

*Accuracy = How well a model correctly
classifies observations*





Confusion Matrix

Test Sample Nicotine Dependence Classification

		Actual		
		Yes	No	Total
Predicted	Yes	184	91	275
	No	32	685	717
	Total	216	776	992

Confusion Matrix

Test Sample Nicotine Dependence Classification

		Actual		
Predicted		Yes	No	Total
	Yes	184 	91 	275
	No	32 	685 	717
	Total	216	776	992

Confusion Matrix

Test Sample Nicotine Dependence Classification

		Actual		
		Yes	No	Total
Predicted	Yes	184	91	275
	No	32	685	717
Total		216	776	992

91 + 32 = 123 incorrectly classified

Confusion Matrix

Test Sample Nicotine Dependence Classification

		Actual		
		Yes	No	Total
Predicted	Yes	184	91	275
	No	32	685	717
Total		216	776	992

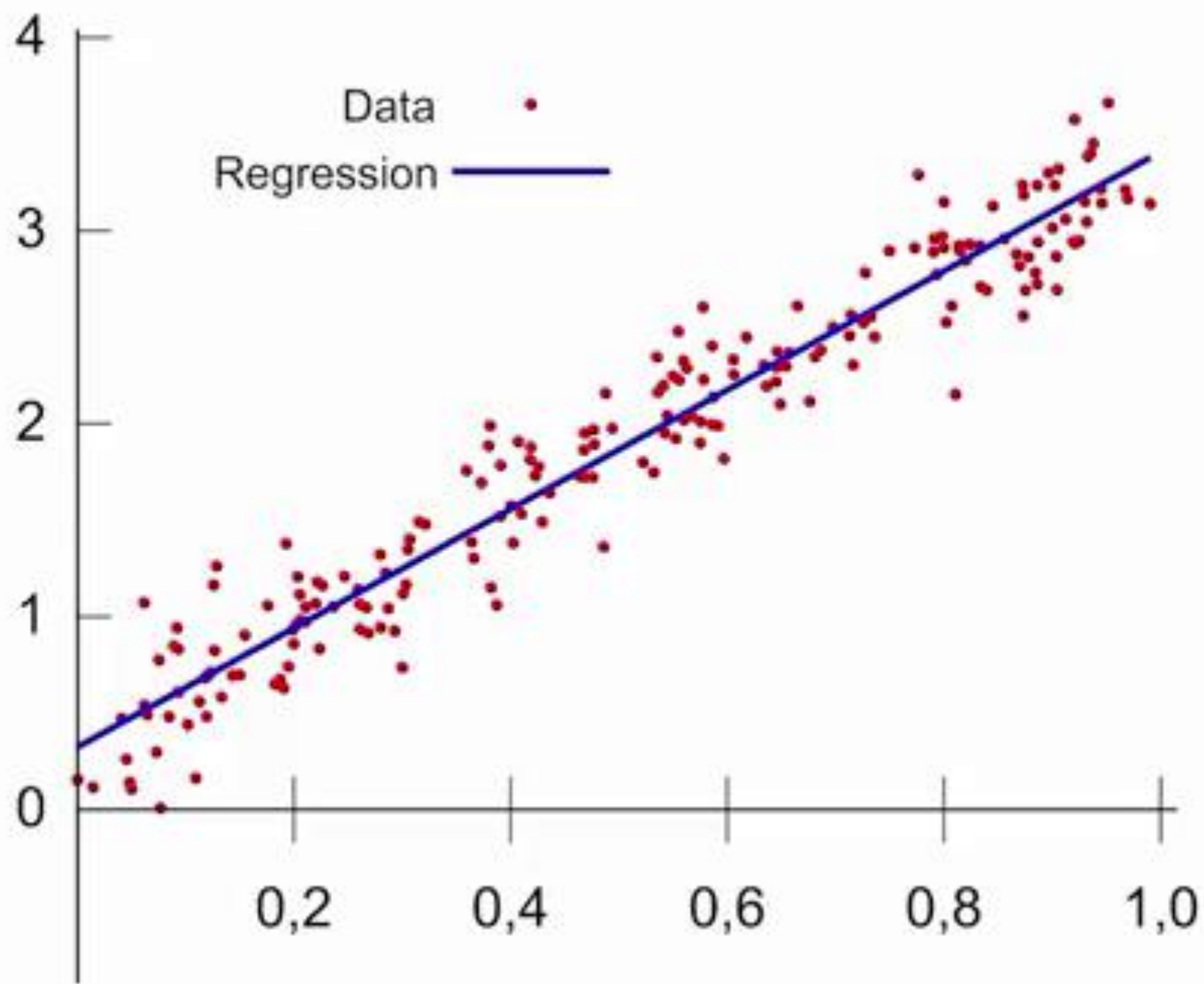
Test error rate = % misclassified = 12%



Decision Trees

What is a Decision Tree?

with Professor Lisa Dierker



Supervised Prediction

linear regression

pattern recognition

discriminant analysis

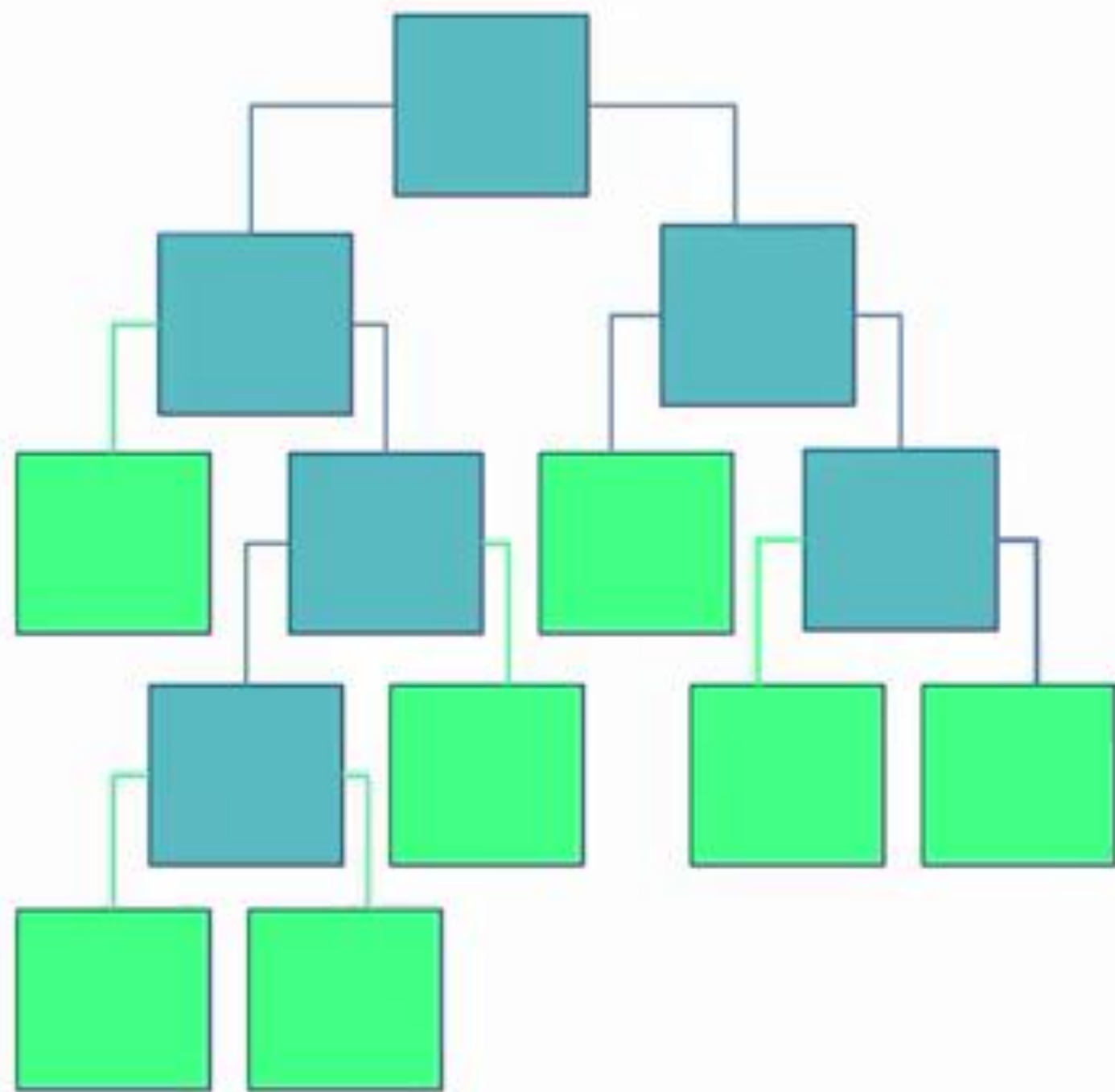
multivariate function estimation

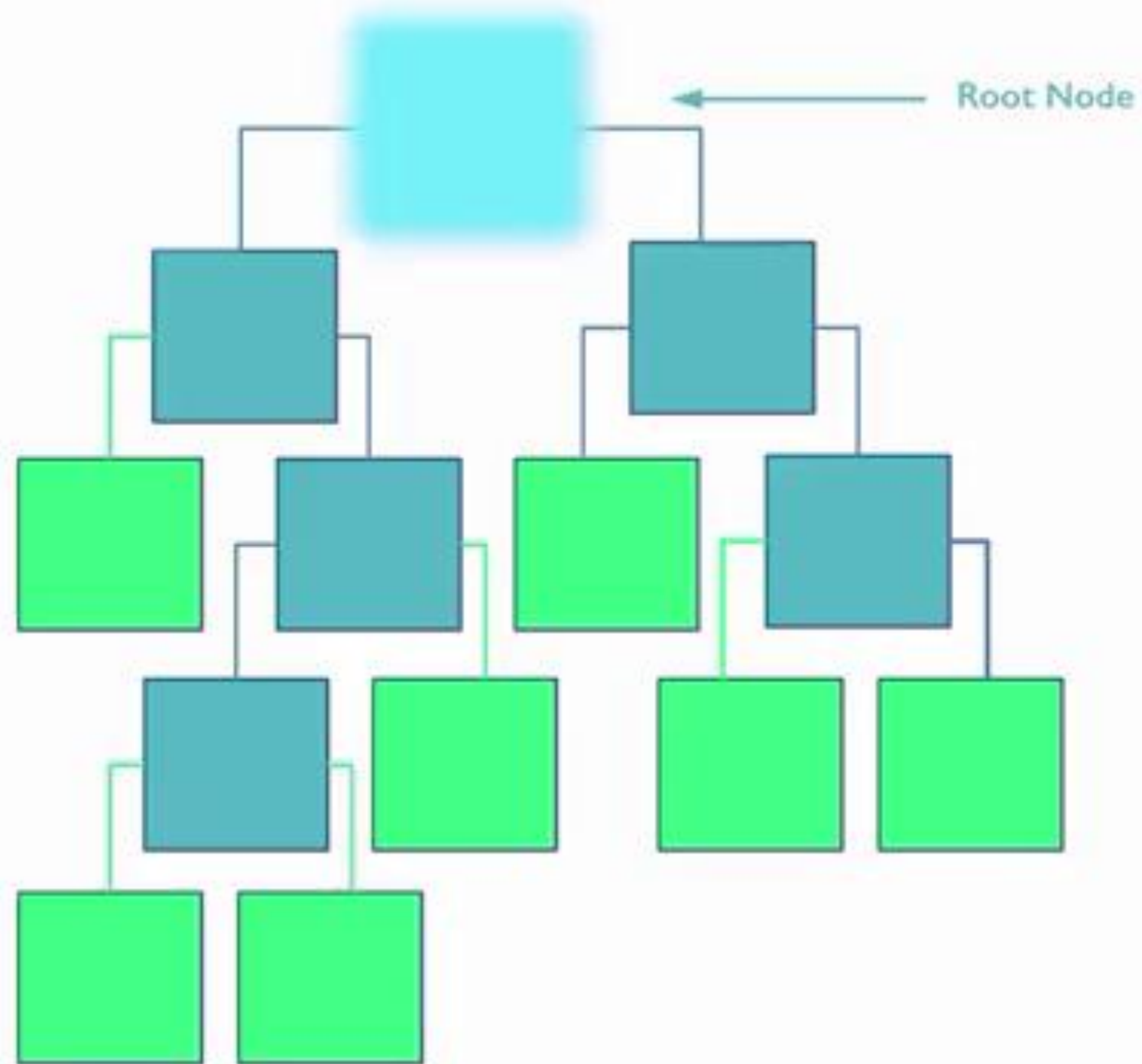
supervised machine learning techniques

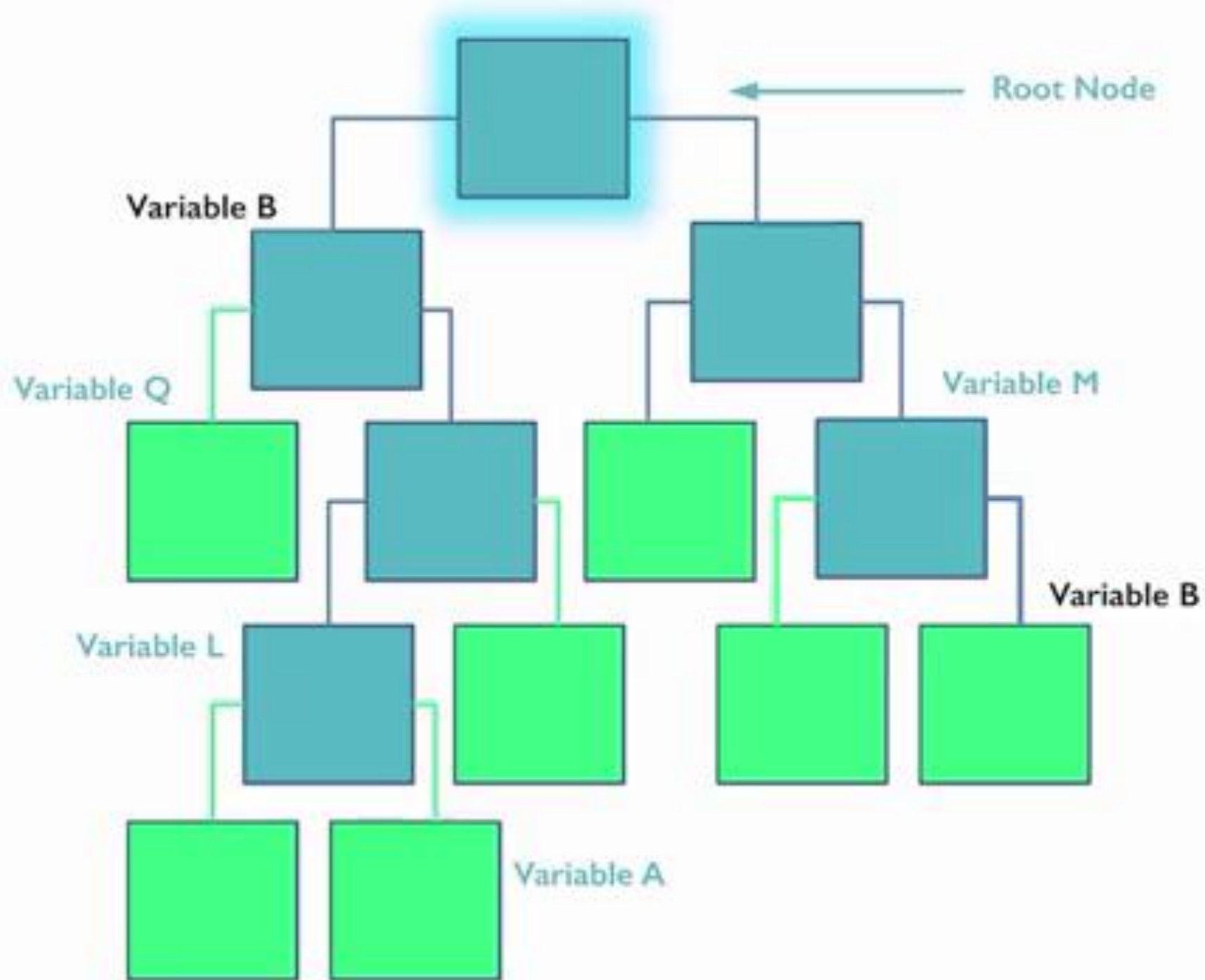
Supervised Prediction

explanatory variables

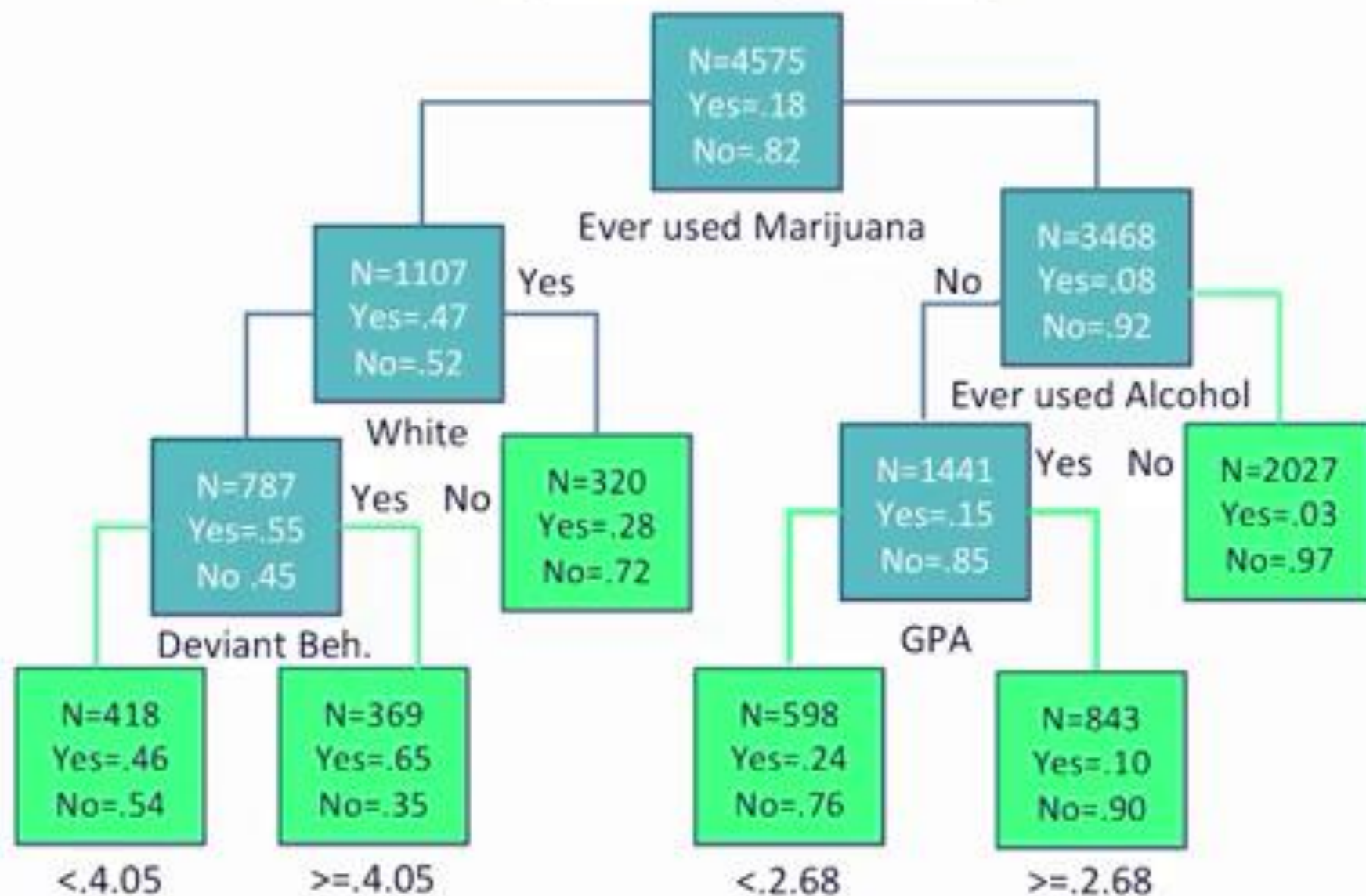
response variable







Target Variable: Regular Smoking



24 possible explanatory variables

target variable - regular smoking

target variable - regular smoking

"Have you (ever) smoked cigarettes regularly, that is, at least one cigarette every day for 30 days?"

24 possible explanatory variables

**associated with regular smoking
behavior in adolescents**

Binary Categorical Variables

Gender

Hispanic

White

Black

Native American

Asian

Substance use measured with individual questions

Alcohol
Marijuana
Cocaine
Inhalants

Additional Categorical Variables

Availability of Cigarettes
Parents on Public Assistance
Expelled from School

Quantitative Variables

Age

Alcohol Problems

Deviant Behavior

Violent Behavior

Depression

Self Esteem

Parental Presence

Parental Activities

Family Connectedness

School Connectedness

Grade Point Average

For more complete details on how these variables were constructed see:

**Dierker, et al., 2004 paper from Prevention Science
SAS program: "Decision Trees Data Management"**

Decision Trees

What is the Process of
"Growing" a Decision Tree?

with Professor Lisa Dierker

Growing the Tree

Binary splits maximize correct classification

All cut-points are tested

**Subgroups showing similar outcomes
are generated**

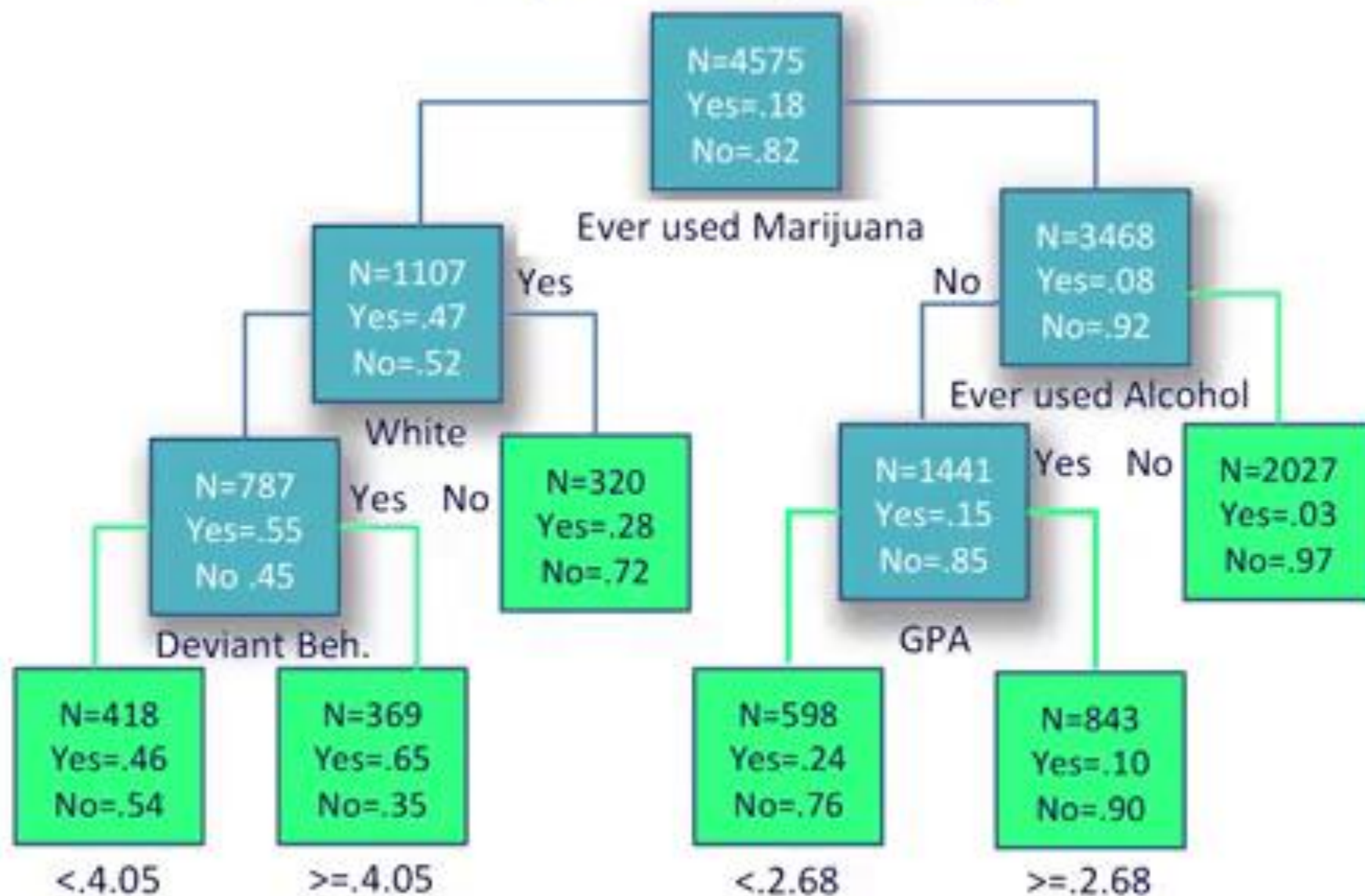
Validating the Tree

Cross-validation guards against overfit

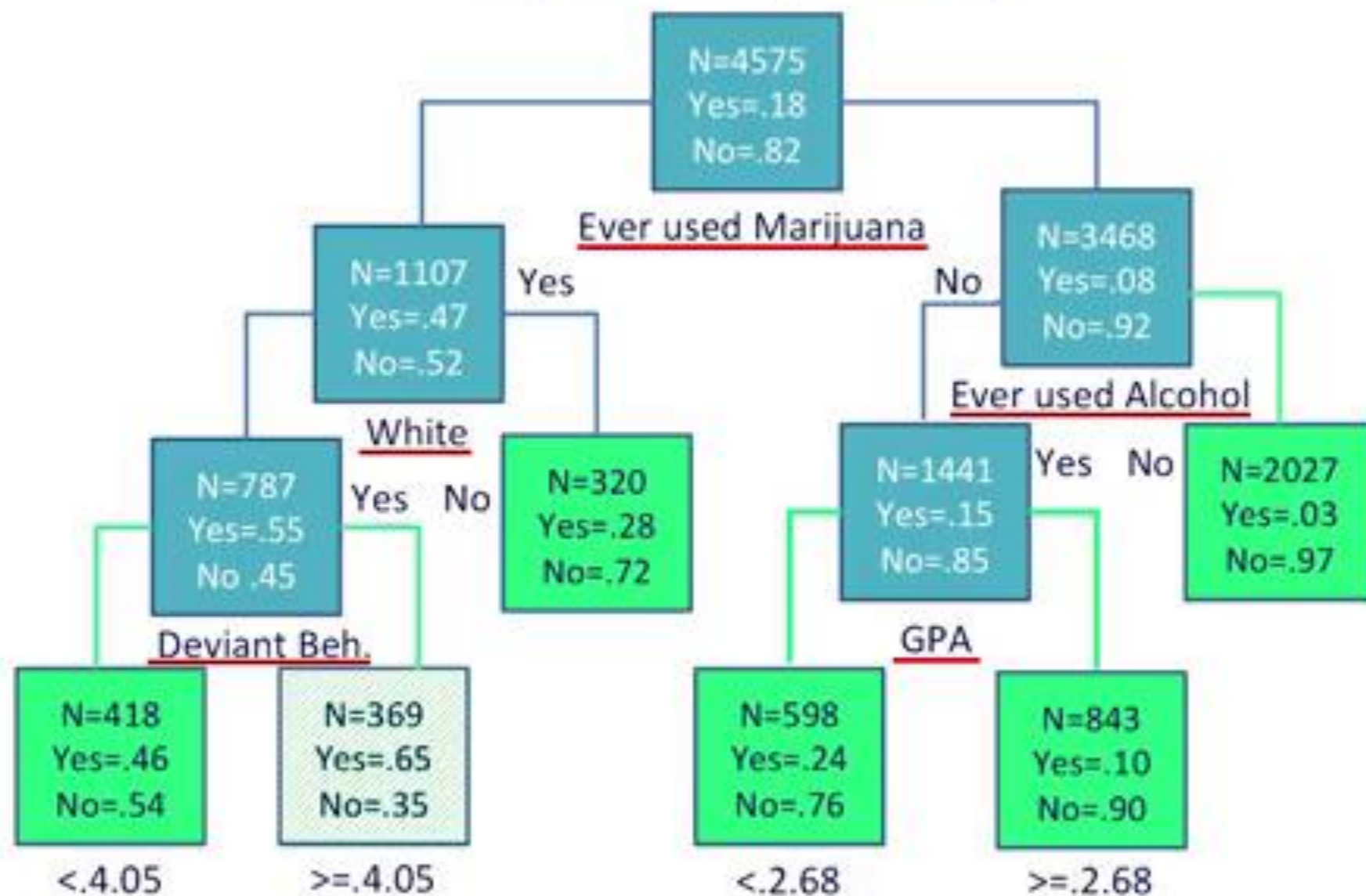
A random subset is tested and only "branches" that improve the classification are retained

Selected sub-tree is the lowest probability of misclassification

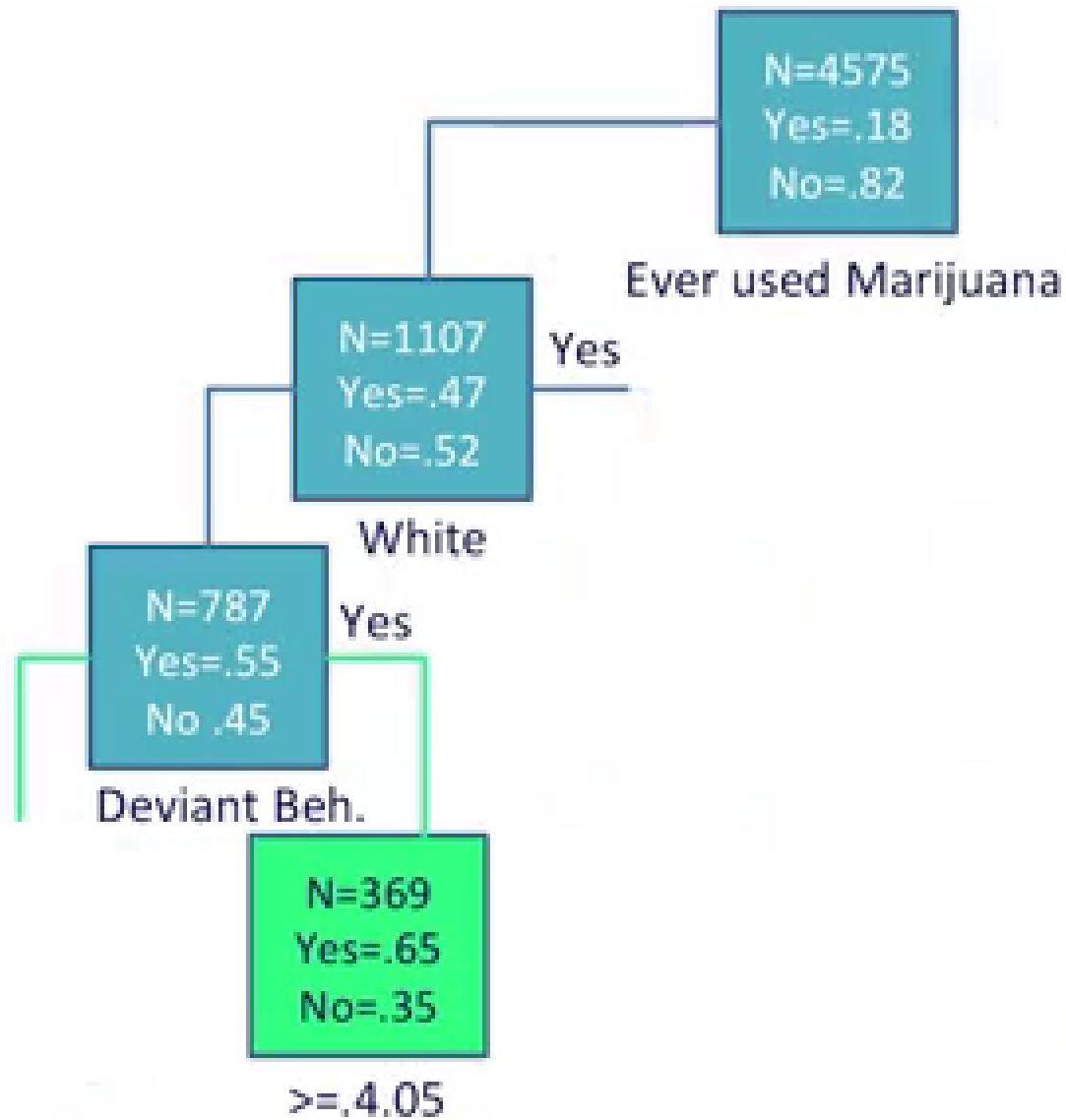
Target Variable: Regular Smoking



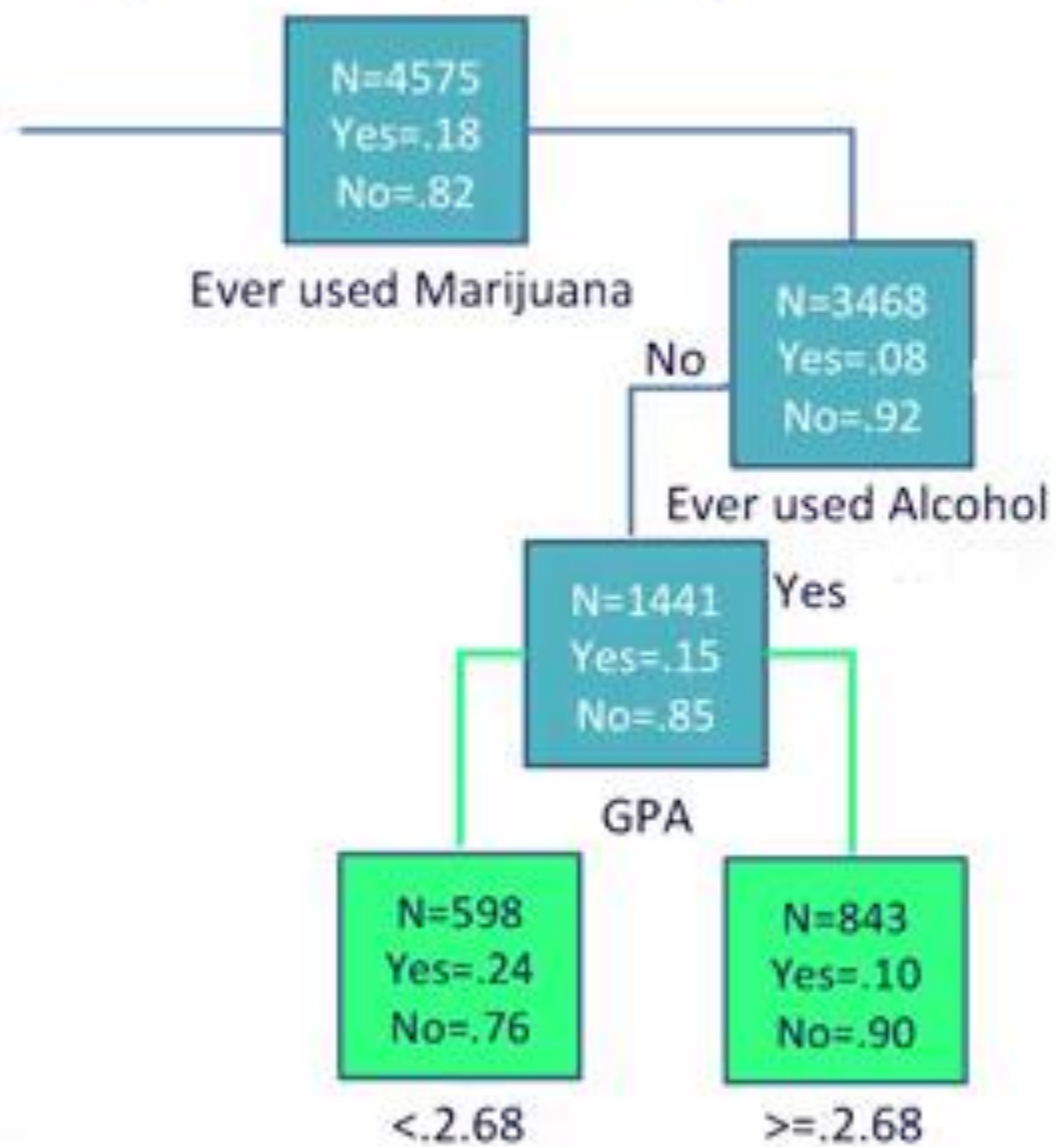
Target Variable: Regular Smoking



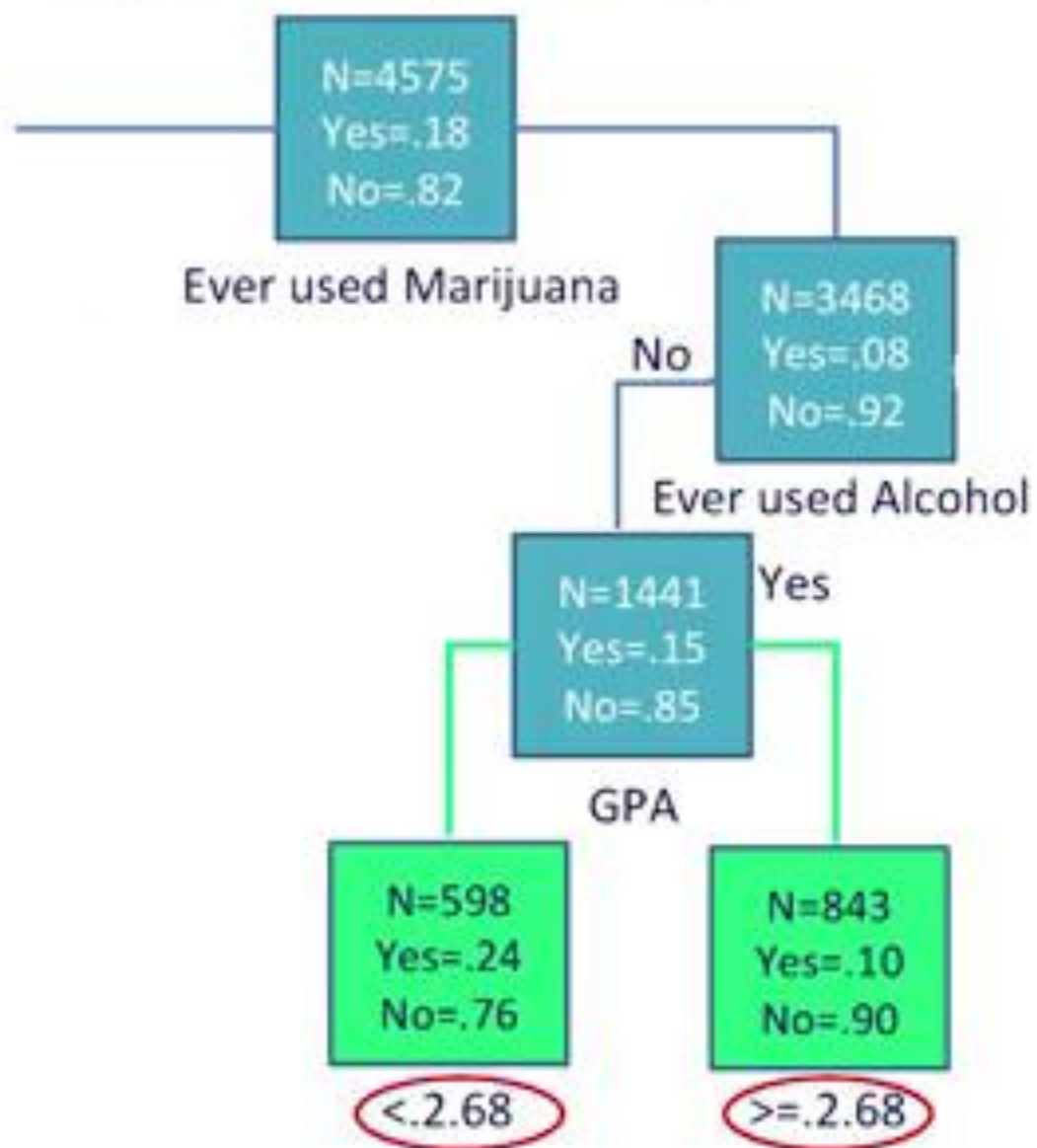
Target Variable: Regular Smoking



Target Variable: Regular Smoking

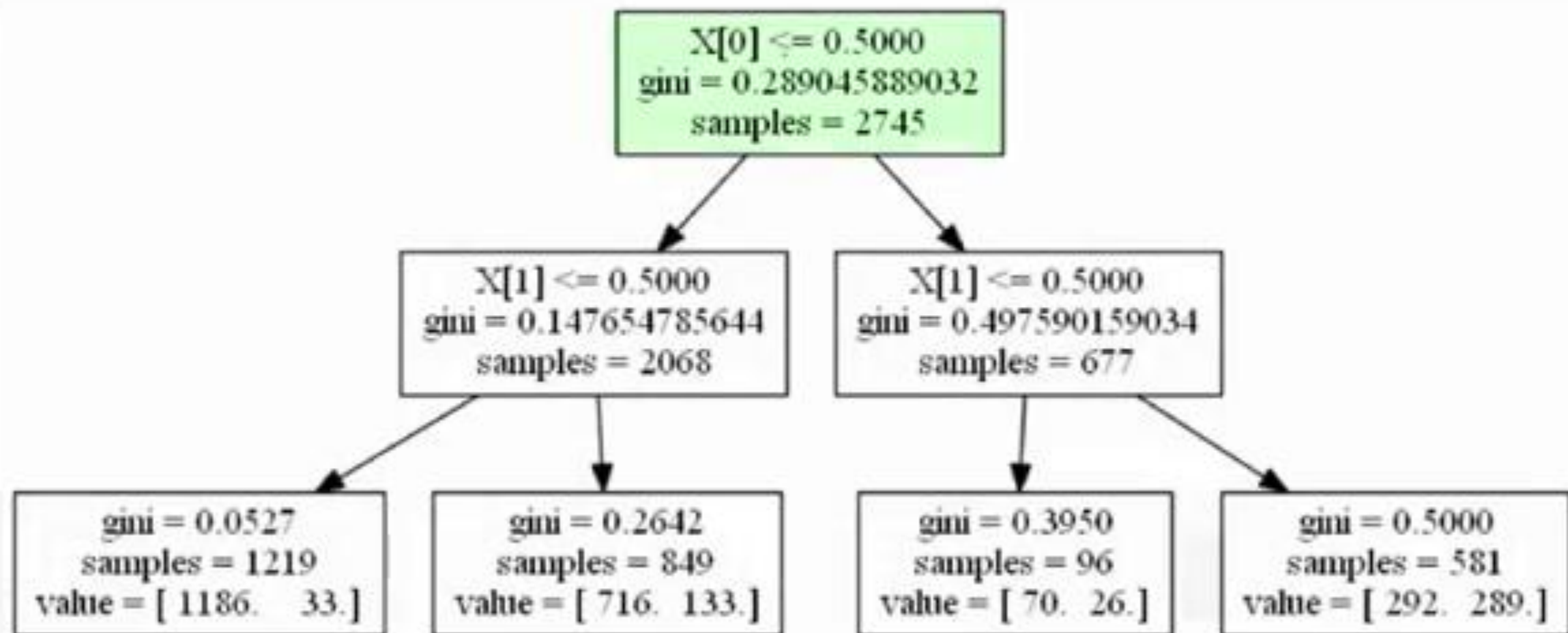


Target Variable: Regular Smoking

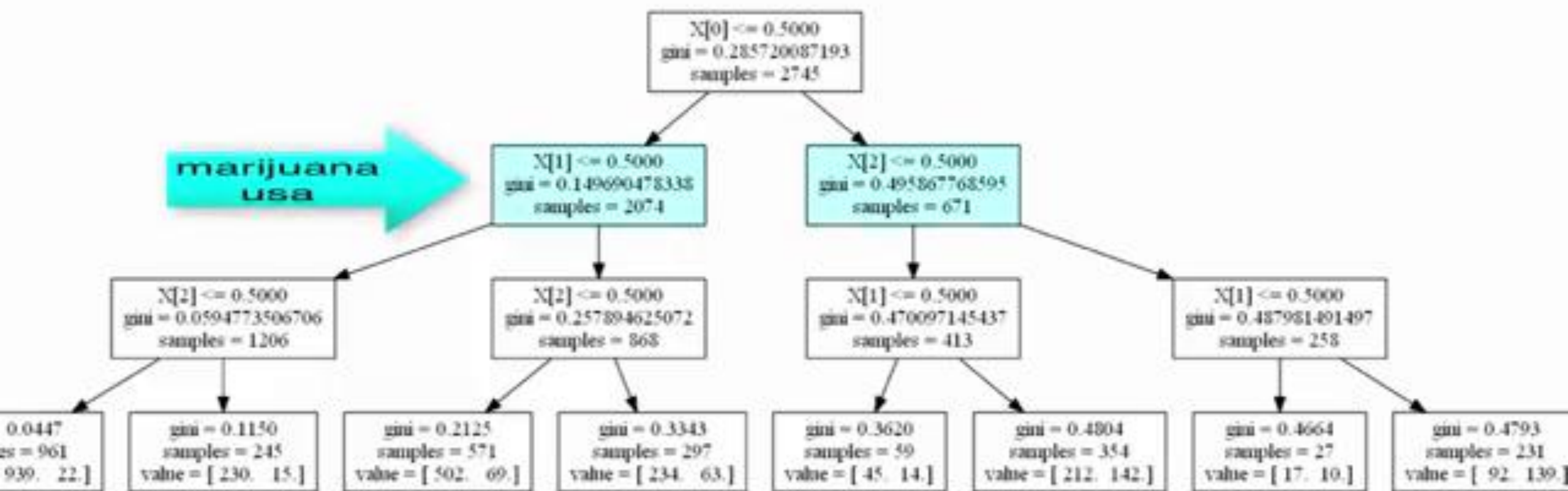


Module 1
Lesson 5 - Building a Decision Tree in Python





marijuana
usa



Strengths of Decision Trees

- Can select from among a large number of variables those and their interactions that are most important in determining the target or response variable to be explained.
- They are easy to interpret and visualize, especially when the tree is small.
- Can handle large data sets well and can predict both binary, categorical target variables (shown in our example) and also quantitative target variables (known as regression trees).

Limitations: Small changes in the data can lead to different splits and this can undermine the interpretability of the model.

Also decision trees are not very reproducible on future data!