

“Intelligence is 10 million rules”
(Doug Lenat)...but Rules are meant to be
generalizable

Reading

What are decision trees?

Carl Kingsford & Steven L Salzberg

Decision trees have been applied to problems such as assigning protein function and predicting splice sites. How do these classifiers work, what types of problems can they solve and what are their advantages over alternatives?

NATURE BIOTECHNOLOGY VOLUME 26 NUMBER 9 SEPTEMBER 2008

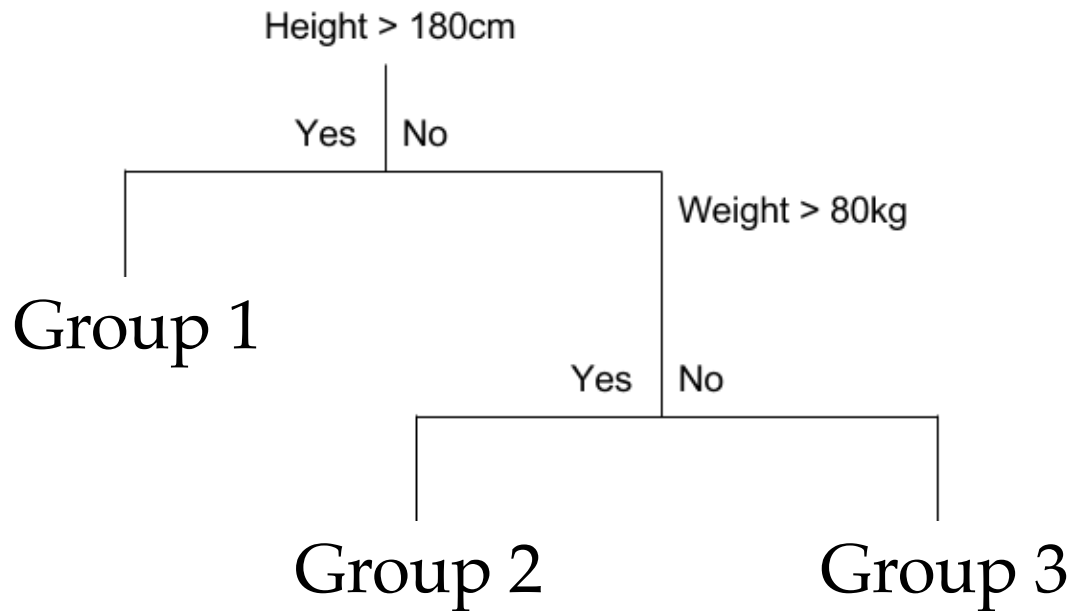
Learning from the data



Pattern recognition

Learning from the data

Machine learning algorithms - Two main types



e.g., Finding number of groups in data and ways to classify (predict) observations based on their characteristics (height / weight)

Learning from the data

Machine learning algorithms - Two main types

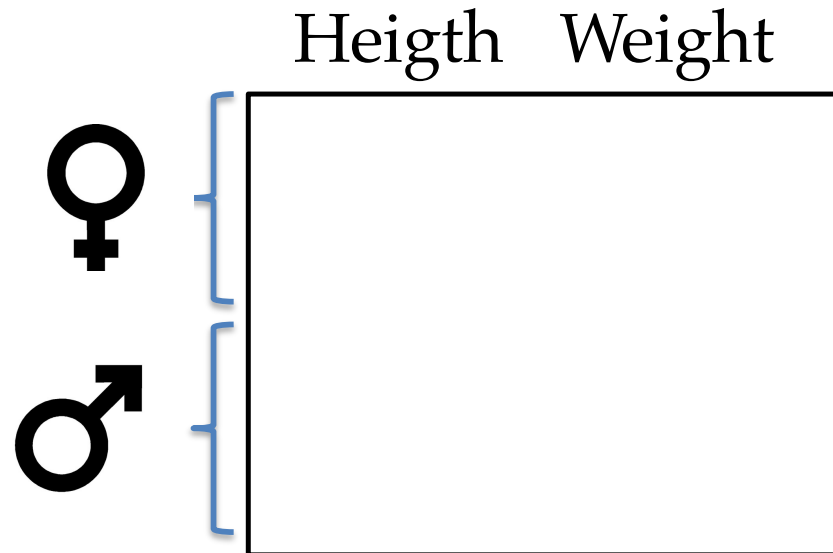


Learning from the data

Machine learning algorithms - Two main types

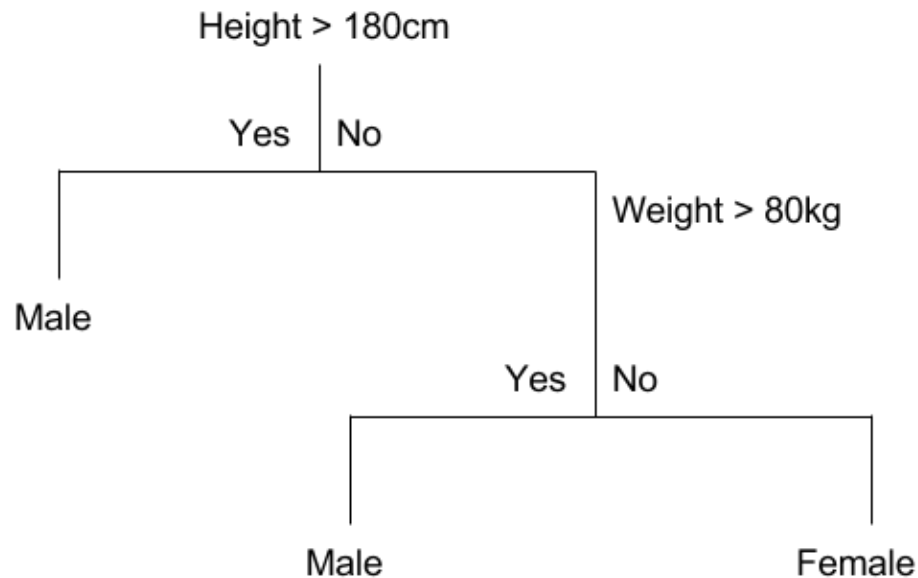


Label =
gender



Learning from the data

Machine learning algorithms - Two main types



Predicting gender on the basis of Height and Weight

Label = gender

CART: Classification and Regression Trees – a powerful (machine learning) yet simple analytical tool for multivariate pattern description

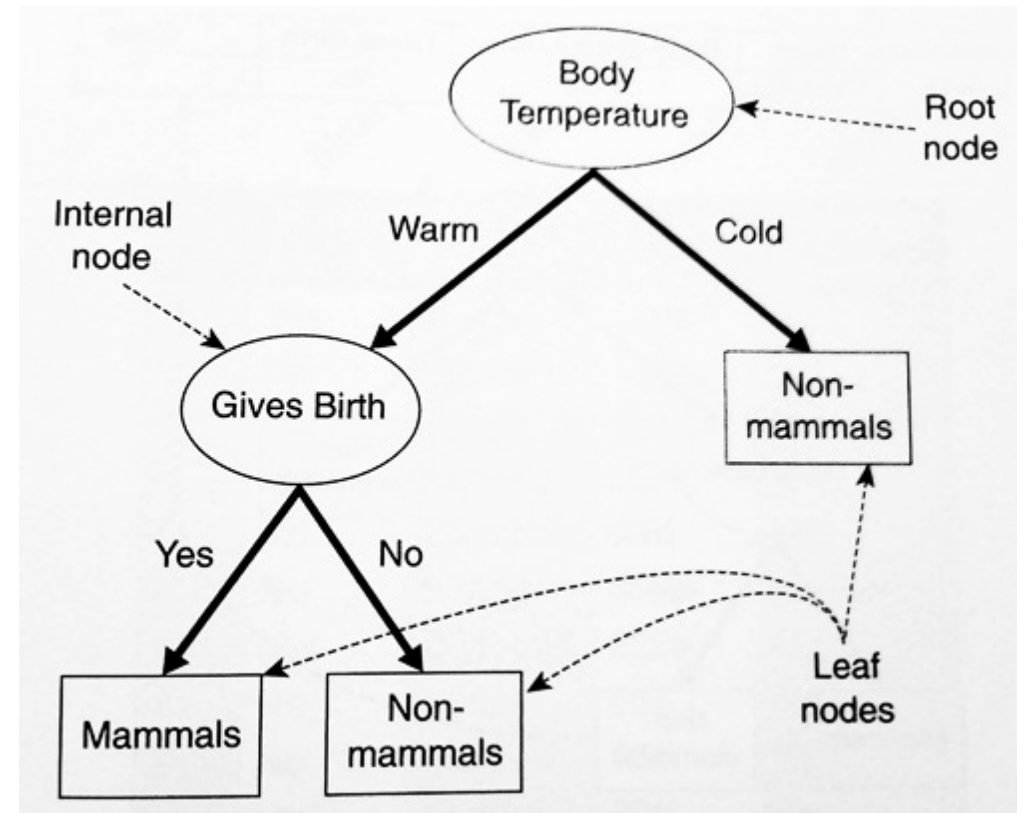
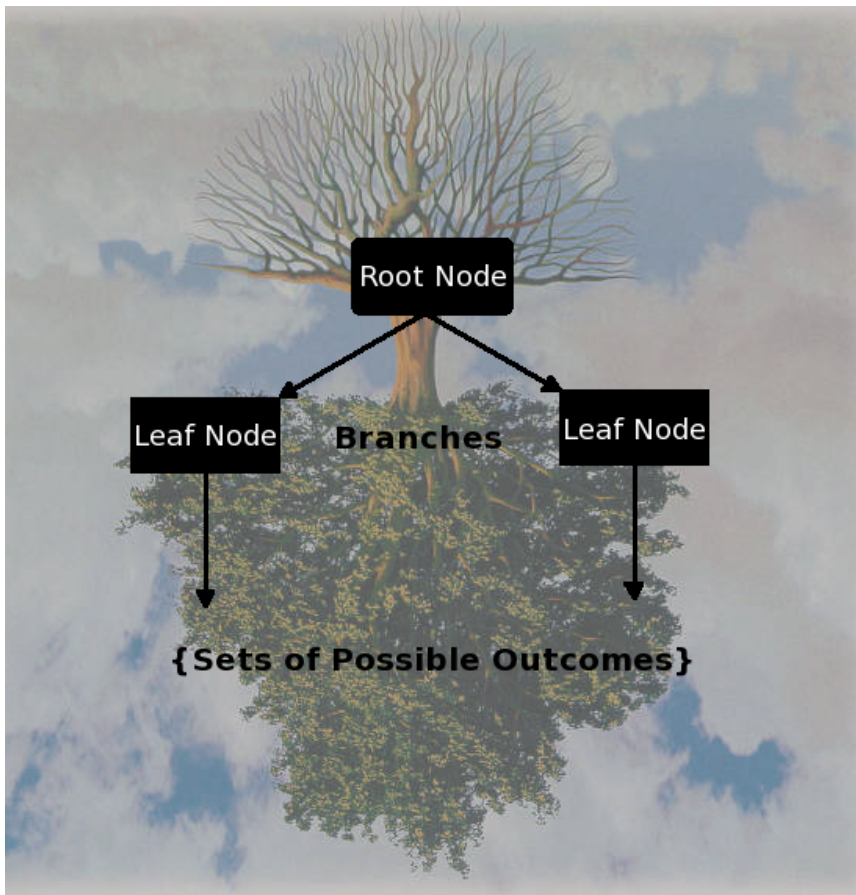


(Leo Breiman and colleagues 1984)

“Decision tree learning is among the most popular machine learning techniques used for ecological modelling. Decision trees can be used to predict the value of one or several (dependent) variables. “ Jopp et al. (2011)

Tree anatomy

“Decision trees are hierarchical structures, where each internal node contains a test on an attribute, each branch corresponding to an outcome of the test, and each leaf node giving a prediction for the value of the class variable.” (Jopp et al. 2011)



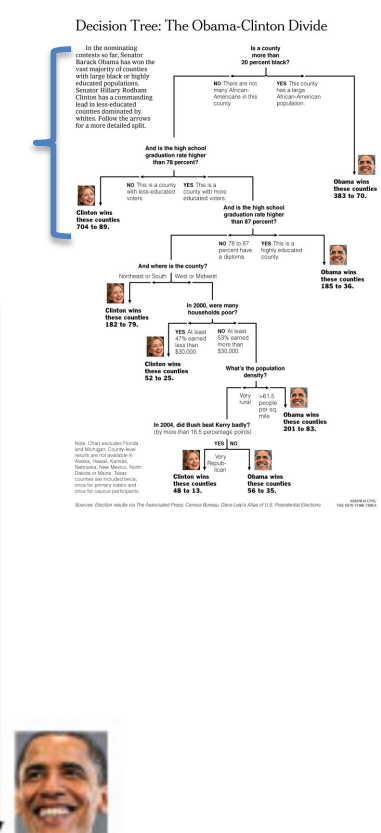
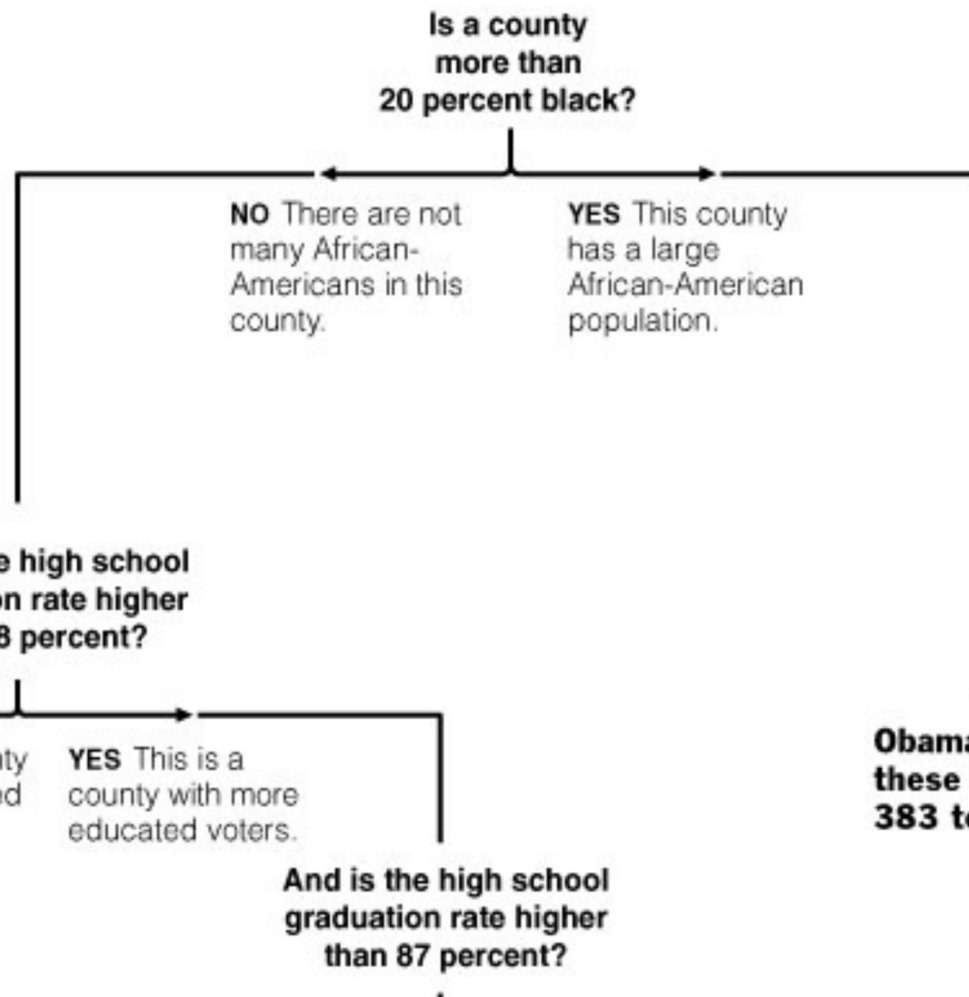
Source
http://mines.humanoriented.com/classes/2010/fall/csci568/portfolio_exports/sdaug_herty/decisiontree.html

Learning from the data – Classification Trees

Deal with complex data but easy to convey results

Decision Tree: The Obama-Clinton Divide

In the nominating contests so far, Senator Barack Obama has won the vast majority of counties with large black or highly educated populations. Senator Hillary Rodham Clinton has a commanding lead in less-educated counties dominated by whites. Follow the arrows for a more detailed split.



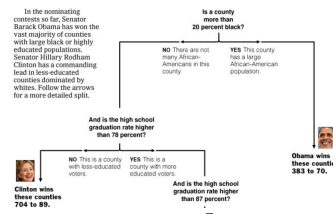
Obama wins these counties 383 to 70.



Clinton wins these counties 704 to 89.

Learning from the data – Classification Trees

Decision Tree: The Obama-Clinton Divide



And is the high school graduation rate higher than 87 percent?

NO 78 to 87 percent have a diploma.

YES This is a highly educated county.



Obama wins these counties 185 to 36.

And where is the county?

Northeast or South

West or Midwest



Clinton wins these counties 182 to 79.

In 2000, were many households poor?

YES At least 47% earned less than \$30,000.

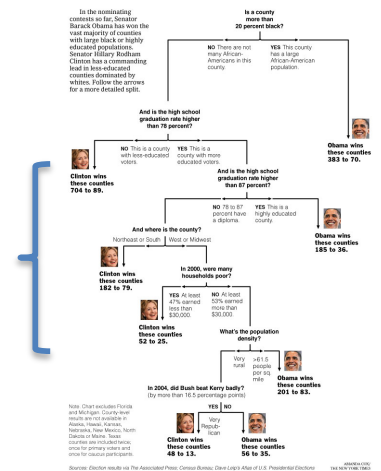
NO At least 53% earned more than \$30,000.



Clinton wins these counties 52 to 25.

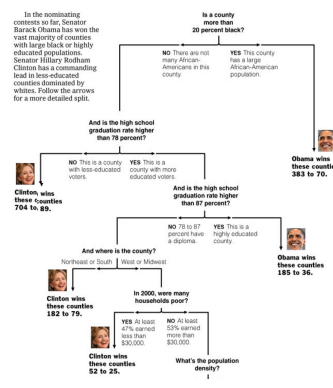
What's the population density?

Decision Tree: The Obama-Clinton Divide

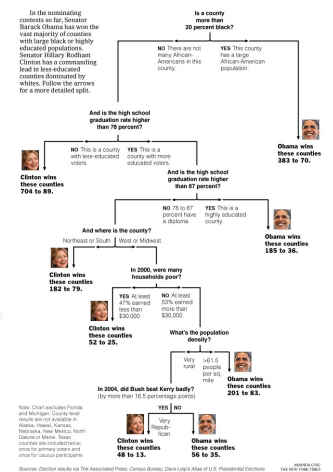


Learning from the data – Classification Trees

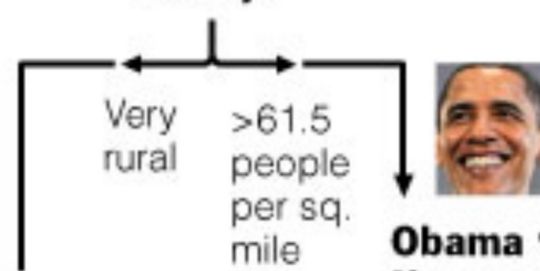
Decision Tree: The Obama-Clinton Divide



Decision Tree: The Obama-Clinton Divide

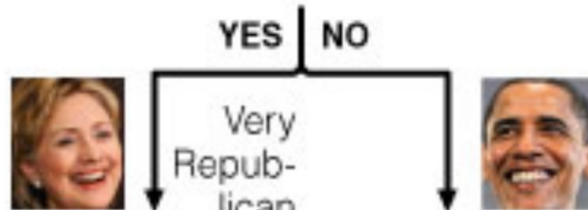


What's the population density?



Obama wins these counties 201 to 83.

In 2004, did Bush beat Kerry badly? (by more than 16.5 percentage points)



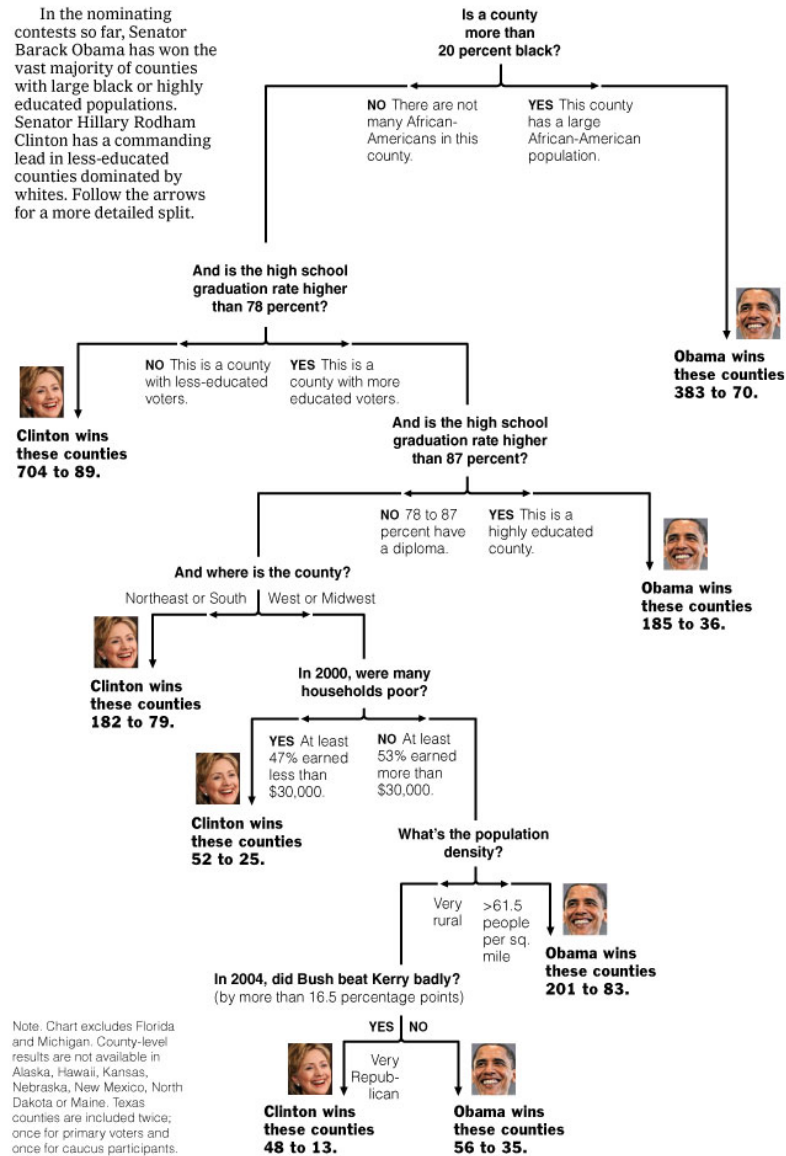
Clinton wins these counties 48 to 13.

Obama wins these counties 56 to 35.

Learning from the data – Classification Trees

Decision Tree: The Obama-Clinton Divide

In the nominating contests so far, Senator Barack Obama has won the vast majority of counties with large black or highly educated populations. Senator Hillary Rodham Clinton has a commanding lead in less-educated counties dominated by whites. Follow the arrows for a more detailed split.

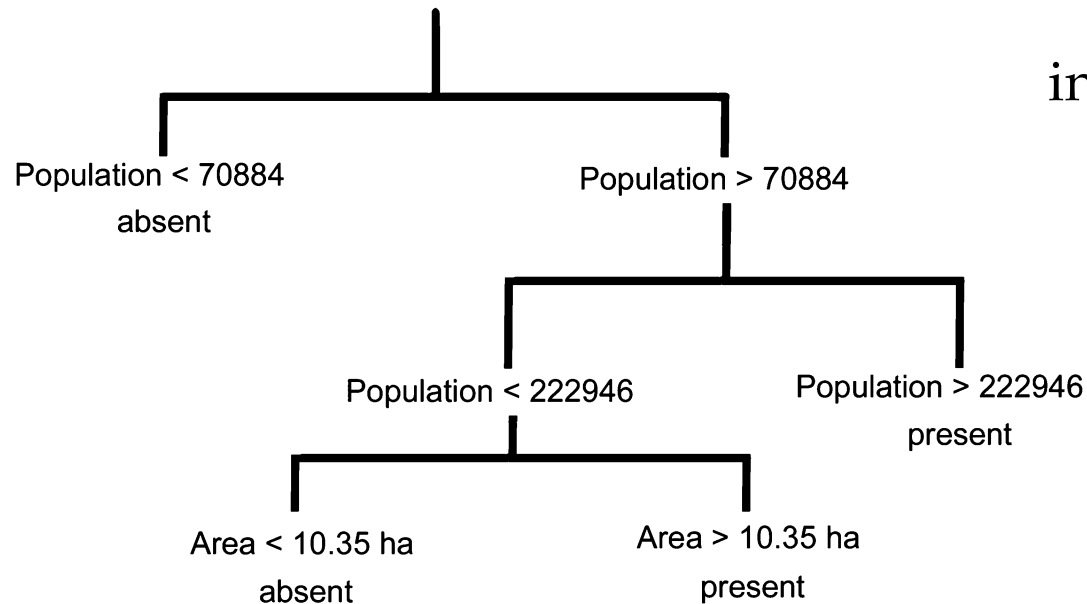


Note. Chart excludes Florida and Michigan. County-level results are not available in Alaska, Hawaii, Kansas, Nebraska, New Mexico, North Dakota or Maine. Texas counties are included twice; once for primary voters and once for caucus participants.

Sources: Election results via The Associated Press; Census Bureau; Dave Leip's Atlas of U.S. Presidential Elections

Non-native species

Smallmouth bass



Classification trees model categories,
including 1s and 0s (male / female,
presence / absence, non-
infected / recovery / infected)

Probabilities of presence are
calculated by the classification
tree, but then transformed into
1 (e.g., >0.50) or absence
(e.g., ≤ 0.50)

Figure 1 Summary of classification tree analysis predicting smallmouth bass occurrence in British Columbia based on lake morphology, distance to road and human population census data.

Columbia based on the classification tree analysis (Fig. 1). Evaluation of the independent validation dataset showed that overall classification success was 93.5%, with 83.1% sensitivity and 100% specificity. Extrapolation of the classification tree

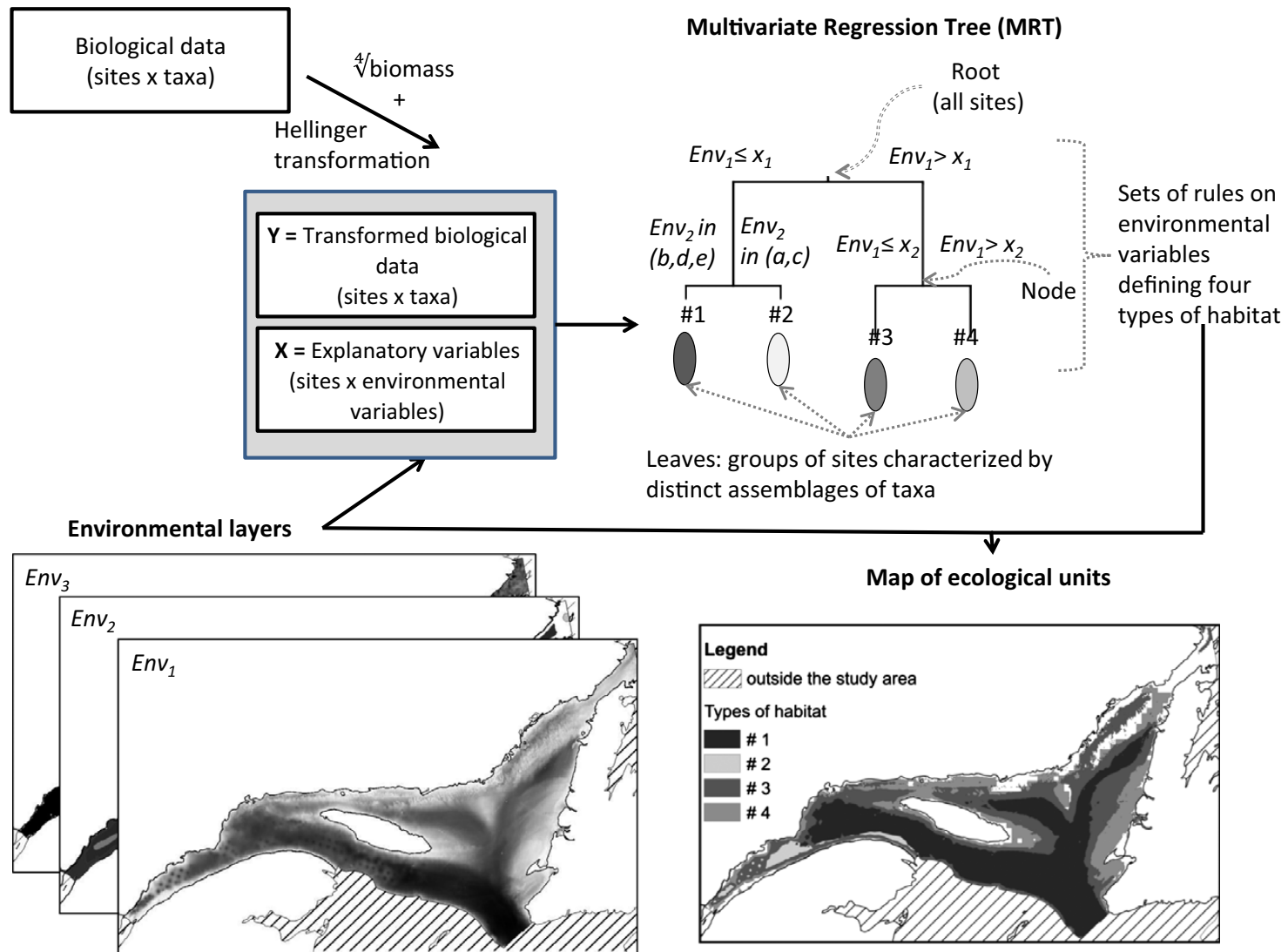
Diversity and Distributions, (Diversity Distrib.) (2009) 15, 831–840



Predicting introduction, establishment and potential impacts of smallmouth bass

Sapna Sharma^{1*}, Leif-Matthias Herborg^{2,3} and Thomas W. Theriault²

Regression trees model quantitative variables (e.g., species abundances)



Classification versus Regression Trees (CART)

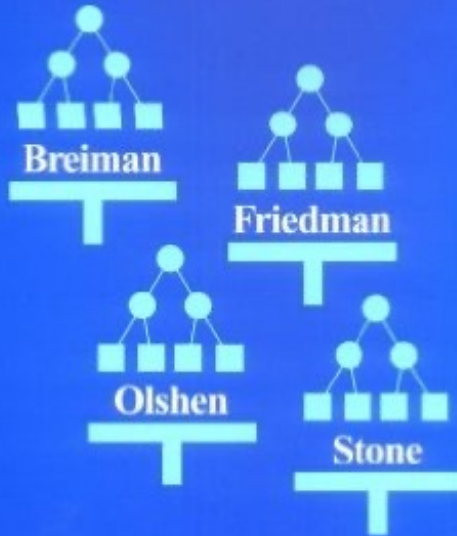
- Classification (sometimes referred as to decision trees) trees model dependent variables that have a finite number of categories (unordered values) - This lecture.
- Regression trees model dependent variables that are continuous.

The classification tree algorithm



Copyrighted Material

CLASSIFICATION AND REGRESSION TREES



Copyrighted Material

Classification and Regression Trees
(Wadsworth statistics / probability series)

Breiman, Leo

Note: This is not the actual book cover

DATA MINING WITH DECISION TREES

Theory and Applications

Lior Rokach ♦ Oded Maimon



SERIES IN
MACHINE PERCEPTION
ARTIFICIAL INTELLIGENCE
Volume 69

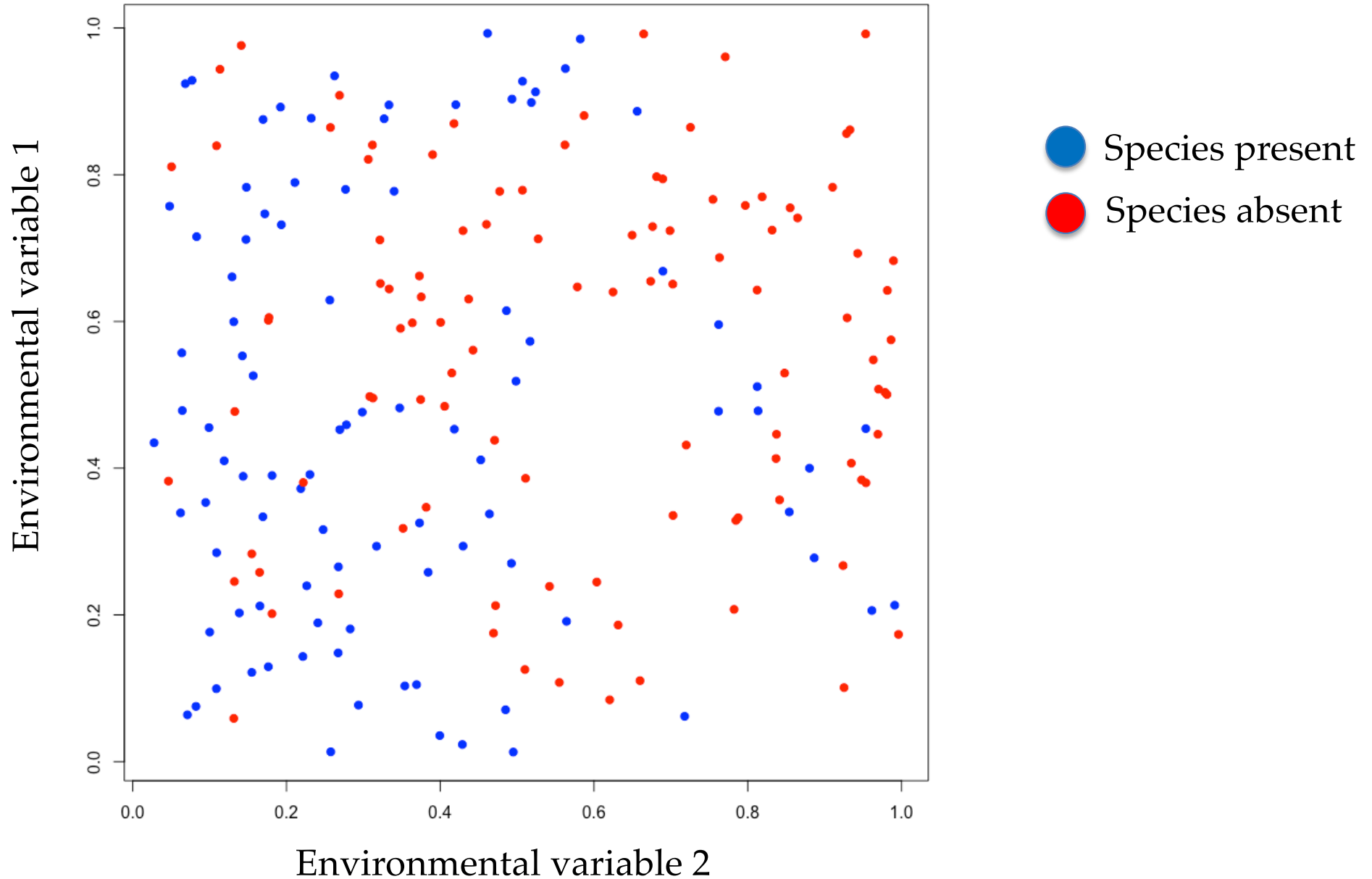


Ecology, 81(11), 2000, pp. 3178–3192
© 2000 by the Ecological Society of America

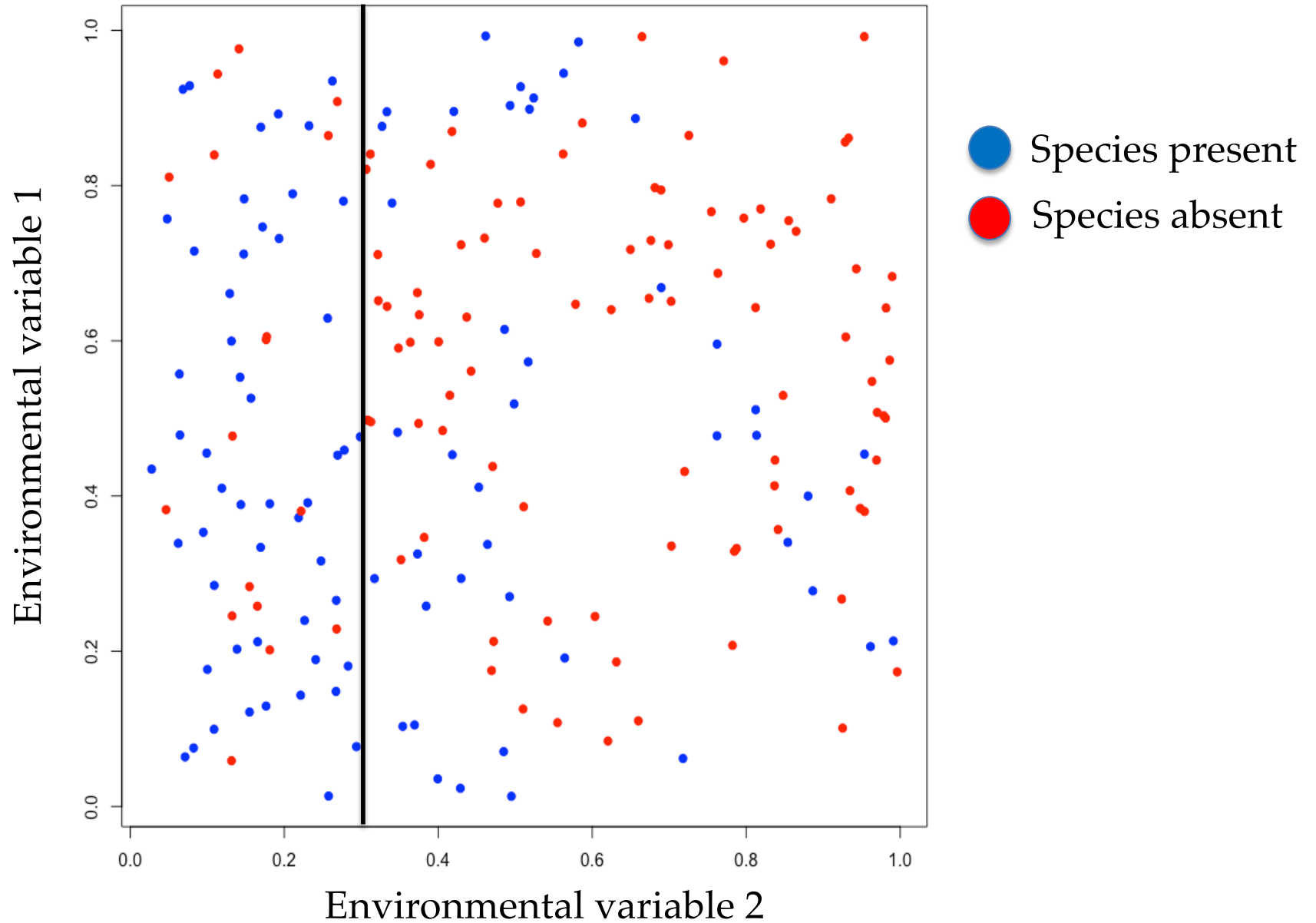
CLASSIFICATION AND REGRESSION TREES: A POWERFUL YET SIMPLE TECHNIQUE FOR ECOLOGICAL DATA ANALYSIS

GLENN DE'ATH¹ AND KATHARINA E. FABRICIUS²

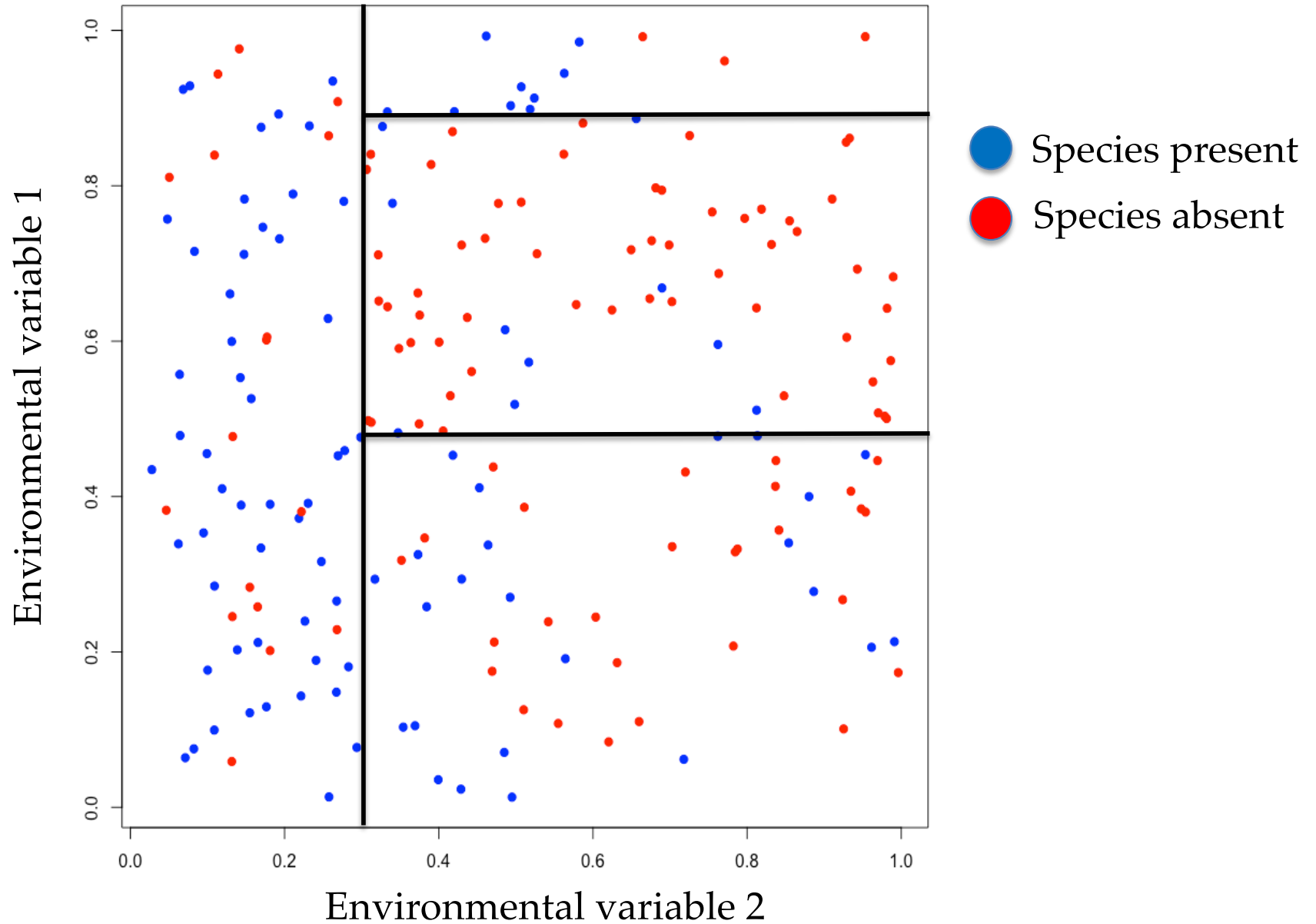
How to model these data?



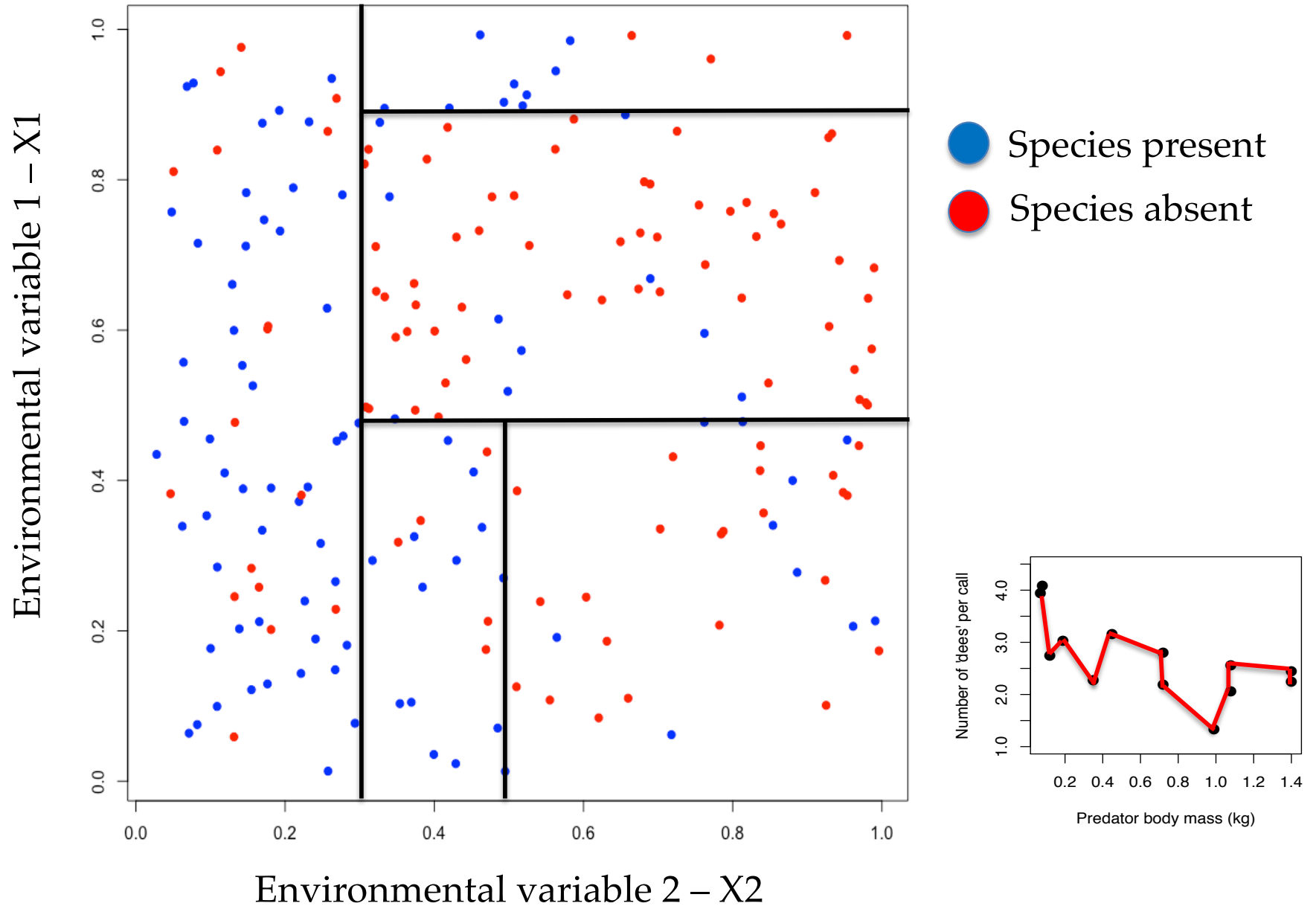
Mosaic or partition plot



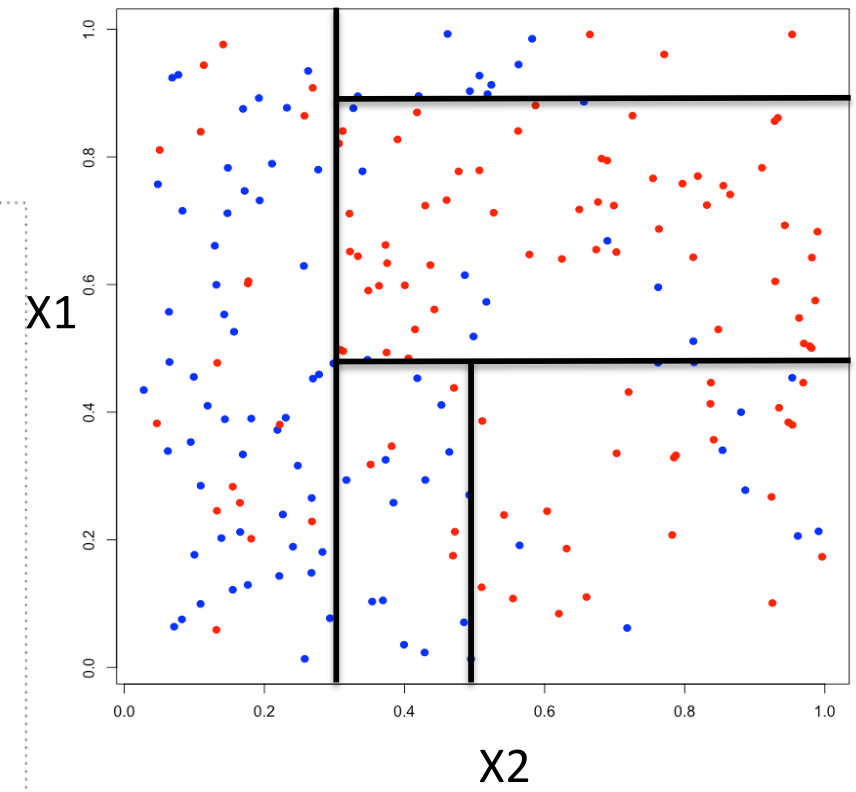
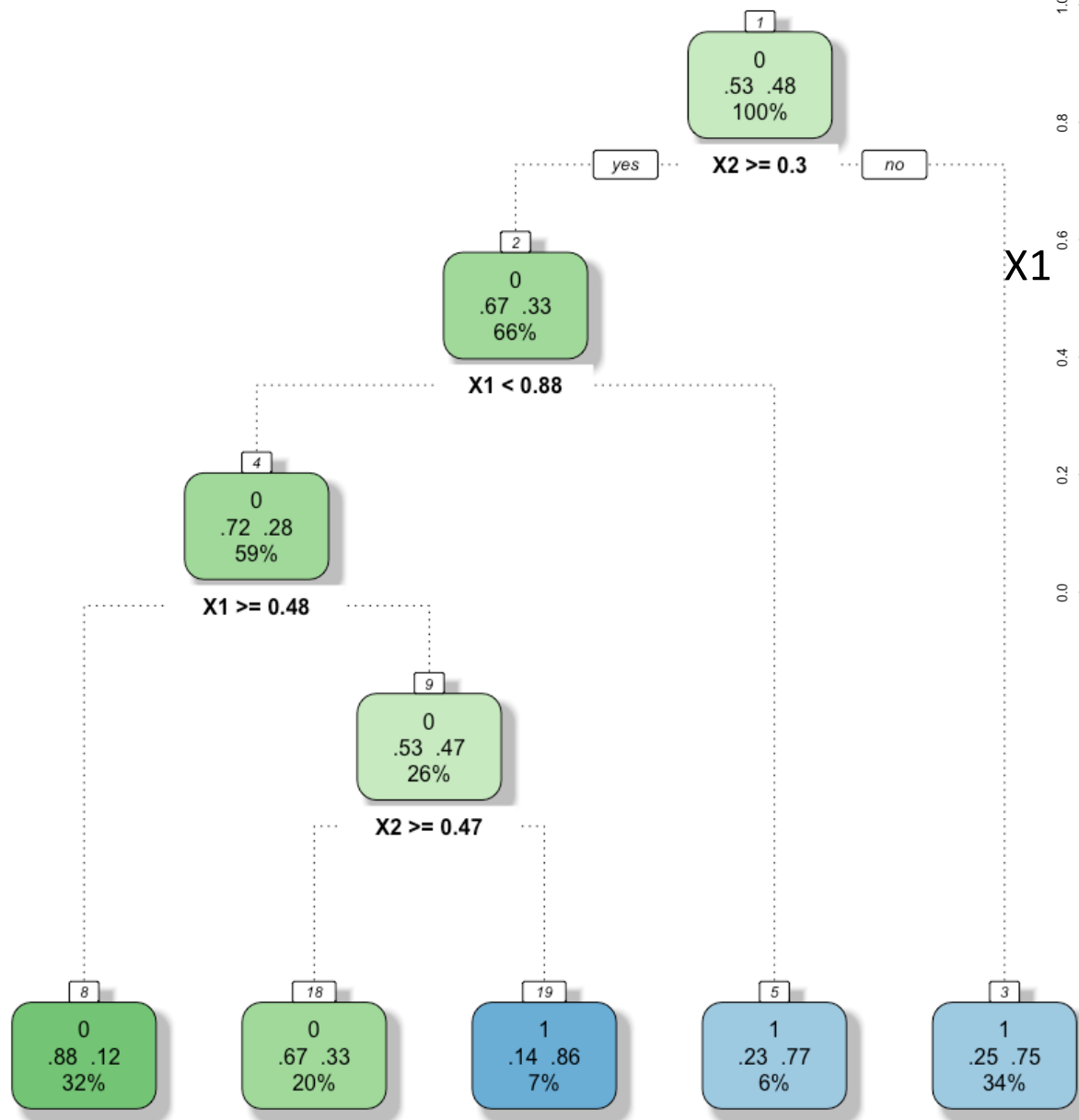
Mosaic or partition plot



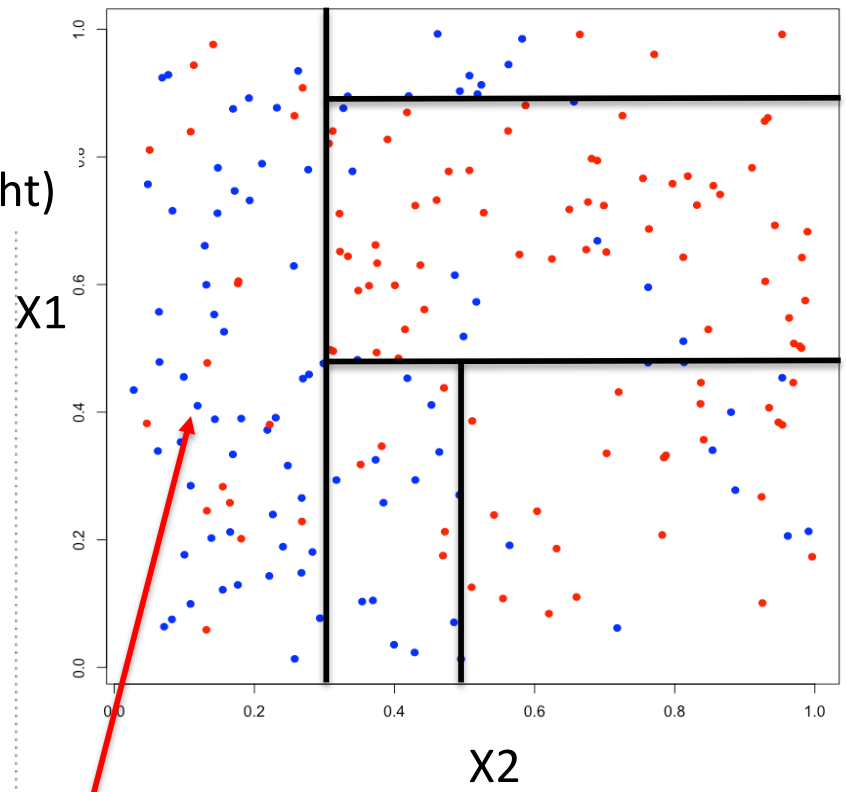
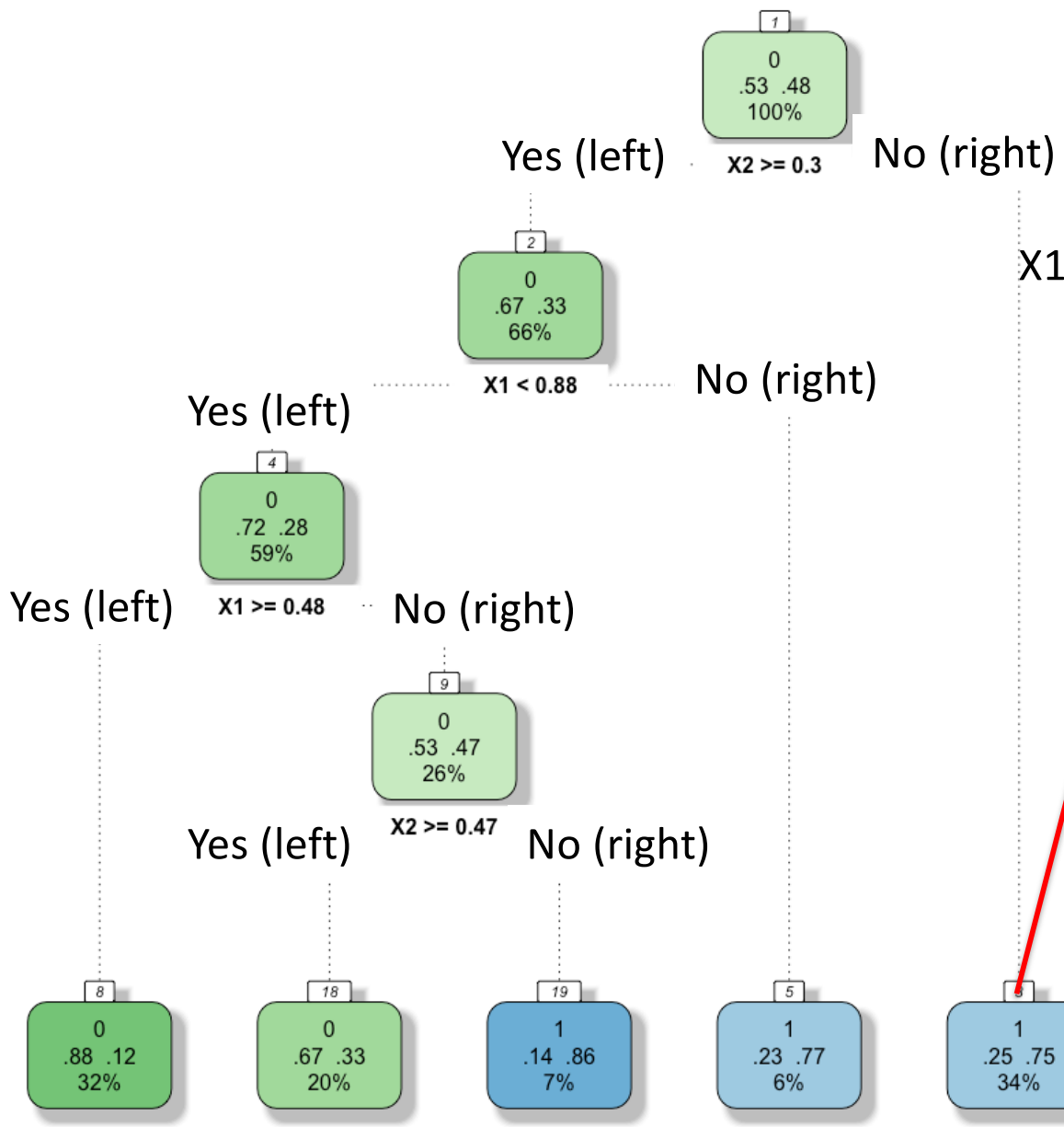
Mosaic or partition plot (best partitioning possible without too much fitting; many ways to determine final model)



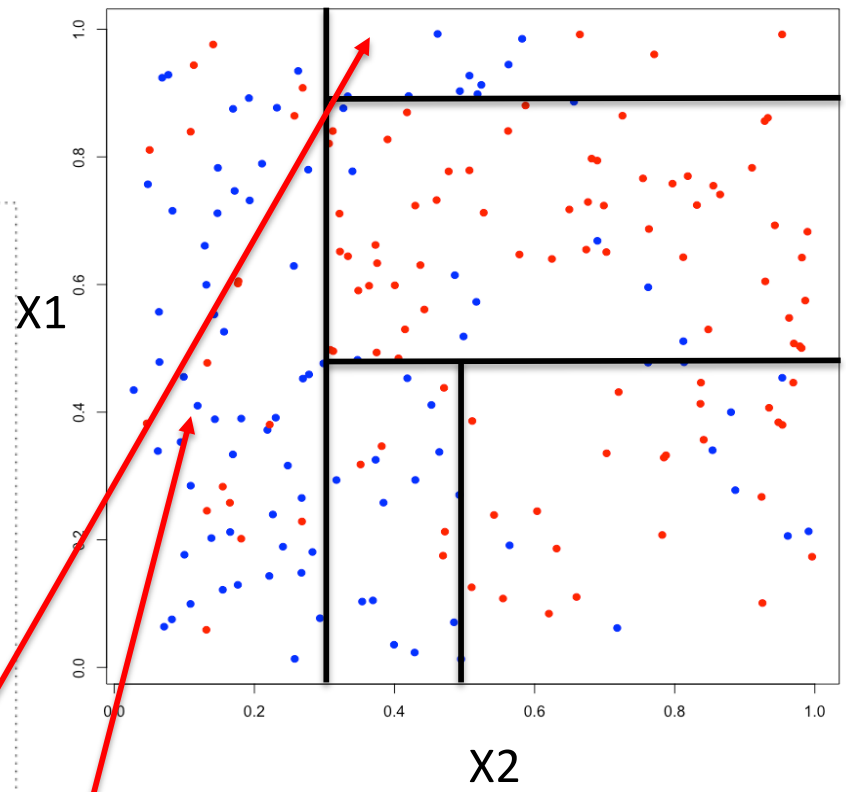
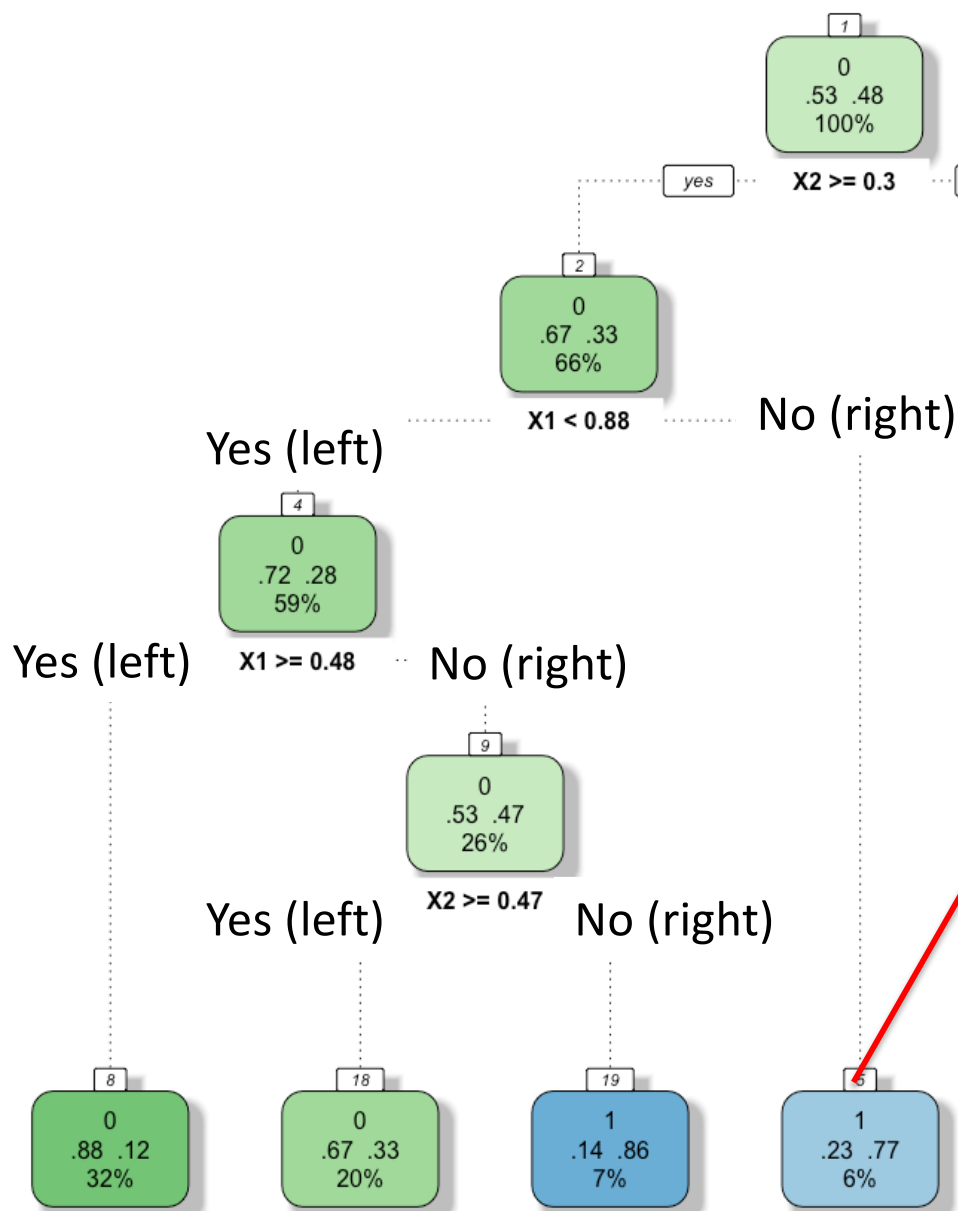
Classification tree for the data



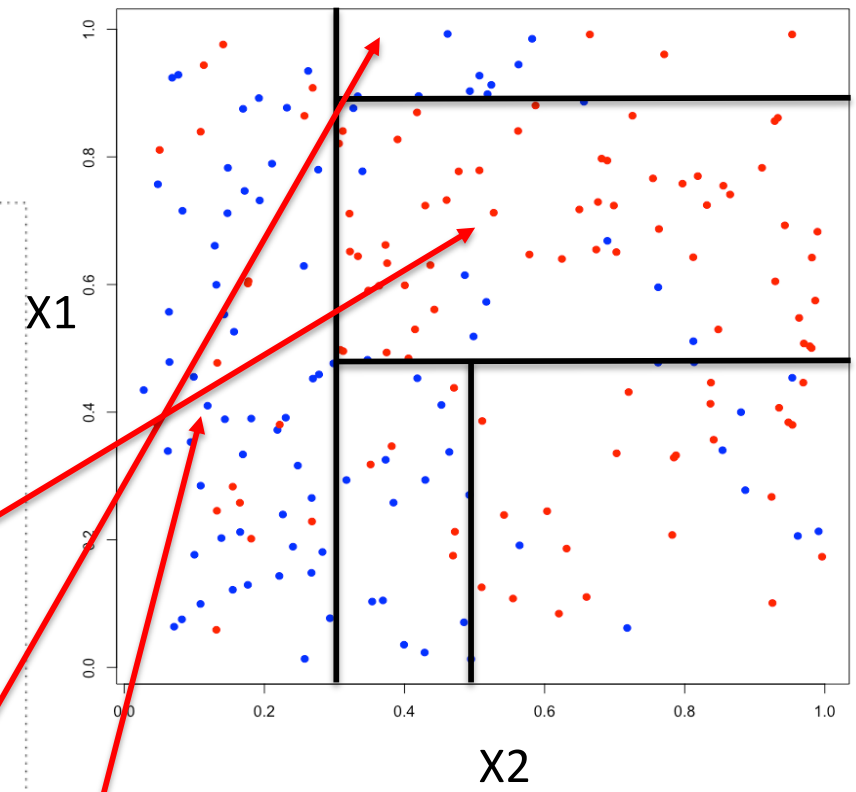
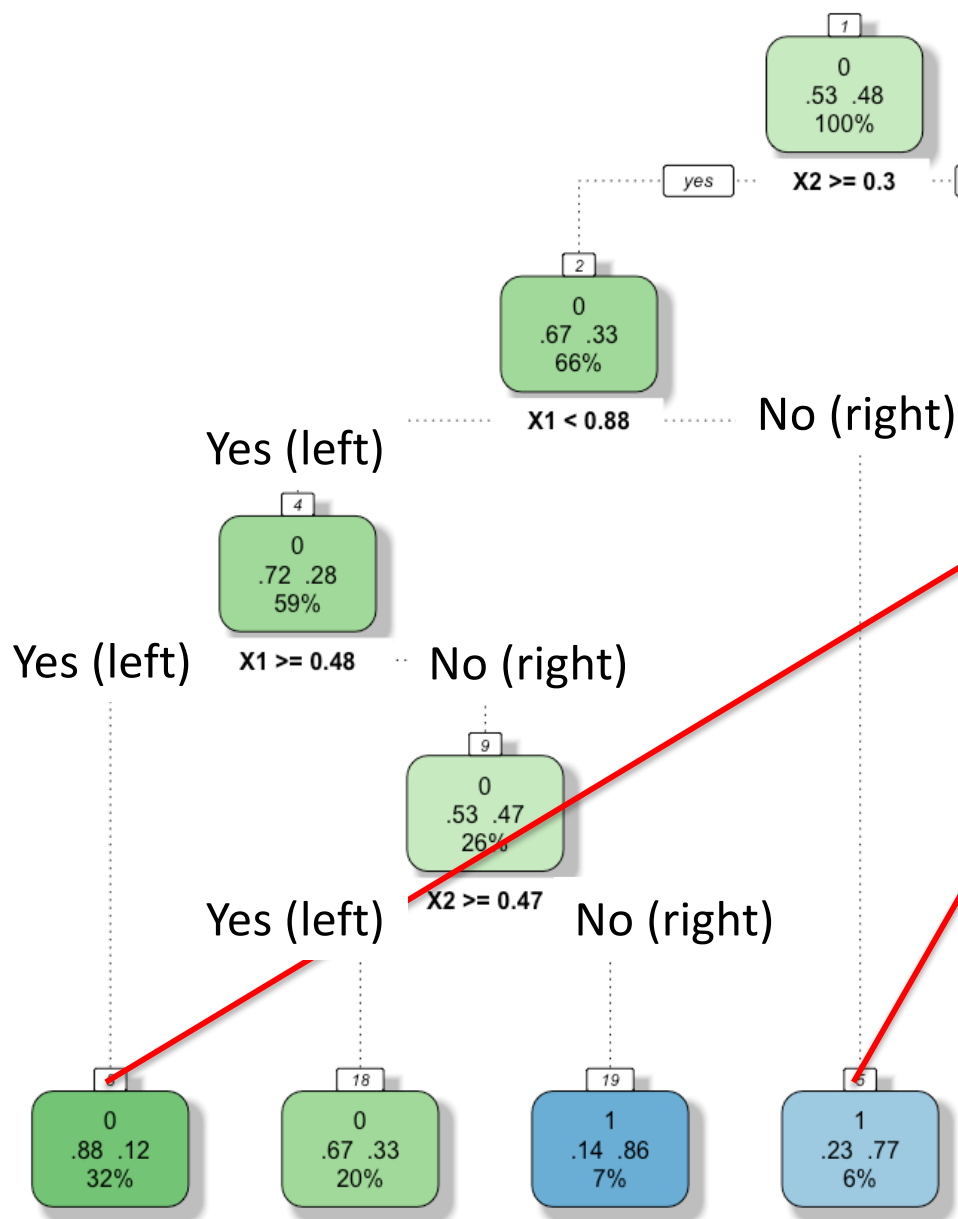
Classification tree for the data



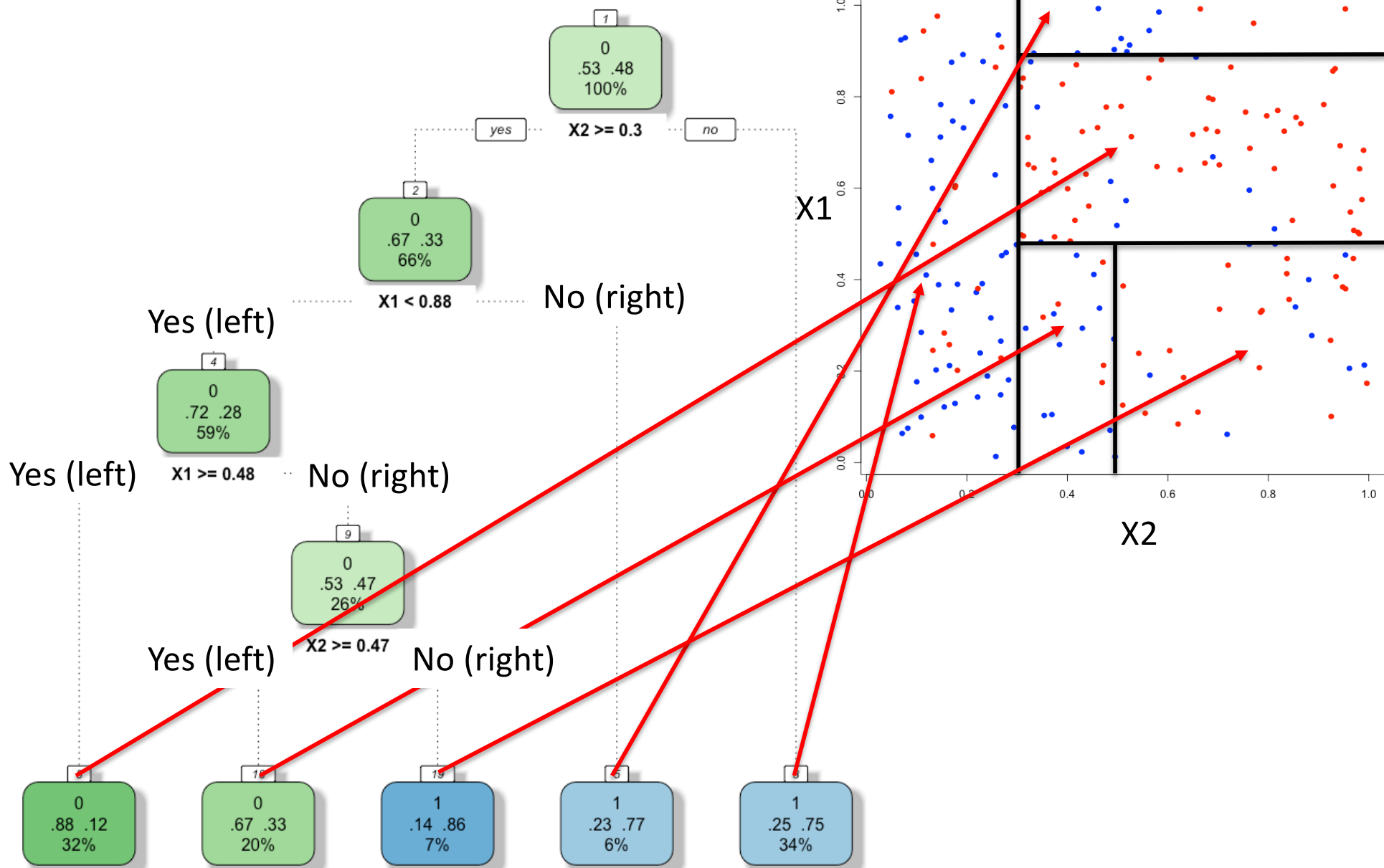
Classification tree for the data



Classification tree for the data



Classification tree for the data



Den Boer et al. 2009

A subtype of childhood acute lymphoblastic leukaemia with poor treatment outcome: a genome-wide classification study



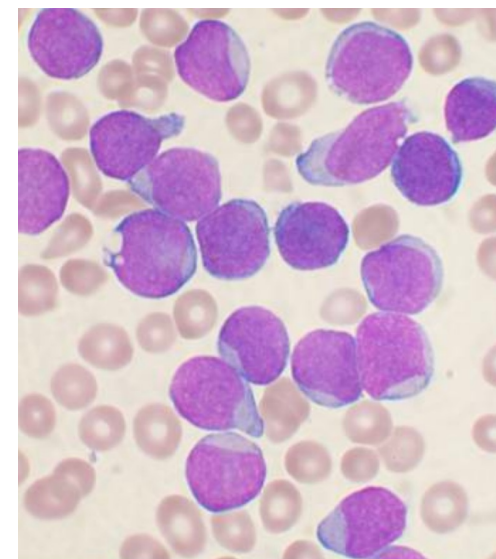
Monique L Den Boer, Marjon van Slegtenhorst*, Renée X De Menezes, Meyling H Cheok, Jessica G C A M Buijs-Gladdines, Susan T C J M Peters, Laura J C M Van Zutven, H Berna Beverloo, Peter J Van der Spek, Gaby Escherich†, Martin A Horstmann†, Gritta E Janka-Schaub†, Willem A Kamps‡, William E Evans, Rob Pieters‡*

Background - In childhood acute lymphoblastic leukemia (ALL) genetic subtypes are recognized that determine the risk-group for further treatment. However, 25% of precursor B-cell ALL (most common type of ALL) are currently genetically unclassified and have an intermediate prognosis. The present study used genome-wide strategies to reveal new biological insights and advance the prognostic classification of childhood ALL.

Lancet Oncol 2009; 10: 125-34

Published Online
January 9, 2009
DOI:10.1016/S1470-

The expression of 22283 genes across 190 patients were considered

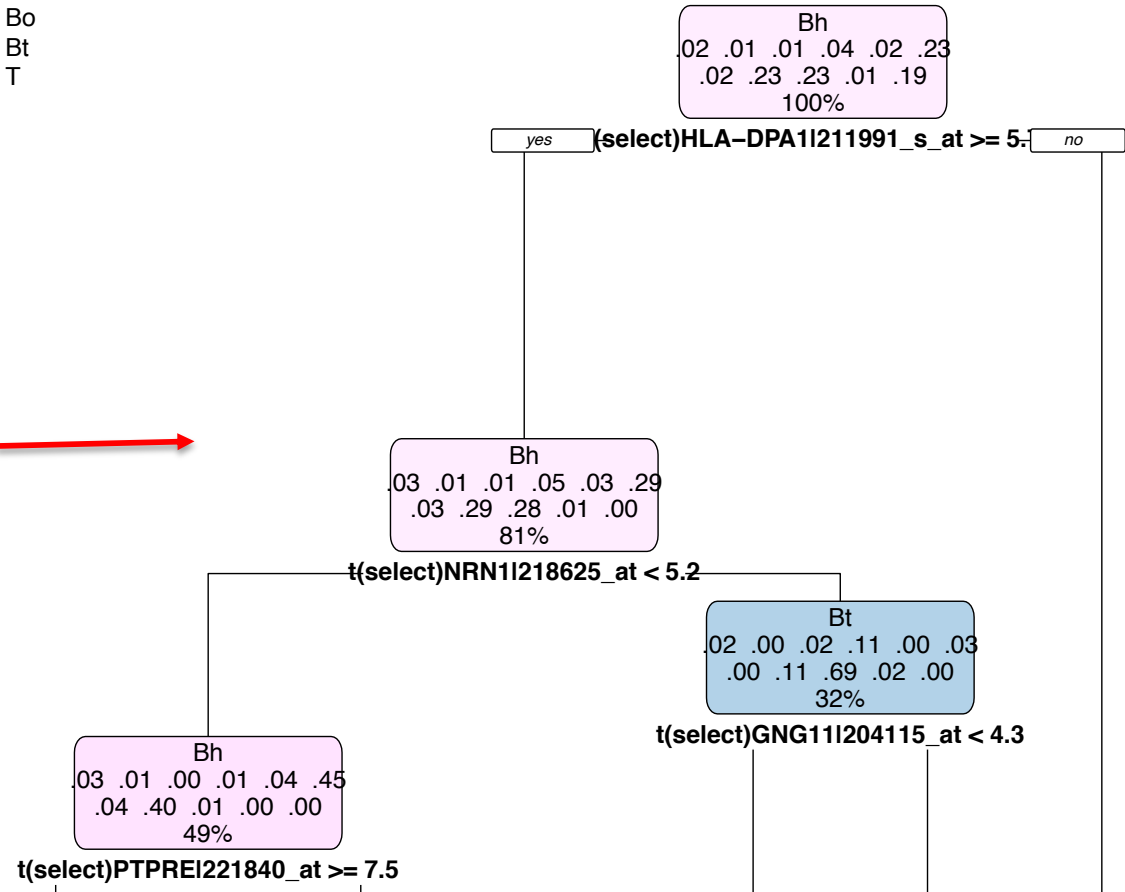
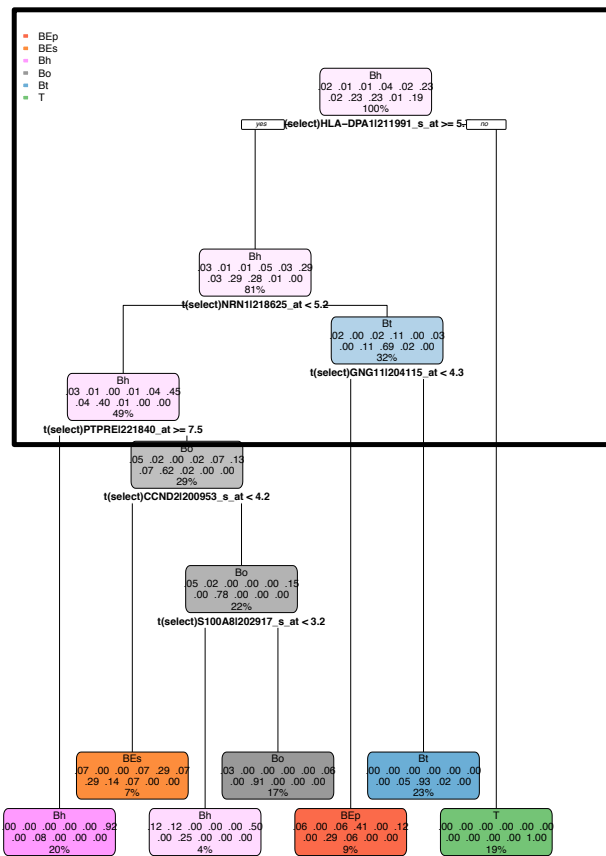


Several subgroups of childhood acute lymphoblastic leukaemia (ALL) have unfavourable prognosis

- BEp
- BEs
- Bh
- Bo
- Bt
- T

Bo (B-other are about 25% of patients and remains unclassified)

Can we improve prognosis based on gene expression?



The expression of 22283 genes across 190 patients were considered to build the model (calibration); 107 independent patients were predicted by the model (validation). The model was 87.7% accurate!

A subtype of childhood acute lymphoblastic leukaemia with poor treatment outcome: a genome-wide classification study

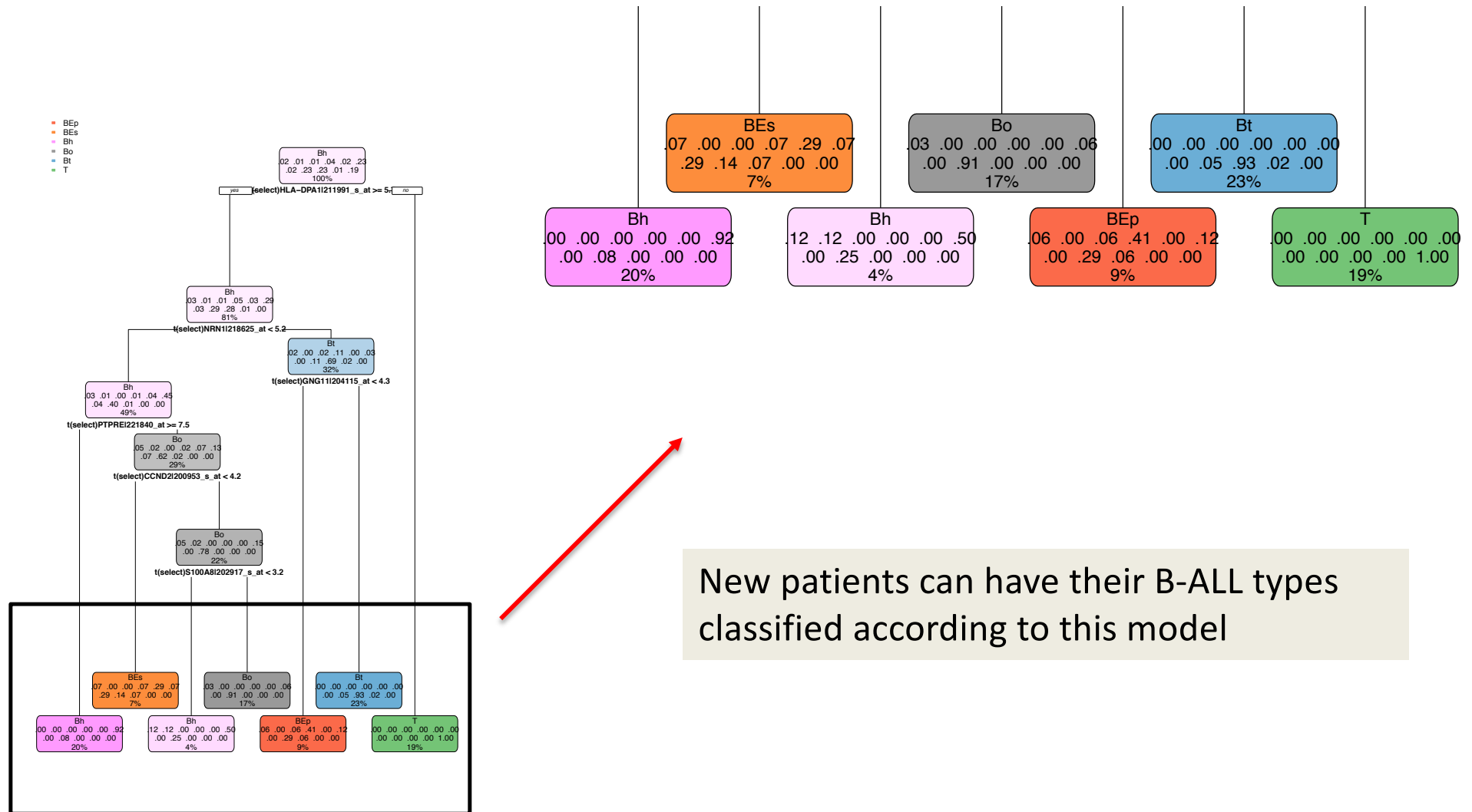


Monique I. Den Boer*, Marjol van Skogemhorst*, Renée X. De Menezes, Mayling H. Cheok, Jessica G.C.A.M. Bujsi-Gladines, Susan T.C.J.M. Peters, Laura J.C.M. Von Zethoven, Willem Bervenloo, Peter J. Van der Spek, Gaby Eschench†, Martin A. Hoestromm†, Gritta E. Junka-Schaub†, Willem A. Kamps, William E. Evans, Rob Pieters†

Summary
 Background Genetic subtypes of acute lymphoblastic leukaemia (ALL) are used to determine risk and treatment in children. 25% of precursor B-ALL cases are genetically unclassified and have intermediate prognosis. We aimed to use a genome-wide study to improve prognostic classification of ALL in children.

Lancet Oncol 2009; 10: 125-34
 Published Online
 January 9, 2009
 DOI:10.1016/S1473-0166(08)70429-0

- BEp
- BEs
- Bh
- Bo
- Bt
- T



New patients can have their B-ALL types classified according to this model

Growing a tree

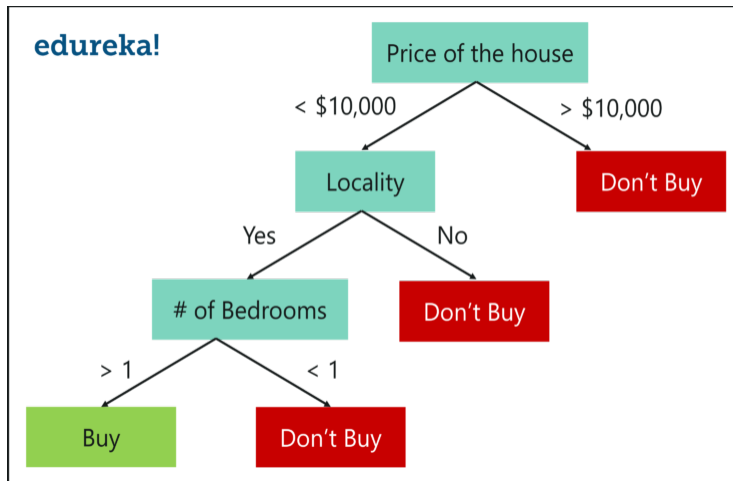
- There are many ways of building CARTs and many complex and advanced ways of doing it.
- Search and establishing hierarchy among variables
Partition values of a variable: $X \leq c$ and $X > c$ for "all" possible c values. Compare fit using (for example) *pseudo* R^2 (correlation between predicted and observed).
- Order of variables are important and may influence the tree – bagging & random forests deal with this issue via building multiple trees (bootstrap) and selecting trees that maximize R^2 or average trees.

More complex models for building trees

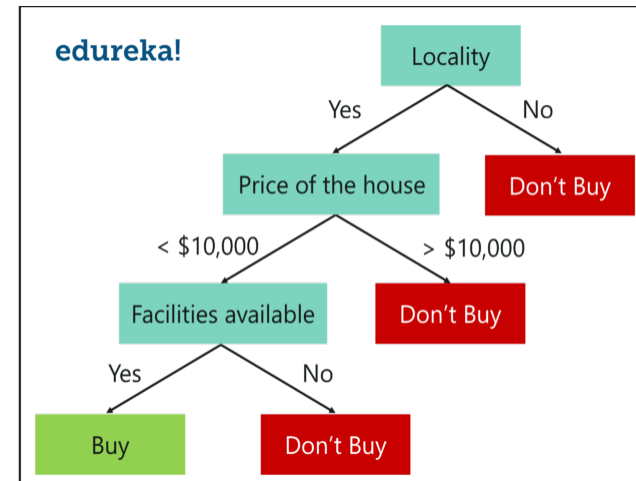
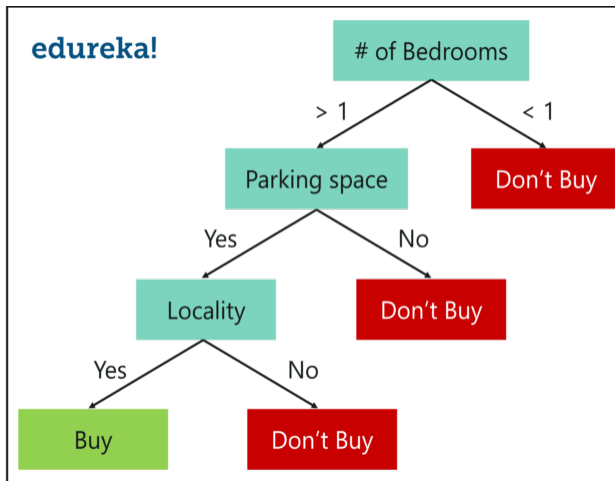
bagging: bootstrapping objects but keeping all predictors

Model for determining factors that influenced house purchasing

Data subsample 1



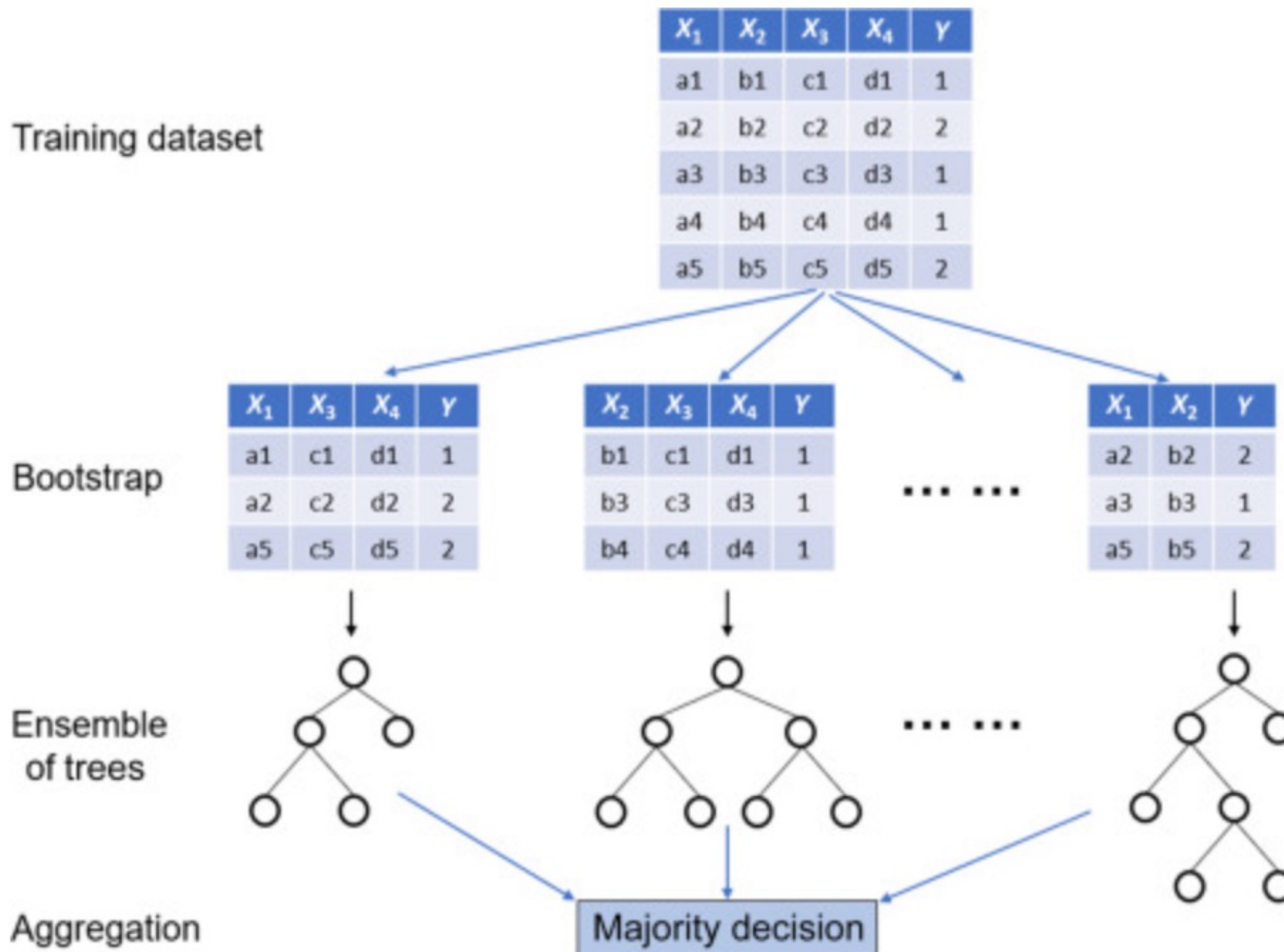
Data subsample 2..... Data subsample 1000



Build separate trees for each subsample (bootstrap) of houses. For each house, make a separate prediction for each tree (buy/not buy). Then make a decision for that house based on the majority rule (if the majority of trees let you the decision to buy that house), then buy it). This is called majority rule. In regression trees (continuous responses), we take the average of the predicted value for any observation of interest.

More complex models for building trees

Random forest: bootstrapping predictors



Classification and Regression trees

- Presenting a complex model as a tree that is easy to interpret is the key why CART became such a popular method.
- “There is no need to understand statistics to fit and interpret CARTs” ...but one should understand the basis to feel comfortable with the method and outputs.
- It treats data without a mechanism (as in OLS regressions, GLMs, etc); the thinking is in the algorithm and not about the mechanism that generated the response variable.