

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

Reading



What are decision trees?

Carl Kingsford & Steven I. 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

1

Learning from the data

Pattern recognition

2

Learning from the data

Machine learning algorithms - Two main types

Unlabeled data

➔

Unsupervised Learning Algorithm

➔

Prediction based on finding patterns in the data

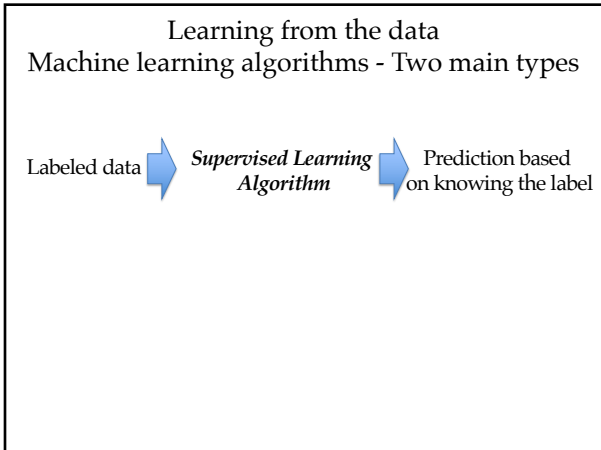
⬇

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

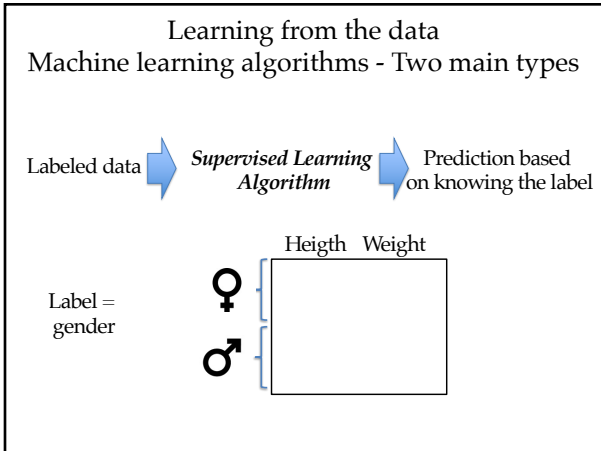
```

graph TD
    A[Height > 180cm] -- Yes --> B[Group 1]
    A -- No --> C[Weight > 80kg]
    C -- Yes --> D[Group 2]
    C -- No --> E[Group 3]
  
```

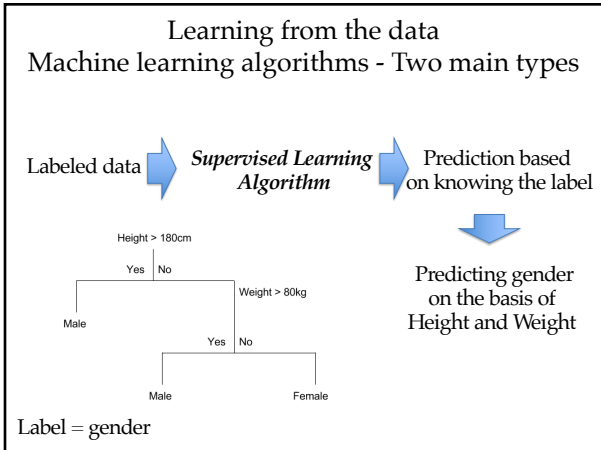
3



4




5



6

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



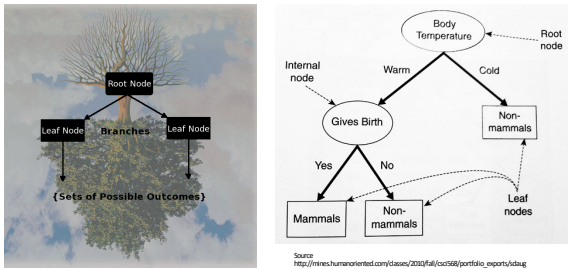
(Leo Breiman and colleagues 1984)

“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)

7

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)



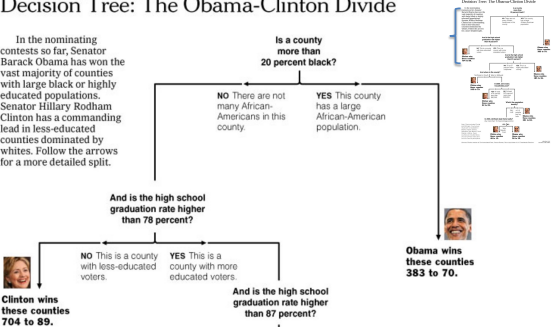
Source: http://rtech.humanoriented.com/issue/2012/fall/oct2012/portfolio_exports/ldag.html#decisiontree

8

Learning from the data – Classification Trees
Deal with complex data but easy to convey results

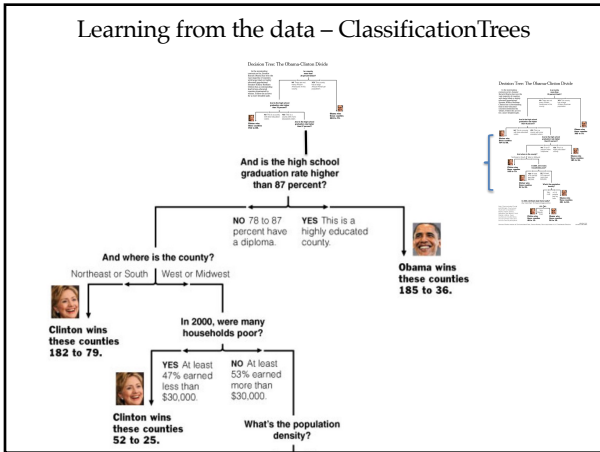
Decision Tree: The Obama-Clinton Divide

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.

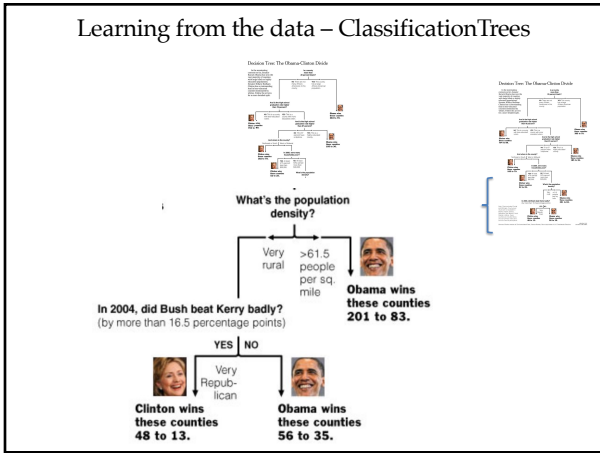


Source: Election results via The Associated Press; Census Bureau; Dave Leip's Atlas of U.S. Presidential Elections

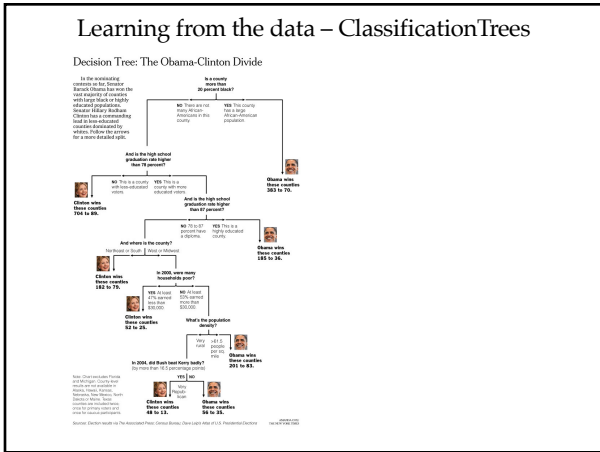
9



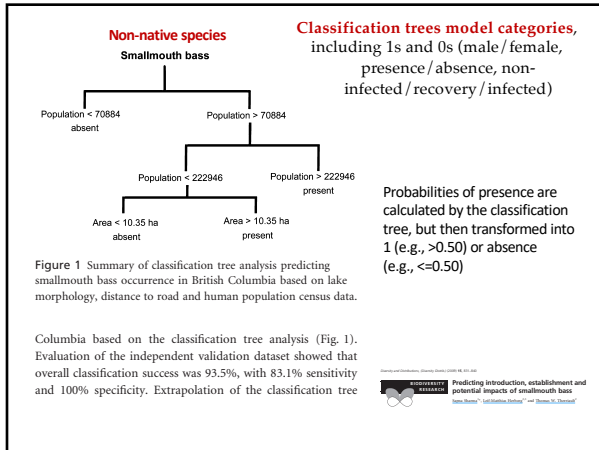
10



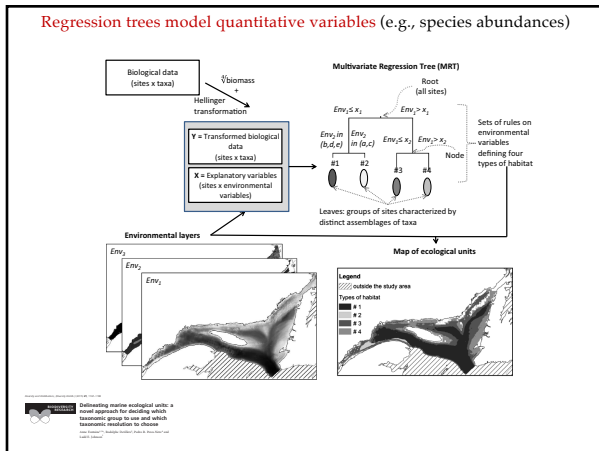
11



12



13




14

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.

15

The classification tree algorithm

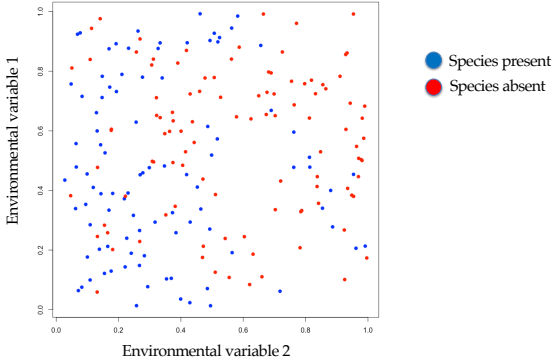


16

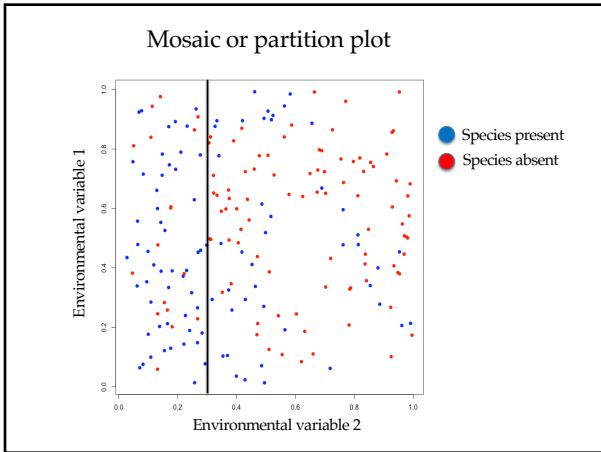


17

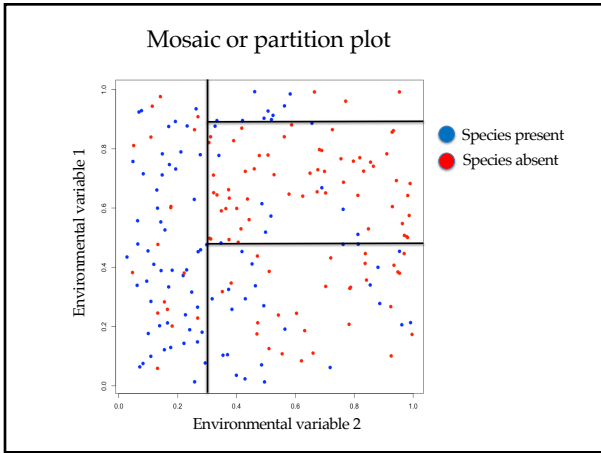
How to model these data?



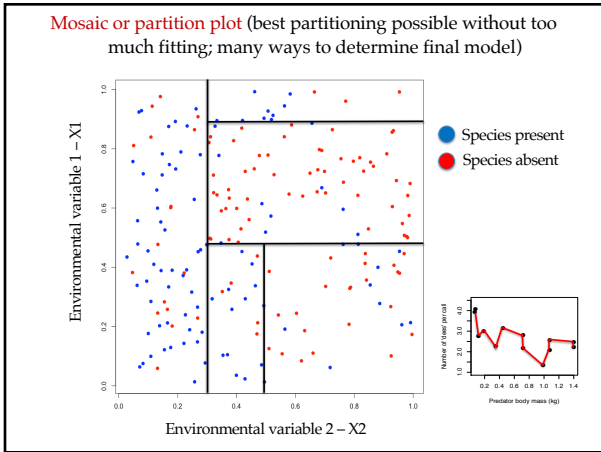
18



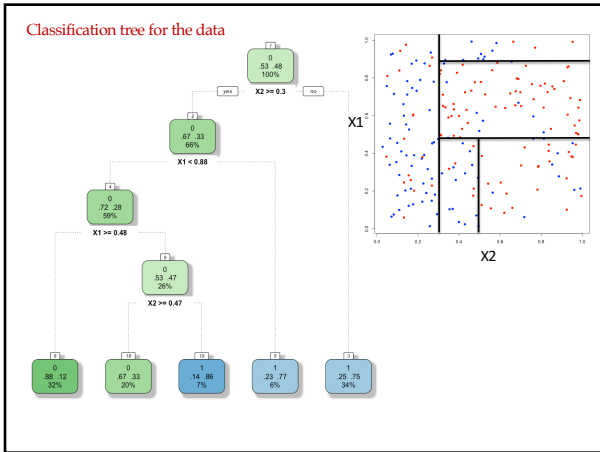
19



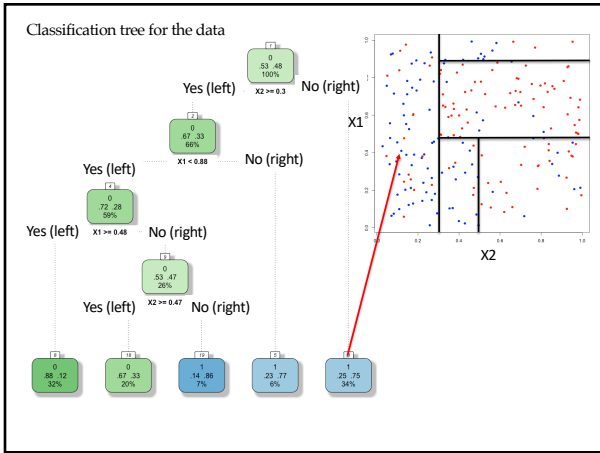
20



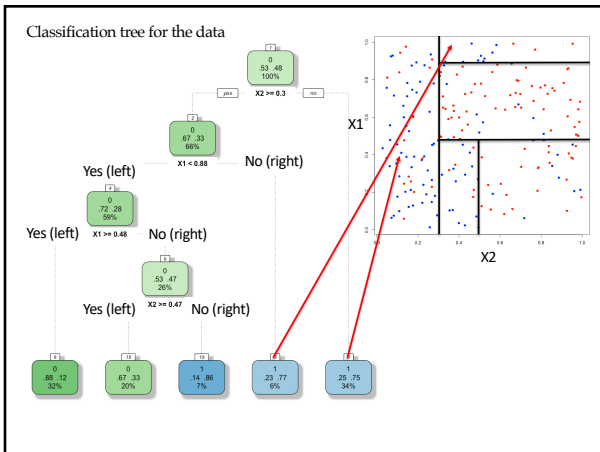
21



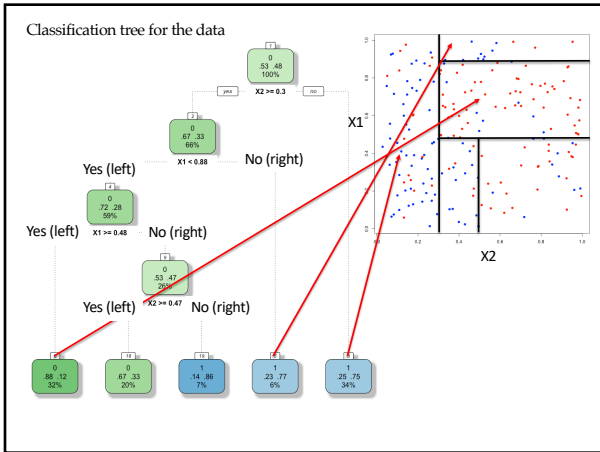
22



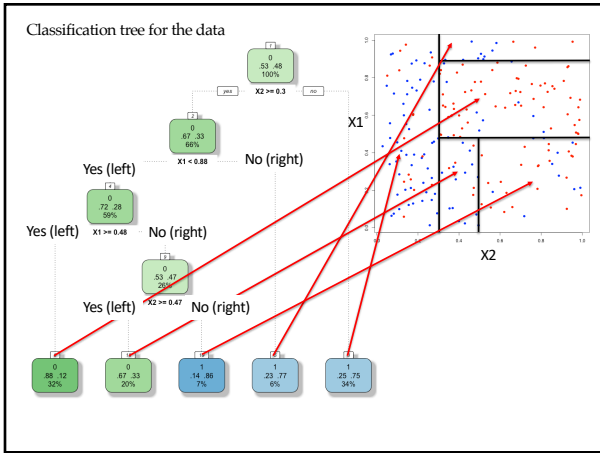
23



24



25



26

Den Boer et al. 2009

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

Marique L Den Boer¹, Marjon van Slagterhorst¹, René X De Menezes, Mayling H Cheuk, Jessica G C A M Bujs, Claudine Susan T C J M Peters, Laura J C M Van Zutven, H Bema Beverloo, Peter J Van der Spek, Gaby Eschschicht, Martin A Horstmann, Gritta E Janku-Schaub¹, Willem A Kamps¹, William E Evans, Rob Pieters¹

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.

The expression of 22283 genes across 190 patients were considered

Lancet Oncol 2009; 10: 135-144
Published Online
January 9, 2009
DOI:10.1016/S1473-0167(09)61270-9

27
