

Slide4

Friday, June 4, 2021 9:36 PM

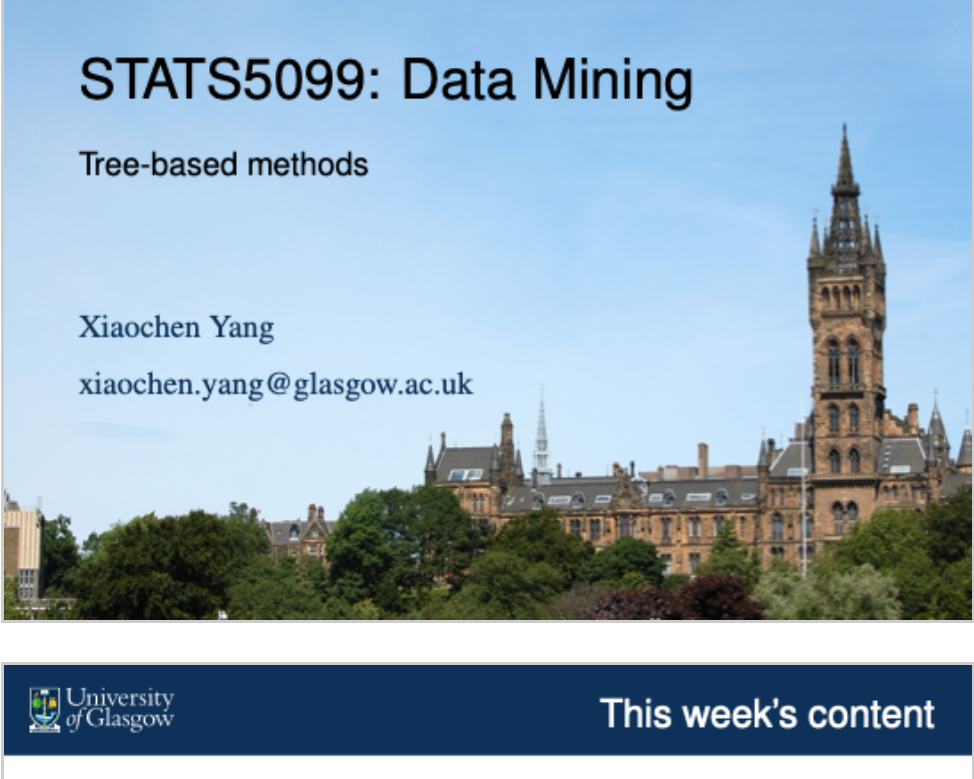
TutorialSlide4

University of Glasgow

STATS5099: Data Mining

Tree-based methods

Xiaochen Yang  
xiaochen.yang@glasgow.ac.uk



University of Glasgow

This week's content

- Classification trees
- Bagging and random forest

2/11

University of Glasgow

Classification trees

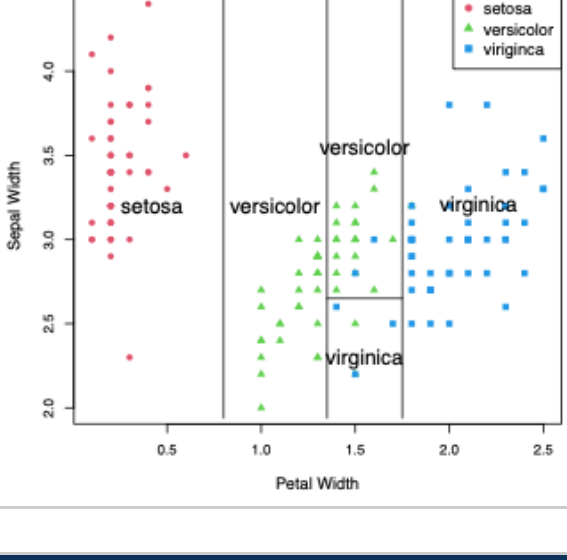
- Non-parametric classifier
- partition the feature space into a number of disjoint and non-overlapping regions
- predict the class of a given observation as the most commonly occurring class of training observations in the region to which it belongs

3/11

University of Glasgow

Classification trees

- partition the feature space into a number of disjoint and non-overlapping regions

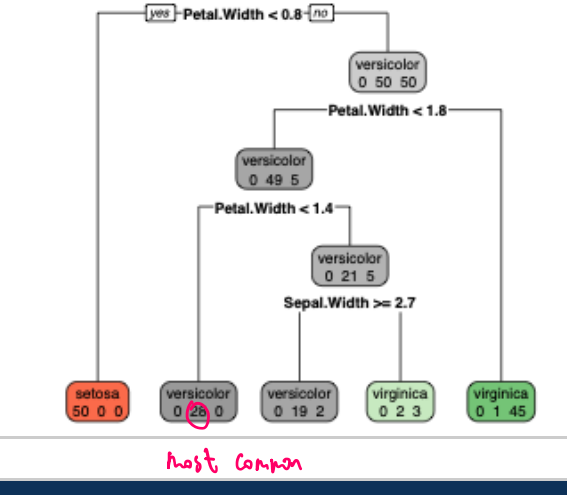


3/11

University of Glasgow

Classification trees

- predict the class of a given observation as the most commonly occurring class of training observations in the region to which it belongs



3/11

University of Glasgow

Step 1: partition the feature space

- Build a tree (recursive binary split)
- top-down approach: begin at the top of the tree, and successively splits the feature space
- binary split at each step
- best split: feature and cut-off point combination that leads to the largest possible reduction in error rate, Gini index or entropy

4/11

University of Glasgow

Step 1: partition the feature space

- Build a tree (recursive binary split)
- top-down approach: begin at the top of the tree, and successively splits the feature space
- binary split at each step
- best split: feature and cut-off point combination that leads to the largest possible reduction in error rate, Gini index or entropy

Avoid overfitting

- minimum number of observations in any terminal node
- maximum depth
- ...
- Pruning

4/11

University of Glasgow

Step 1: partition the feature space

- Build a tree (recursive binary split)
- top-down approach: begin at the top of the tree, and successively splits the feature space
- binary split at each step
- best split: feature and cut-off point combination that leads to the largest possible reduction in error rate, Gini index or entropy

Prune a tree

- Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of  $\alpha$ .

$$C_{\alpha}(T) = \sum_{m=1}^{|T|} Q_m(T) + \alpha|T|$$

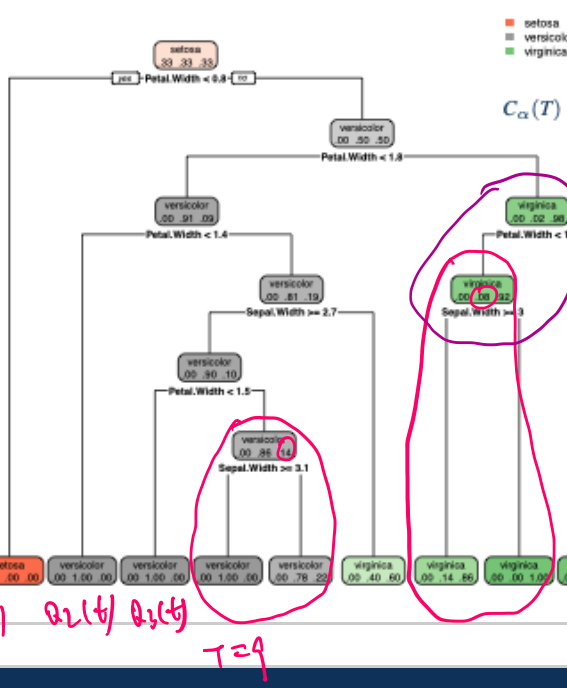
*Handwritten notes: Gini index / entropy rate ↓, # of terminals ↓, fit to data, i.e. (balance)*

- Use the validation set or K-fold cross-validation to choose  $\alpha$ .

4/11

University of Glasgow

Tree pruning

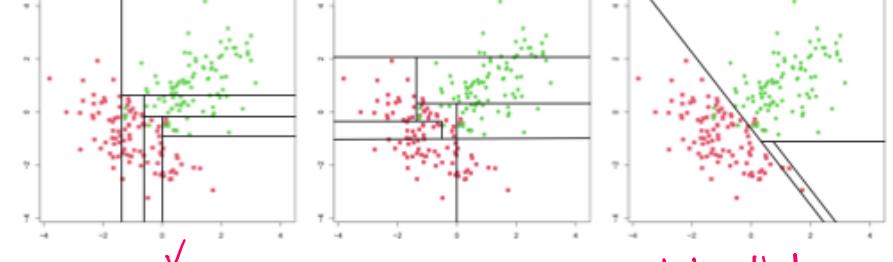


5/11

University of Glasgow

Conceptual question

Which of the following plot(s) can be generated by using recursive binary split?



6/11

University of Glasgow

Classification trees: summary

Pros:

- easy to explain, highly interpretable
- easily handle categorical variables and missing data

Cons:

- predictive accuracy is not very competitive
- very non-robust: a small change in the data can cause a large change in the final estimated tree

7/11

University of Glasgow

Bagging and random forests

- ensemble methods: combine the predictions from multiple models to make more accurate predictions than any individual model
- underlying statistical idea: averaging a set of observations reduces variance

Given a set of  $n$  independent observations  $Z_1, \dots, Z_n$ , each with mean  $\mu$  and variance  $\sigma^2$ , the mean of  $Z$  is still  $\mu$  and its variance is reduced to  $\sigma^2/n$ .

- bagging: average a set of models

8/11

University of Glasgow

Bagging and random forests

Bagging (bootstrap aggregation)

- generate  $B$  different bootstrapped training data sets
- build a prediction model  $f_b(x)$  on the  $b$ th bootstrapped training set
- average across all the models:

$$f_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B f_b(x)$$

*Handwritten notes: (without penalty), we large tree to bagging*

9/11

University of Glasgow

Bagging and random forests

Random forest

- generate  $B$  different bootstrapped training data sets
- build a prediction model  $f_b(x)$  on the  $b$ th bootstrapped training set using a random sample of features
- average across all the models:

$$f_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B f_b(x)$$

*Handwritten notes: (models are not correlated), make sure all models are independent*

9/11

University of Glasgow

Summary: considerations when choosing classifiers

- parametric vs nonparametric
- linear vs nonlinear decision boundary
- model complexity (tend to under- or over-fit?)
- computational cost (training, test)
- data: missing data, categorical features, correlated features  $\Rightarrow$  random forest

10/11

University of Glasgow

Tree-based methods in R

- classification tree:  
`class.tree`, `varImpPlot`, `printcp`, `prune`  
output: `variable.importance`
- bagging and random forest:  
`randomForest`, `varImpPlot`  
arguments:  
`mtry`: number of variables randomly sampled as candidates at each split  
`ntree`: number of trees to grow

11/11