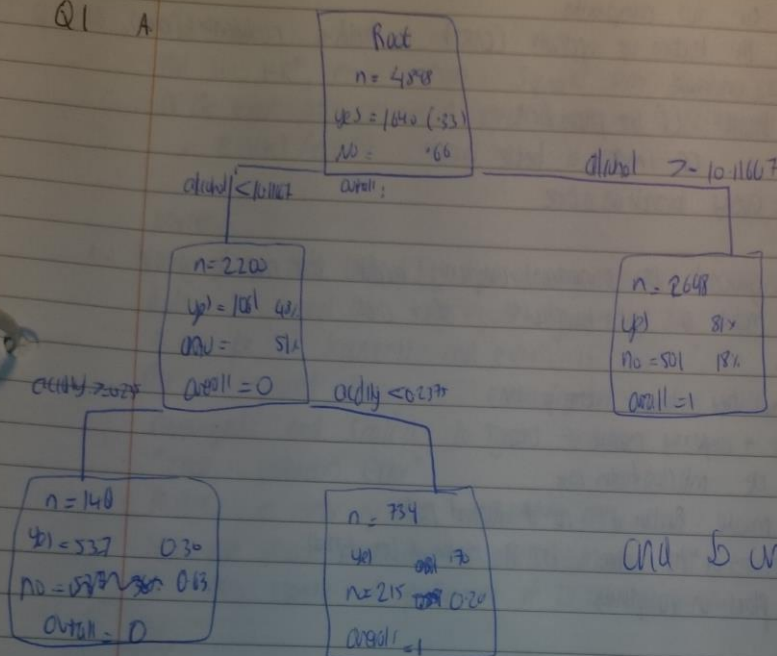


DA 2015 EXAM PAPER

Q1 A



and 5 on

First node - record the node number used for position in the tree
 Split - highlight the rule used for each split at each node
 n = number of cases in the particular node
 loss = number of misclassified cases in each node (similar category)
 y0 = true overall distribution attribute for that particular node
 y1 = true probability of being assigned to class 0 and 1 respectively
 * = denotes terminal nodes of the CART tree

B Explain CP, red error, xerr and xstd

CP

- CP stands for complexity parameter
- CP measures the complexity which is a penalty placed on

- each node for its complexity
- It is used in the bottom up approach (CART) to building classification tree, to pure branches or node
- Node with smaller CP are pure first.
- Bigger values of CP implies a better branch
- Calculated at every branch in a tree
- The value governs the minimum complexity benefit that must be gained in order to make a split worthwhile (default 0.01)
- Grow maximal tree - largest tree possible
- $R_{\alpha} = \text{cost} + \text{complexity measure of tree } T$
- Define cost as the misclassification rate
- Complexity measure: function of no. of terminal nodes
- $R_{\alpha} = R(T) + \alpha * |T|$ where $|T|$ is number of terminal nodes
- α penalty placed on complexity
- For a single node t : $R(t)_{\alpha} = R(t) + \alpha$
- For a subtree T_t : $R(T_t)_{\alpha} = R(T_t) + \alpha |T_t|$
- When α increases
 - * Both $R(t)_{\alpha}$ and $R(T_t)_{\alpha}$ increase but $R(T_t)_{\alpha}$ increases faster.
 - * Value of α when $R(t)_{\alpha} = R(T_t)_{\alpha}$ price paid for complexity - adding on to tree
$$\alpha = \frac{R(t) - R(T_t)}{|T_t| - 1}$$
- Pure tree with lowest value of α and recalculate
- CP is α except it is divided by $r(t)$ - misclassification rate for the next node
- Set a stopping size on CP, possibly 0.01 to determine final tree

29/04/16

DA 2015 EXAM

3

Q 10.

Rel Error

This is $1 - R^2$, root mean squared error, similar to linear regression this is the error on observations used to estimate the model

$$= RSS(n) / RSS(1)$$

Xerror

- Each of these trees is examined using 10-fold cross-validation, in which the data are divided into 10 equal segments; the tree is built using 9 of the 10 segments and error is assessed on the 10th segment. This is repeated leaving off each segment in turn and errors are then averaged and scaled to give Xerror
- "cross validation error"
- Multiply rel error by our misclassification rate
- Xerror rate of 0.9 means that the misclassification rate is 0.9 and misclassification of not nodes
- More realistic estimate of performance of the tree on new samples of data

Xstd

- Xstd is the variance between the 10 subsample estimates
- Ideally we pick the tree with the lowest Xerror $\pm 1.5 \times Xstd$
- rel error is estimated with the training data - the sample used for estimating the tree - and thus it decreases as the tree increases, because the tree becomes more and more adjusted to the data.
- This apparently better performance should not be taken for "real" when predicting for a new sample of data because larger trees tend to overfit on training samples and will hardly generalize well on new data samples
- Choose tree with 4 splits
- CP at 0.01
- Xerror and rel error at 0.01
- Xstd small value too

1c. Primary Split

- A primary split at a node is the rule which determines where the data should flow \rightarrow i.e. to left or right subtree
- This rule is used on all rules where the variable necessary to split on is available.
- A split is evaluated and determined as primary based on the concept of impurity
- A perfectly 'pure' split will assign all cases to one class, while a 'unpure' node will assign cases in a 50:50 split between classes
- The goodness of fit is defined to be the decrease in impurity
- Commonly used impurity functions are Gini or entropy
- It looks at the largest class in a dataset and strives to isolate it from the other classes.
- Look at impurity of parent node
- Each split gives two children
- Measure impurity of the two children
- Choose the split with the largest decrease in impurity - R calls this improvement.

Surrogate Split

- For each split node, the "primary splitter" is the variable that best splits the node maximizing the purity of the resulting child nodes
- When the primary splitting variable is missing for an individual observation, that observation is not discarded, instead a surrogate splitting variable is sought.
- A surrogate splitter is a variable whose pattern within the dataset, relative to the outcome variable, is similar to the primary splitter.
- Thus, the program uses the best available information in the face of missing values.
- The surrogate may have a different cut off point from the primary splitter but the number of cases the surrogate split sends into left and right nodes should be close to that of the primary split.
- By default, CART analysis produces 5 surrogate variables as part of its standard output.

10/04/16

2015 EXAM - DA

7

2 C Show how ROC curve is constructed

This is a plot of the true positive rate vs false positive rate for different possible cut off points of a diagnostic test.

- Need to calculate the TPR and FPR (Sensitivity vs 1-Specificity)
- Draw up a confusion table, using X as a cut off point
- Using the confusion table, derive Sensitivity = $\frac{TP}{TP+FN}$ and Specificity = $\frac{TN}{TN+FP}$
- The TPR = Sensitivity and FPR = 1-Specificity
- Repeat for multiple values of X and draw the curve

Shows:

- Trade off between sensitivity and specificity (Sensitivity increases, specificity will decrease)
- Closer the curve follows the left and upper border, means a more accurate model
- Area under the curve is a measure of test accuracy

D Error rate vs $P(+)$

- This is a plot of Error versus $P(+)$
- $P(+)$ is the proportion of + in population

- TPR and TNR are not dependent on prior probabilities but accuracy is:

$$Acc = P(+)*TPR + P(-)*TNR$$

$$Err = 1 - Acc = 1 - (P(+)*TPR + P(-)*(1-FPR))$$

$$= P(+)*FNR + P(-)*FPR$$

$$= (FN-FP)*P(+) + FP$$

- Point in ROC space are going to be mapped and lies in "error/ $P(+)$ " space
- Allows us to see very clearly how error rate varies according to $P(+)$
- Diagram the response curve when test always positive and goes to (1,1) line (grey) and when always negative goes to (0,0) line (grey)

Rule: $Err = FP, P(+)=0$ $Err = 1-TP, P(+)=1$

- Cost curve plot has error rate versus as a function of the prevalence of positive examples
- Generally this idea when misclassification costs are not equal
- There is a point/duality between ROC space and cost space meaning that a point in ROC space is represented by a line in cost space and vice versa
- A classifier's operating range can be immediately read off a cost curve - it is defined by the RFE value where the cost curve intersects the diagonal line representing trivial classifiers
- Costlier perform worse closer to middle band

Q3 A What is an ensemble? Advantages?

- Ensembles (or committees) of machine learning methods that we use to power multiple models to achieve better prediction accuracy than any individual model can on its own
- Eg. 100 classifiers with an individual error rate of 0.4 have 0.03 error rate overall
- Made up of 1000+ etc less accurate models to best predictive power
- The output of an ensemble is determined by voting system from model in an ensemble provides an output and the most popular answer is the overall output of the ensemble

Advantages

- Higher accuracy
- Many types of variable output
- More outputs than just misclassification rate (classification etc)
- Can be proximal for MPS
- Less overfitting
- Unlikely that all classifiers will make same mistake
- So long as each error is made by a minority of classifiers, optimal classifiers will be chosen
- Greater predictive power compared with individual model
- Random Forest, Boosting, Bagging, Rulefit

30/04/16

EXAM 2015 ON

9

Q3 B. Difference between Boosting and bagging

Bagging - Bootstrap aggregation

- Generate B bootstrap samples of the training data: random sampling with replacement
- Train a CART tree using each bootstrap sample
- For classification - majority vote
- For regression - average of predicted values
- Bagging minimises error for bias and variance
- Way to decrease the variance of your prediction by generating different data
- By increasing size of your training set, you can improve the model predictive but can decrease variance, normally being the prediction to expected value
- PARALLEL aim to decrease variance not bias

Boosting

- Iteratively learning weak classifiers
- Technique for combining multiple base classifiers whose combined performance is significantly better than that of any of base classifiers
- Sequential training of weak learners
- Each base classifier is based on data that is weighted based on the performance of the previous classifier
- Each classifier votes to obtain final outcome
- SEQUENTIAL - Aim to decrease bias, not variance

Boosting Adv/Disadv

- Powerful classification algorithm
- Can achieve similar classification results with much less knowledge of parameters or settings
- Can be sensitive to noisy data and outliers
- Less susceptible to overfitting problem
- Feature selection resulting in relatively simple classifier
- Instead of re-sampling, uses training set re-weighting
- Usually a sub-optimal solution

- Lose the simple interpretability of classification trees
- Computation more difficult, slower
- No prior knowledge needed about weak learners
- No parameters to the extent of - number of trees
- From empirical evidence, AdaBoost is particularly vulnerable to overfitting
- Boosting performs an exhaustive search for best predictor to split on, whereas RF model only searches a subset of data
- Boosting grows trees in series, with later trees dependent on the result of previous trees, RF grows in parallel, independent of one another

3 C Explain How Random Forests are constructed

- We assume that we know about the construction of single classification trees
- RF is an ensemble method based on classification trees
- Numerous trees are built using different training sets so as to ensure different results
- The ensemble can predict accurately to majority vote

Each tree grown as follows:

- 1- If # of cols in training set is N , Sample N cols at random but with replacement, from original dataset. This sample will be dataset for growing the tree
 - 2- If there are M input variables, a number $m \leq M$ is specified such that at each node, m variables are selected at random out of M and the best split on these m is used to split the node. The value of m is held constant during the forest growth
 - 3- Each tree is grown to the largest possible extent - No pruning
- Reducing m reduces both the correlation and the strength
 - Using out of bag error rate to get value of m
 - Out of bag error rate, around 36% of data not used - used to calculate error rate

DA - 2015 EXAM

1.C.

Improvement:

- Related to CP
- R call to decide in "improving improvement"
- How good to additional Split is correctly dividing the data
- Measure sum of GINI or Entropy

Adj:

- Called Association
- $(\text{default mismatch} - \text{surrogate mismatch}) / \text{default mismatch}$
- Simply a measure of similarity
- Relative reduction in error obtained by using surrogate Splitter to predict or instead of using data % for left and right
- Percentage reduction in error obtained by using surrogate to predict the primary variable

Agreement

- The amount that both Splitters sent cases in to same direction
- Used for evaluating quality/accuracy of surrogate Splitter against primary Splitter.

2 A Confusion Table

- Also known as an error matrix, is a specific layout that allows visualization of the performance of an algorithm
- Each column of the matrix represents the instance in an actual class
- None seems from the fact that it makes it easy to see if the system is confusing two classes
- Gives misclassification rates of a given model
- Constructed by running the test data through the model then comparing the predicted class to the actual class, creating the 4 box
- Often used to determine the performance of a classification model on a set of test data for which the true values are known.

	Predicted Yes	Predicted No
Actual Yes	TP	FN
Actual No	FP	TN

2.6 Assessing A model's performance

- A list of statistics which can be calculated from the confusion matrix.
- Accuracy: Overall, how often is the classifier correct?

$$= (TP + TN) / (TP + TN + FP + FN)$$
- Misclassification rate / Error rate: Overall how often is it wrong?

$$= (FP + FN) / (TP + TN + FP + FN)$$

$$= 1 - \text{accuracy}$$
- True Positive Rate: When it's an actual yes, how often does it predict yes?

$$= TP / (TP + FN)$$
 True Positives / Actual yes
 • Also known as sensitivity or recall.
- False Positive rate: When it's actually no, how often does it predict yes?

$$= FP / (FP + TN)$$
 FP / Actual no's
- Specificity - When it's actually no, how often does it predict a no?

$$= TN / (FP + TN)$$
 TN / Actual no's
 • equivalent to $1 - \text{FPR}$
- Precision - When it predicts a yes, how often is it correct?

$$TP / (TP + FP)$$
 TP / Predicted yes
- Prevalence - How often does the yes condition actually occur in our sample?

$$\text{Actual yes / total} = \frac{TP + FN}{TP + TN + FP + FN}$$
- Balanced Accuracy

$$(\text{sensitivity} + \text{specificity}) / 2$$
- Run a McNemar test which follows a Chi-Square distribution.
- Kappa Statistic - agreement between raters/analysts
 - how do agreement taking into account to accuracy that would be generated by chance

$$\text{Kappa} = \frac{O - E}{(1 - E)}$$
 O: observed accuracy
 E: expected accuracy
 Value between 0 and 1, larger is better

2015 EXAM DA

11

3D PARTIAL DEPENDENCY PLOT

- Shows how each predictor is related to response, holding other variables constant
- Let x be the initial predictor of interest with v different values
- Construct v datasets for each of the v values of x leaving all other variables constant
- For the v datasets, predict the response using random forest
- Calculate a single value averaged over all observations for each dataset
- Average the predictions over trees
- Enable us to visualize interaction between 2 variables
- If each explaining variable contributes additively to the target y , by some unknown function, this method is great to show that estimated hidden bias
- For 2 categories - plots will be mirror images of each other

3E RF vs Rulefit

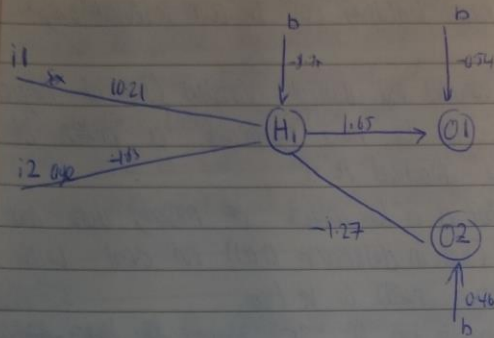
- A mixture of regression and trees. RF only uses trees
- RF are good at picking up linear relationships, but rulefit does (trees for non-linear relationships and linear regression for detecting linear relationships)
- Rulefit grows a forest of trees like RF but then uses individual from the variables to form a logistic regression model.
- The important model that rulefit uses includes a combination of the importance from trees and linear regression. RF just uses trees
- First component produces 'rules' and the second component fits a linear model with all the inputs
- Highly interpretable, decision rules have an easily understandable form

05/05/16

DA 2015 EXAM

12

Q4 A. Draw neural net, probability of adult male lizard



Age no ye) gender = 1
0 1 Age = 0

ADULT MALE

$$A1 = -8.75 + 10.21(1) + (-1.85)(0) = 1.46$$

$$F(O1) = \frac{e^{1.46}}{1 + e^{1.46}} = 0.8115$$

$$O1 \Rightarrow 1.65(0.8115) - 0.54 = 0.799$$

$$F(O1) = \frac{e^{0.799}}{1 + e^{0.799}} = 0.681$$

$$O2 \Rightarrow -1.27(0.8115) + 0.46 = -0.57065$$

$$F(O2) = \frac{e^{-0.57}}{1 + e^{-0.57}} = 0.361$$

- Category no male with larger value $\Rightarrow 0.799$

ADULT FEMALE

$$A1 = -8.75$$

$$F(O1) = \frac{e^{-8.75}}{1 + e^{-8.75}} = 0.000158436$$

$$O1 \Rightarrow 1.65(0.00015) - 0.54 = -0.54$$

$$F(O1) = \frac{e^{-0.54}}{1 + e^{-0.54}} = 0.362$$

$$O2 \Rightarrow -1.27(0.00015) + 0.46 = 0.454$$

$$F(O2) = \frac{e^{0.454}}{1 + e^{0.454}} = 0.6124$$

4.3 Missing Data

Raw Data will almost always need to be cleaned before analysis to find and deal with erroneous and missing data. To find data that will be potentially damaging to model building can be done in several ways:

- Counts of NA's, blank variables and illogical variables (such as negative age etc) can give an indication as to how significant a problem of missing data is and also how localized it is.
- Could be the case that there is a low amount of missing data and the model being constructed (such as decision trees) can deal with such a problem so no direct action needs to be taken.
- Trees can theoretically handle up to 25% missing data and remain effective. Several rules exist on how data might be dealt with.
- One way of dealing with missing data is to impute a value to fill in the gap. Replacing a missing value with a mean or median value can neutralize the effects of missing data but also can bias the results towards the center. If there is a lot of data that has been imputed to the mean, the standard deviation of the data will be falsely low and the variance may no longer be realistic.
- If there are a large number of cases of missing the same variable, it may be possible to remove variables entirely.
- It may be preferable to use a surrogate variable, or in the case of CART a Surrogate Split. These are splits that give similar results and should be used when the information is difficult or too expensive to gain from the original variable (which could explain why there is so much missing).
- If a case is missing data, the entire case could be deleted. Instance deletion with sufficiently large data sets, this will have a negligible impact on the ultimate data. Removing too many observations (and reducing sample size) will reduce the confidence we can place in our model and ultimately its predictions.
- NN cannot deal with missing data, and as such more extensive data such as feature or case deletion may be necessary to ensure the model can be accurately trained.

01/05/16

DA: 2015 EXAM PAPER

14

CART vs LOGISTIC REGRESSION

Data

- CART analysis allows you to perform analysis on a highly skewed messy dataset
- LR can be affected by messy data and outliers
- CART can handle both categorical and quantitative data, whereas LR can only handle binary categorical data
- CART tends to work better with local datasets, whereas LR tends to work better globally (parametric vs non-parametric)

Detection

- CART excel at the detection of local structure. Each half of the tree is analysed separately. CART automatically disregards insignificant variables
- CART automatically identifies interactive terms
- CART tend at detecting linear structure and cannot represent it effectively
- LR good for linear relationships. Many non-linear structures can be reasonably approximated with a linear structure
- LR need to be adjusted manually with different combinations of variables - no detection of interactive terms

Assumptions

- The decision tree assumes the splits are axis parallel and will become more complex with the increase in number of features and multidimensional hypothesis are possible
- LR assumes there is only one decision boundary that is smooth and non-linear.

Overfitting

- Complex decision trees may overfit the data and trees will become unstable could pose a problem
- LR is intrinsically simple, it has low variance and so is less prone to overfitting

Speed

- While both algorithms are fast, LR has a faster run time on large data sets than decision trees

Interpretability

- CART and LR - nice probabilistic interpretation
- CART easy to visualize and interpret for non-statisticians

Transformations

- CART does not require transformation of log or square root etc
- LR uses log or probit function

How can they be combined

- Run a shallow tree
- Assigned each case to a terminal node
- Treat terminal nodes as categorical variables
- Feed data into LR model