

06/10/15

## DATA ANALYTICS

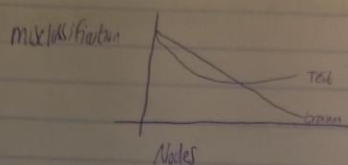
### CLASSIFICATION AND REGRESSION TREES

- like Logistic regression?
- issue : how to derive splits  
grow the tree?  
what size tree?
- Do it by looking at all possible splits
- Deciding the split:
- Split with higher classification rate the "best"

8/10/2015

13/10/15

## DATA ANALYTICS



Drop in sample of data into each tree and get estimate of misclassification  
Split data into say 10 parts and calculate misclassification rate, choose lowest.

Tree opposite to Logistic reg. trees good for interaction  
Logistic reg good for linear relationships

### COMPARISON

- Simple to capture linear relationship with logistic regression
- Trees not good at linear relationships
- Tree  $\rightarrow$  Can use one variable many times.

### INTERACTION

When level depend on effect of another variable - level of  $y$  from  $x_1$  depends on  $x_2$   
will have something like  $y = x_1 \beta_1 + x_2 \beta_2 + x_1 x_2 \beta_3$

$\uparrow$  captured interaction

- Trees automatically include interactions

No interaction - just look at the marginal.

CASE 2 example: No interaction no relationship

No relationship/interaction: tree will give just the root node

CASE 3:  $x_1$  has an effect, doesn't depend on level of  $x_2$  no interaction  
tree: just one split

16/10/15 DATA ANALYTICS EXTRA NOTES

$N_j$  Number in class  $j$  overall

$N(t)$  Total number of cases in node  $t$

$N_j(t)$  Number of class  $j$  cases in node  $t$

Proportion of class  $j$  cases in node  $t = N_j(t) / N(t)$

$\pi_j$  = prior probability

$$\sum \pi_j = 1$$

Probability case will be both in class  $j$  and in node  $t = p(j, t) = \pi_j \frac{N_j(t)}{N_j}$

Probability of falling into class  $j$  given that a case is in node  $t$

$$p(j|t) = p(j, t) / p(t)$$

Probability any case fall into node  $t = p(t) = \sum_j p(j, t)$

Two common impurity measures

when two categories in target variable:

$$\text{Gini} = 2 * p(1|t)p(2|t)$$

How to Choose Split?

- Look at impurity of parent node  $t: i(t)$
- each split gives 2 children  $t_L$  and  $t_R$
- Measure impurity of the two children

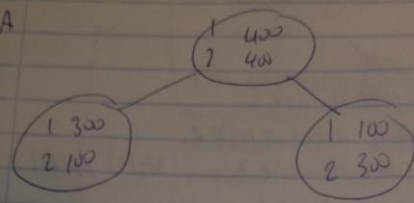
$$i(t) = p_L * i(t_L) + p_R * i(t_R)$$

Smaller value of  $i(t_L)$  and  $i(t_R)$  are better.

- Choose Split with largest  $\Delta(t, s) \rightarrow R$  could be important

$R$  multiplied  $\Delta(t, s)$  by  $N$  - total no. of cases

A



$$G_{\text{root}} = 2 \times p(1|t) p(2|t)$$

$$= 2 \times (0.5)(0.5)$$

$$= 0.5$$

$$\text{Impurity} = i(t) - p_L \times i(t_L) - p_R \times i(t_R)$$

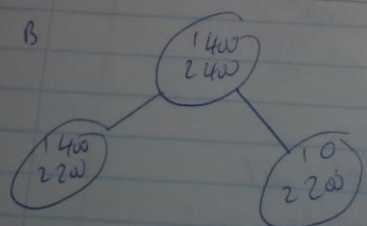
$$0.5 - 0.5 \times (2 \times 0.75 \times 0.25) - 0.5 \times (2 \times 0.25 \times 0.75)$$

$$0.5 - 0.1875 - 0.1875$$

$$= 0.125$$

CORRECT

B



$$G_{\text{root}} = 2 \times p(1|t) p(2|t) = 2 \times 0.5 \times 0.5 = 0.5$$

$$\text{Impurity} = i(t) - p_L \times i(t_L) - p_R \times i(t_R)$$

$$= 0.5 - 0.75 \times (2 \times \frac{4}{6} \times \frac{2}{6}) - 0.25 \times (2 \times 0 \times 1)$$

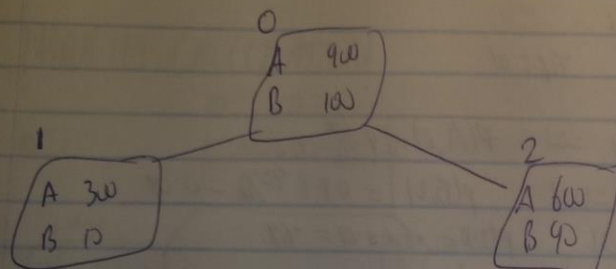
$$0.5 - \frac{1}{3}$$

$$S(t,s) = \frac{1}{6} \quad (0.1666)$$

Pick <sup>high</sup> vote if  $\Delta(t,s) \rightarrow$  ~~first~~ Second split

CORRECT





$P(A|0) = 9/10$   
 $P(B|0) = 1/10$   
 $P(A|1) = 300/310$   
 $P(B|1) = 10/310$   
 $P(A|2) = 600/690$   
 $P(B|2) = 90/690$

Equal prior = multiply by 0.5 to make fair

$P(A|1) = 0.5 \cdot 300/900 = 0.17$   
 $P(B|1) = 0.5 \cdot 10/100 = 0.05$   
 $P(A|2) = 0.5 \cdot 600/900 = 0.33$   
 $P(B|2) = 0.5 \cdot 90/100 = 0.05$   
 $P(1) = 0.17 + 0.05 = 0.22$   
 $P(2) = 0.33 + 0.05 = 0.38$

$P(A|1) = 0.17/0.22 = 77$   
 $P(B|1) = 0.05/0.22 = 0.23$   
 $P(A|2) = 0.33/0.38 = 87$   
 $P(B|2) = 0.05/0.38 = 13$

Assign A or B whichever has higher value

4

Data of priors

$$\pi_A = 0.4 \quad \pi_B = 0.1$$

$$P(A,1) = 0.4 \cdot \frac{30}{90} = 0.3$$

$$P(A,2) = 0.4 \cdot \frac{60}{90} = 0.6$$

$$P(B,1) = 0.1 \cdot \frac{10}{90} = 0.01$$

$$P(B,2) = 0.1 \cdot \frac{90}{90} = 0.09$$

$$p(1) = 0.3 + 0.01 = 0.31$$

$$p(2) = 0.6 + 0.09 = 0.69$$

$$P(A|1) = \frac{0.3}{0.31} = 0.97$$

$$P(A|2) = \frac{0.6}{0.69} = 0.87$$

$$P(B|1) = \frac{0.01}{0.31} = 0.03$$

$$P(B|2) = \frac{0.09}{0.69} = 0.13$$

Node 1 and 2 to A (higher values)

MISCLASSIFICATION

$$r(t) \text{ prob of misclass at node } t = 1 - \max(p(i|t))$$

$$R(t) = r(t) \cdot p(t)$$

$$\text{for the tree } R(T) = \sum r(t) p(t) \quad T \text{ is terminal node}$$

Score for with data of priors

$$\text{Node 0: } r(t_0) = 1$$

$$R(t_0) = r(t_0) \cdot p(t_0) = 1 \cdot 0.7 = 0.7$$

$$\text{Node 1: } r(t_1) = 1 - 0.97 = 0.03$$

$$R(t_1) = r(t_1) \cdot p(t_1) = 0.03 \cdot 0.31 = 0.0093$$

$$\text{Node 2: } r(t_2) = 1 - 0.87 = 0.13$$

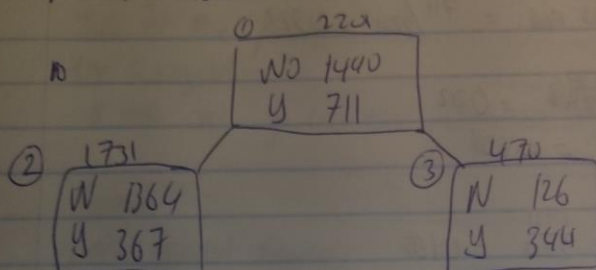
$$R(t_2) = r(t_2) \cdot p(t_2) = 0.13 \cdot 0.69 = 0.0907$$

$$R(T_c) = 0.0093 + 0.0907 = 0.099$$

16/10/19. DATA ANALYTICS EXTRA NOTES

5

## IMPURITY CALCULATIONS



$$\text{Impurity root} = 2^x p(N|\text{root}) * p(Y|\text{root})$$

$$2^x \frac{1440}{2201} * \frac{711}{2201} = 0.437$$

$$\text{Node 2: } i(t) = p_2 * i(t_2) - p_R * i(t_R)$$

$$i_{\text{node 2}} = 2^x \left( p(N|\text{node 2}) * p(Y|\text{node 2}) \right)$$

$$2^x \left( \frac{1364}{1731} \right) \left( \frac{367}{1731} \right) = 0.33451$$

$$\text{Node 3} = 2^x \left( \frac{126}{470} \right) \left( \frac{344}{470} \right) = 0.3924$$

$$p(1) = 1$$

$$p(2) = \frac{1731}{2201} = 0.786$$

$$p(3) = \frac{470}{2201} = 0.2135$$

$$\text{Decide in impurity} = i(\text{root node}) - [p(\text{proportion left}) * i(\text{Left}) + p(\text{proportion right}) * i(\text{R})]$$

$$0.437 - [0.786 (0.334) + 0.214 (0.392)]$$

$$= 0.091$$

### (CALCULATIONS FOR CP

$$R(0) = \text{missed} / \text{or not read} = 711 / 2201 = 0.323$$

$$r(1) = 367 / 1731 \text{ } \cancel{NV} \text{ } \cancel{1731} = 0.212$$

$$r(2) = 126 / 473 = 0.268$$

$$R(1) = \frac{r(1)}{0.212} \times \frac{1731}{2201} = 0.167$$

$$R(2) = \frac{r(2)}{0.268} \times \frac{473}{2201} = 0.057$$

$$R(\text{Tree}) = R(1) + R(2) \quad 0.167 + 0.057 = 0.224$$

$$CP = \frac{R(0) - R(\text{Tree})}{\# \text{ nodes} - 1} = \frac{0.323 - 0.224}{2 - 1} = 0.99$$



20/10/15

## DATA ANALYTICS

1(1)

### Missing Data

- look at missing packages
- Use Surrogates -
- if most cases go left  $\rightarrow$  send missing case to left.

### Surrogate

- possible equivalent information
- A split that splits exactly like primary split
- Finding a split that splits in a similar way to primary split.
- Done at every node in the tree
- What is probability both splits send data in same way - "agreement" in R
- Default rule - send down side with majority of data.

Large MSE -  $\sigma$  factor being treated as quantitative - error.

Example: improvement in order  $\rightarrow$  less # of missing cases

If missing the variable in normal split it will split by the extra variable in extra split.

Adjustment is same as association

Surrogate node dependent.

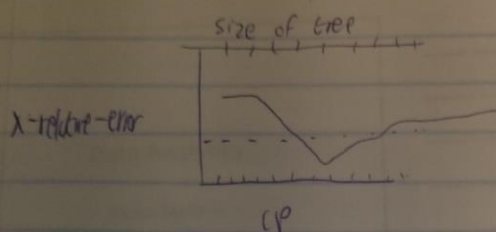
Be Skeptical of variable importance - different packages have different values

## 15 DATA ANALYTICS

Sometimes you might want a high sensitivity value and lower specificity value  $\Rightarrow$  i.e. want more true positives than true negatives

May want to weight wrong classifications

23/10/15 DATA ANALYTICS



- Plot of complexity parameter CP - values go down
- dotted line - top end of standard deviation of  $\lambda$ -rel-err
- cross validation error - relative to the root node
- Put CP=0 for biggest value

RCC for random model is a 459 1/2

How to compare 2 ROC curves? - area under curve?

- Draw curve on top of each other

why no splitter - may not be a good enough variable available  
convert area to factor summary(dataframe) will give frequencies  
not min or max

27/10/17

# DATA ANALYTICS

1(1)

## Model Evaluation

looking at "operating condition" - prior probability etc.

RCC curve rarely actually curved - Assume were looking area under jagged curve  
 (1+) prior probability of a plus one

When overall class % is very small - model tends to over emphasise it.  
 Convex hull - looking for the slope of tangent being the same

Graph - Two models, A and C.

FN  $\rightarrow$  FNR - both the same - write wrong on slides

FP  $\rightarrow$  FPR

New Graph - got an outline at bottom of graph with  
 lower possible error rates

How model behave for prior probability.

Costs important

		Predict	
		+	-
Actual	+	✓	x
	-	x	✓

So far we have model

X's the same  $\rightarrow$  may

want to weight it.

want to know ratios of x

Plot cost v cut off point for different cut off points



2 (1)

### COSTS

- Total cost = numbers \* cost
- Can incorporate cost into growing or pruning trees
- The two graphs - opposite of each other cost of false pos and false negative
- Only the ratio of costs is important.
- GINI - measure impurity
- Cost is basically a weighting variable

Ecov - expected cost used in ROCs as well

Max value - FN=1, FP=1.

Telling is ranges in which model will work

## 01/15 DATA ANALYTICS

CAGEL CAGEL

Predicted values in two vectors (yes, no)

The graph of trees v regression

- Should expect something around a diagonal

- (no) model (not) better, most of the time or high

Decides - RFP in the R graph, different to my own way

What are you going to use for classification?

- Depends what question you ask / are answering

- Purpose of study

- Classification tool?

- What do we want to optimise?

03/11/15

## DATA ANALYTICS

When fitting logistic regression - assuming a linear relationship

Seeing if there is a linear relationship in the variables

1 2 3 4 5 → split data variable into 5 sectors or 10 etc

yes  
no

### Corat Package Output

predicted v actual (reference) values

Accuracy: main diagonal over total

95% CI - for accuracy, 'good idea' CI:  $p \pm 1.96 \sqrt{\frac{p(1-p)}{n}}$  <sup>we shall get</sup> <sub>naughty</sub>

or bootstrap and take percentiles

wide b/c sample size is small

No information rate: how good do without model if 'predict everything as yes'

D-value: how better you did model than without it.

Kappa:  $\frac{\text{correct} - \text{actual}}{\text{total}}$  agreement controlling for chance

McNemar Test p-value: paired chi-squared test, suggests <sup>looking at 6 and 5</sup> no difference in example probability of 1.

Sensitivity:

Specificity: need to know which one is the event M or R or T or F or Y or N

$H_0: A_1 \leq NIK$  v  $H_1: A_1 > NIK$

evidence against  $H_0$ , accept alternative hypothesis.

Pos pred value: percentage of predicted yes over all yes.

Neg pred value: how at risk instead of not at risk

Prevalence: Prevalence of event 27/151

} better for medical view

03/11/15

## DATA ANALYTICS

When fitting logistic regression - assuming a linear relationship

Seeing if there is a linear relationship in the variables

1 2 3 4 5 ← split data variable into 5 sector or 10 etc

yes  
no

### Corat Package Output

predicted v actual (reference) value

Accuracy: main diagonal over total

95% CI - for accuracy, 'good idea' CI:  $p \pm 1.96 \sqrt{\frac{p(1-p)}{n}}$  <sup>we should get</sup> <sub>naughty</sub>

wide b/c sample size is small

No information rate: how good do without model if predict everything as yes

D-value: how better you did model than without it.

Kappa: <sup>(correct)</sup> actual yes / total agreement controlling for chance

McNemar Test p value: paired chi-squared test, suggests no difference in example probability of 1. <sup>looking at 6 and 5</sup>

Sensitivity:

Specificity: need to know which one is the event M or R or T or F or Y or N

$H_0: Acc \leq NIR$  v  $H_1: Acc > NIR$

evidence against  $H_0$ , accept alternative hypothesis.

Pos pred value: percentage of predicted yes over all yes

Neg pred value: look at rows instead of columns

Prevalence: Prevalence of event 27/51

} better for medical uses



McNemar - two models, cross tabulation of classifier  
matched - some observations in both models  
model A better than model B? look at yes, no and no, yes  
if  $\text{size} = n$   
 $df = 1$

27/11/15

## DATA ANALYTICS

Predicting probabilities from neural net

```
predict (fitnn, data.frame(class="", sex="", age=""))
```

Ensemble Methods Foundation, algorithm Zhou, Zhuohua