DA ROC Cost Curves

## ROC and Priors

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

$$Acc = p(+)^* TPR + p(-)^*(1-FPR)$$

$p(+)$ = proportion of + in population written as $\pi$, up to now

Can draw lines of iso accuracy on ROC curve

$$TPR = \frac{Acc - p(-)}{p(+)} + \frac{p(-)}{p(+)} \cdot FPR$$

- We are looking for the point on the convex hull where the slope of the tangent
  $= p(-)/p(+)$
- Operating conditions $p(+)/p(-)$
- Difficult to imagine slopes of tangents

## A new Space

- Going to plot Error vs $p(+)$
- $Err = (FN-FP)^* p(+) + FP$
- New space  Y-Axis: Error    X-Axis = $p(+)$
- Points in ROC Space are going to be mapped onto lines in "error/$p(+)$" space
- Allow us to see how error rate varies according to $p(+)$
- How the model will behave under various conditions.
- Easier to see when $p(+)=0$ and $p(+)=1$
- Gives points $(0, FP)$, $(1, (-TP)$ for each TP and FP combination
- Always negate rule (classify all as negate) $FP=0=TP$  $(0,0)(1,1)$
- Always positive rule $TP:FP=1$  $(0,1)(1,0)$
- look at the lower envelope of the graph
  the ratio between diagonals shows range by operable in
  If lines are above diagonals → always negate/positive will be ideal → the model
  is useless

## Costs

- Incorporate cost information.
- Total cost $= TP^* C(+|+) + FN^* C(+|-) + FP^* C(-|+) + TN^* C(-|-)$
- Objective: minimize cost.
- $C(+|-)$ cost of misclassify $+$ as $-$, vice versa
- Assume $C(+|+)$ $C(-|-) = 0$
- Can incorporate costs into growing and pruning tree as well!

## Growing Tree using cost information

- Can use/nor use cost in growing a tree
- GINI without cost: $g(t) = \sum_{i=1} \sum_{j \neq i} p(j|t) p(i|t)$

- GINI with cost: $\sum \sum C(j|i) p(i|t) p(j|t)$  (costs $j$ as $i$)

- Without costs $r(t) = 1 - \max p(j|t)$   (node level)
- With cost    $r(t) = \min \sum C(j|i) p(j|t)$   (node level)
- $R(T)$ (misclass of tree) $= \sum_{|T|} r(t) p(t)$    (tree level)
- Therefore costs do alter pruning

- Difficulty to estimate cost.
- Ratio of costs is important.

## Incorporating Costs into Model

- $Err = p(+)^* FN + p(-)^* FP$
- $Cost = p(+)^* FN(+|-) + p(-)^* FP(-|+)$  $= $ Expected cost (ECost)
- Define our performance line in terms of cost and prior
- Maximum value of ECost: All cases incorrectly classified.
   $= p(+)^* C(+|-) + p(-)^* C(-|+)$
- Norm (ECost) $= $ ECost / MAX COST
- Now plot Nor ECost vs a function of $p(+)$ and cost

## DA : COSTS

- $p(+)$ is also redfied to include misclassificun costs
- Multiply $p(+)$ by $C(+|-)$ and normalize it so axis go from $0 \to 1$

$$PC(+) = \frac{p(+)^* C(+|-)}{p(+)^* C(+|-) + p(-)^* C(-|+)}$$

- For equal misclassifium cost $PC(+) = p(+)$

- $PC(+) = 0$ when $p(+) = 0$ or $C(+|-) = 0$
- $PC(+) = 1$ when $p(-) = 0$ or $C(-|+) = 0$

- Expected Removal (or Norm EGH = $FN^* PC(+) + FP^* PC(-)$

Recall: ratio of number of relevant records removed to the total number of relevant records in the database - measures how well a search system finds what you want

Precision: ratio of the number of relevant records removed to the total number of irrelevant and relevant records removed - how well it weeds out what you don't want.

13/04/16   OA

## MODEL EVALUATION
- Result make sense with preliminary analysis?
- Make sense with background of data?
- Splits and terminal node make sense?
- # of terminal nodes
- Test data
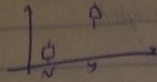- Applies to many different types of models

Test r Training Set:
- Split data into training and test data
- Model built on training data
- Use test set for evaluation
- Typically 80:20 or $\frac{2}{3} / \frac{1}{3}$
- Make sure you have enough case(s) for each target category - use stratification
- Has to be done Randomly
- Test data should reflect future data

Methods:
- Simple plots and summary statistics
- Confusion matrices
- ROC - Receiver Operating curve
- lift chart
- Role of cost and prior probabilities

- Use test data as input
- Calculate probabilities or value

- Trying to evaluate how close predicted probabilities are to be true probabilities

- In a good model should be a clear difference between yes and no outcome ie



- A bad model will have plots that are basically the same → no distinction between groups

## Briers Score
- N cases, 2 groups

$\hat{p}_i$ - predicted probability used to assign to group

$t_i$ - target value 0 or 1

$$= \frac{1}{N} \sum (t_i - \hat{p}_i)^2$$

- like mean square error
- Small when probability estimate is small and $t_i = 0$ and probability estimate is big and $t_i = 1$

## Create Confusion table.
- Pick a threshold / cutoff based on the predicted probabilities, ie. 0.5
- If $p_i < 0.5 \to 0$     $p_i > 0.5 \to 1$
- cutoffs can vary
- Classify into event and non-event

|        |     | Predicted | |
|--------|-----|-----|-----|
|        |     | +   | −   |
| Actual | +   | TP  | FN  |
|        | −   | FP  | TN  |

Want TP and TN big

Sensitivity $= \dfrac{TP}{TP + FN}$   Measure of accuracy for predicting target events (1's)

TPR

Specificity $= \dfrac{TN}{TN + FP}$   Measure of accuracy for predicting non target events (0's)

False Positive Rate $= (1 - \text{True Negative rate})$

% cases correctly classified: $\dfrac{TP + TN}{TP + TN + FP + FN}$  ← opposite is misclassification

Sometimes call accuracy (ACC) or Misclassification rate $= 1 - ACC$

- Could incorporate a cost associated with wrongly classified

- Can create plots such as True positive rate v.s cutoff value
  ↳ At 0 cutoff, all positives are correctly classified
  ↳ At 1 cutoff, all positive are incorrectly classified

- True negative rate v.s cutoff
  ↳ At 0 all classified as positive -
  ↳ At 1 all classified as negative - 1 x rate

- Plot accuracy v.s. cutoff or all graphs together
- Models may have same accuracy but different sensitivity and specificity
- Rare events ⇒ the accuracy calculation will be swamped by larger event number

ROC Curve
- Plots true positives v.s False positives for a selection of cutoff
- TPR vs FPR $= (1 - TNR)$ for a selection of cutoff
  Sensitivity vs $1 - $ Specificity
- Calculate at different cutoff
- Good curve   (0,0) (0,1) and (1,1)

- Look at Area Under the Curve AUC (want as close to 1)
- Bootstrap and create a number of ROC curves and calculate a CI
- 45° line for random model
- Check model parameters, do they make sense?
  - calculate predicted probabilities
  - Draw boxplot    look at RoC curves etc

- Sometimes more important to get a high Sensitivity or Specificity

Model Evaluation - Alternative Approach
- Two outcomes with $N = 3333$     # Yes = 483
- Run model, logistic regression or tree
- Produce predicted values using test set $P_i$   (probability of churning)
- Sort the data into 10 parts - deciles

- best case    is

| Predicted | No | Yes |
|---|---|---|
| | 0 | N |
| | 0 | $\checkmark$ |
| | $\vdots$ | $\vdots$ |
| | N | 0 |
| | N | 0 |
| | N | 0 |

(all yes values at top, all no values at bottom)

- Cumulate approach - cumulate value in cell downward and can then calculate % captured

- $N$ = # in sample

$N_r$ = total no. of responses where = 1 (defined as an event)

$N_d$ in each decile $N/10 = N_d$

Non cumulative % captured response in decile i : $\dfrac{N_{ri}}{N_r} \times 100$

cumulate rank $S$ $\Sigma nk$ as percent

Cumulative: % Captured response ≤ decile i    $\frac{\Sigma\,Nej}{Ne} \times 100$

% Captured response: column %'s for yes) "Recall" in R

% Respondents - Row % for yes or "Precision" in R

Lift -    $\dfrac{\% \text{ Captured Responses}}{\% \text{ Random capture}}$    Recall


devel

Lift for decile i = $\dfrac{\% \text{ Captured resp}}{10}$    for non cumulative    — expect 10% in each decile here to

Lift for decile i = $\dfrac{\% \text{ Captured resp}}{\sum\limits_{j=1}^{i} 10}$    For cumulative approach
i.e 30 for 3rd decile

- Want lift high in top deciles then low    - how much better than chance model
- % response/precision    - want it to decrease to zero quick.

Possible Scenarios
- Can plot for model results, random results and best possible results.
- Performance chart in Rattle:
   Adjustment = cumulative % captured response / by % decile
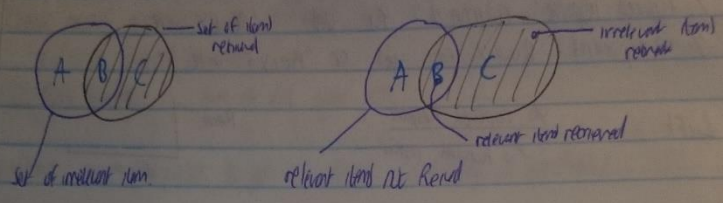   Strike rate = cumulative % Response/by % decile
   Black line = Cumulative random % captured response

- Confusion matrix
- Misclassification rate
- ROC
- Lift           - Recall
- Precision      - Cut off
- Accuracy

Which tool? - Dependant on study
- What you want to optimize



— Set of items retrieved

— irrelevant items returned

relevant item retrieved

Set of irrelevant item

relevant items not Retrieved

$Recall = \dfrac{B}{A+B}$

A: number of relevant records not retrieved

B: number of relevant records retrieved

$Precision = \dfrac{B}{B+C}$

B: # of relevant records retrieved

C: # of irrelevant records retrieved

RECALL - Ratio of number of relevant records retrieved to the number of total # of relevant records in the database

PRECISION: ratio of # of relevant records retrieved to the total number of irrelevant and relevant records retrieved

| actual | | + | − |
|---|---|---|---|
| | Relevant + | TP | FN |
| | Irrelevant − | FP | TN |

$Precision = \dfrac{TP}{TP + FP}$

$Recall = \dfrac{TP}{TP + FN}$

$TP + FP = $ # of items retrieved

AS recall ↑ precision ↓     or    recall ↓ Precision ↑

- Categorizing items as relevant or irrelevant
- Determining the # of relevant events in the database
- Recall measures how well a search system finds what you want, precision measures how well it weeds out what you don't want

14/16  DA

## CARAT PACKAGE OUTPUT

Accuracy: predicted v actual (reference) values

Accuracy: $\frac{TP + TN}{TP + TN + FP + FN}$

95% CI: CI for accuracy   $CI = p \pm 1.96 \times \sqrt{\frac{p(1-p)}{n}}$ ← t-value

or bootstrap and take percentile

– wide because sample size is small

No information Rate: How you would do without model, if I predict everything as ayes.

P-Value: How much better you did with model (accuracy) than without it (NIR)

Kappa: First row $\frac{actual\ yes}{total}$  agreement controlling for chance

McNemar's Test P Value: paired chi-squared test, suggests no difference in example (probability of 1)

looking at the 6 and 5 in example table

Sensitivity: Measures accuracy of predicting "yes" correctly   $\frac{TP}{TP + FN}$

Specificity: Measures accuracy of predicting "no" correctly   $\frac{TN}{TN + FP}$

$H_0: Acc \leq NIR$    v   $H_1: Acc > NIR$, evidence against $H_0$, accept alternative

Pos Pred Values: Percentage of predicted yes over all yes

Neg Pred Value: Percentage of predicted no over all no's

Prevalence: Prevalence of event  27/51  total left column / Total

Detection rate = $\frac{TP}{TP + TN + FN + FP}$

Detection Prevalence = $\frac{TP + FP}{TP + FP + TN + FN}$
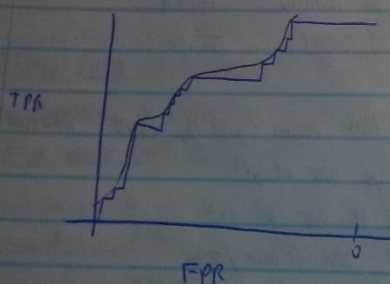
Balanced Accuracy = $\frac{(sensitivity + specificity)}{2}$

Example of McNemar

|   | A correct | misclass B |
|---|---|---|
| B miscl(ass) | $n_{00}$ | $n_{10}$ |
| N | $n_{01}$ | $n_{11}$ |

- Information only in off diagonal elements $n_{01}$ $n_{10}$
- No difference between the classifiers would expect $\frac{n_{01}+n_{10}}{2}$ in each cell

$$\chi_1^2 = \frac{(O-E)^2}{E} = \frac{(n_{01} - \frac{n_{01}+n_{10}}{2})^2}{(\frac{n_{01}+n_{10}}{2})} + \frac{(n_{02} - (\frac{n_{01}+n_{02}}{2}))^2}{(\frac{n_{01}+n_{02}}{2})} \quad d.f = 1$$

Reduces to : $\frac{[|n_{01} - n_{10}|-1]^2}{n_{01} + n_{10}}$  $\chi^2$ with d.f = 1



TPR

FPR

ROC with (concave) and ROC (use convex hull)

# ROC and PRIORS

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

$$Acc = p(+) \cdot TPR + p(-) \cdot (1-FPR)$$

$p(+)$ = proportion of + in population written as $\pi$, up to now
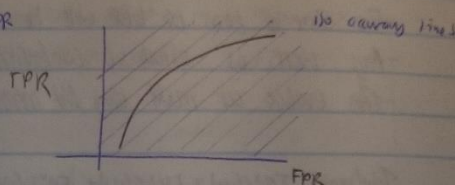
$p(-)$  "    of -    "    "

$p(+) + p(-) = 1$

13/04/16   DA

Can draw lines of ISO accuracy on the ROC curve

$$TPR = \frac{Acc - p(-)}{p(+)} + \frac{p(-)}{p(+)} * FPR$$



## Convex Hull

- All points on convex hull dominate
- looking for point on convex hull where slope of tangent = $p(-)/p(+)$
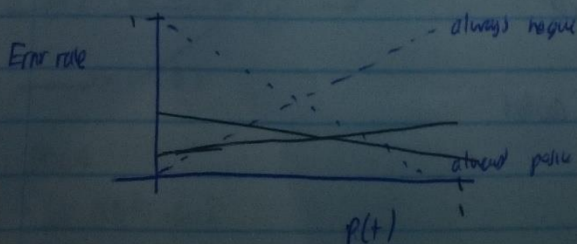- Operating condition $p(-)/p(+)$

## A new Space

- Plot error vs $p(+)$
- $Err = 1 - Acc = 1 - [p(+) * TPR + p(-)(1 - FPR)]$
- $Err = p(+) FNR + p(-) FPR$
  - $= (FN - FP)' p(+) + FP$
- New Space: Y Axis: Err      X-Axis: $p(+)$

- Point in ROC curve are going to be mapped onto lines in "error/$p(+)$" Space
- Allow us to see clearly how error rate varies according to $p(+)$
- How the model will behave under various conditions.



Take minimum convex hull!

- look at the envelop
- look at values for each $p(t)$
- For some a item the bet will be always positive or always negative
- May want to attach a cost/profit/loss to each misclassification
- Then evaluate our model using this information


- Total cost: $TP^* C(+|+) + FN^* C(+|-) + FP^* C(-|+) + TN^* C(-|-)$
- Objective to minimize cost
- Can incorporate cost into growing and pruning a tree also


Growing Trees with Cost
- GINI with cost: $g(t) = \sum_{i=1}^{c} \sum_{j \neq i}^{c} p(j|t) p(i|t)$

$$\sum_{j=1}^{c} \sum_{j \neq i} C(j|i) p(i|t) p(j|t) \quad \text{with cost } j \text{ vs } i$$

- Without cost misclassification $= r(t) = 1 - \max P_j (j|t)$  → node level
- With cost $r(t) = \min \sum_i C(j|i) p(j|t)$  → node level

  $R(T) = \sum_{|T|} r(t) p(t)$  → tree level
  ∴ costs after pruning regrow

Difficulties
- Choosing model of cost
- Difficulty to estimate cost
- What is important is ratio of cost
- Objective changed
- Many more complicated models oscillate