09/04/16

# DATA ANALYTICS NOTES

## What is data mining?

- Exploration and analysis by automatic or semi-automatic means of large quantities of data in order to discover meaningful patterns and rules
- "A tool for extracting the jewel of truth from the slurry of data"
- "All models are wrong but some are useful" - George Box
- "We live in a data rich world, most of us stand on the shore of a vast sea of available data, suited up with the latest diving gear and equipped with the slickest tool and gadgets but with hardly a clue what to do".
- Need to sort data and build predictive models

## Steps in problem solving Process

- Recognise the problem or question
- Review previous findings
- Model the Solution
- Collect the data
- Analyse the data
- Present and act on results

## Time

- Time horizon for prediction
- Time window of relevant behaviour
- Time base of the population

## Questions to ask:

- Source of your data?
- How well does Sample data represent population?

- Why did you decide on that particular approach?
- What alternatives did you consider?
- How likely do the independent variables cause to dependent variable?

Stakeholders
- Who has stake in areas of your project?
- Have they been briefed on problem and outline for solutions?
- Resources and charges?
- Review and revisit model
"Data mining on it's own will not provide the best models; these will be created by the interplay between the knowledge extracted from the data and the experience of specialist staff"
- Only an aid to decision making, not the decision itself
- Neither sophisticated software nor statistical techniques can overcome the inherent limitations of the raw data that goes into them
- GARBAGE IN ⇒ GARBAGE OUT

- Do we have enough data both cases and variables?
- Is it legal and available for use?
- How easy is it to collect and process?

Quality of Data
Accurate? complete? current? consistent?
Categorical, ordinal, quantatative

First Steps
Explore data graphically and produce descriptive stats
One variable - dot plot etc
Two vars - Scatterplot

9/04/16  DA

Many var - Scatterplot matrix
One quan, One cat :  Piechart, box plot
Two cat :  Frequency distribution, tables

## PROBLEMS WITH DATA

Outliers
- How are we sure it is an outlier?
- do nothing, ignore the case of variable, replace the values, transform the data

Missing Data
- May reduce number of cases included  (reduced data size)
- Remove rows of case    or else variable?
- Usually coded as NA

- Look at % of each var that is missing
- % of each case that is missing
- if one missing, likely to miss another? - look for patterns
- Visualisation of Missing Values
- Assume missing data is at random

What can we do?

- nothing
- list wise deletion - only included if it has data for all cases
- Pair wise deletion - data included if the case contains certain variables like AC or $D_iF$
- Omit variables or cases with high % (10% +) of missing data
- Weighting for non response
- Imputation - Substitution of data.
- Create new variable
- Models are unique in dealing with missing data

Imputation
- Method for substitution of missing data: - regression model
    - copy from left
    - bootstrap
    - mean → sd, median, shape of distribution.
    - Midrange
    - Distribution-based replacement calculated on percentiles of variable's dist.
    - Hot deck - divide sample into groups and select value at random within group
    - Most frequently occurring variable
- Disadvantage: - Alters relationship between variables
    - May increase biases in survey estimates.
    - Researchers may falsely treat data as a complete dataset
    - Imputed values should be flagged

Multiple imputation - Create M datasets with imputed values (complete datasets)
    - Combine results which should reflect missing data uncertainty.

09/04/16   DA   Classification & REGRESSION TREES

- Supervised learning with a target variable
- Categorical for classification tree, continuous for regression tree
- Need observations on other variables
- Goal: predict or classify an outcome
    ↳ construct a rule to apply in the future
    ↳ To see what variables are important/related to target variable

Heart attack example
- 215 cases, 37 died, 100 variables screened

Classification tree
- Upside down tree, root node - all patients
- Two branches defined by a question - Terminal nodes ⇒ classification

Outcome Variable definition
- Situation dependent or blindingly obvious
- Consensus, prescribed, empirical, time frame

- coding: focus on two category outcome variable Yes/No dead/Alive

Splits for continuous Variables
- At each node independently:
    - Examines all possible splits
    - looks at each value for each variable
    - Is it a good split? Which is best split?
    - Importance of screening
    - Outliers? → don't affect outcome

Splits for categorical data
- Consider all possible splits    A  B,C    AB,C  etc    n-values $\Rightarrow 2^{n-1}$ splits

Evaluating Splits
- will have a high class %. i.e.   5 dead 300 alive  etc
- Impurity function
  ↳ node which contains only one class is perfectly pure
  ↳ Equal proportion of each class is least pure
  ↑ need a measure to distinguish these 2 cases

Prior probabilities
- Chance that a case will be presented to a tree
- Expert Knowledge
- Data
- Assume prior
- Rare cases

$N_J$ - Number of in class J overall
$N(t)$ - total number of cases in node t.
$N_J(t)$  # of class J cases in node t.
Proportion of class J cases in node t = $N_j(t) / N_J$
$\pi(j)$ = prior probabilities.   $\sum_j \pi(j) = 1$

Probability that a case will be both in class j and fall into node t
$$p(j,t) = (\pi_j \, N_j(t)) / (N_j)$$

P the of a case falling into class j given that a case is in node t
$$p(j|t) = \frac{p(j,t)}{p(t)}$$

7/04/16   DA   CLASSIFICATION  & R TREES

P that any case falls into Node $t$ = $p(t) = \sum_j p(j,t)$

Common Impurity functions

$c$ = # categories in target variable

Entropy = $-\sum_{j=1}^{c} p(j|t) \log_2(p(j|t))$

Gini = $\sum_{i=1}^{c} \sum_{j=1 \atop j \neq i}^{c} p(j|t) p(i|t)$  = $1 - \sum_{j=1}^{c} p^2(j|t)$

When $c=2$   Gini = $2^* p(1|t) p(2|t)$  | Smaller the better → means probability one big, one small

How to use this to choose split

- look at impurity of parent node $t$  : $i(t)$

- Each split at any node gives us 2 children $t_L$ and $t_R$

- We measure the impurity of the two children

  $i(t) - P_L^* i(t_L) - P_R^* i(t_R)$

  $P_L$ = probability going left    $P_R$ = Prob going right

Choosing nodes

- $\Delta(t, s) = i(t) - P_L^* i(t_L) - P_R^* i(t_R)$

- Smaller values for $i(t_L)$ and $i(t_R)$ are better (higher classification rate)

  $P_L + P_R = 1$

- Could be 5000+ splits at each node → evaluate for each one:

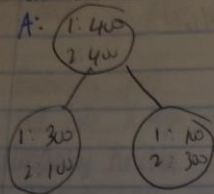- $i(t)$ will be the same for all

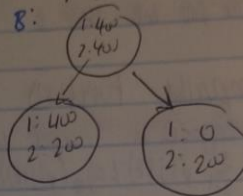- choose split with largest $\Delta(t, s)$

- R calls this improvemen.

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

Use GINI to calculate decrease in impurity for each split

A:



B:

Improvement: $\Delta(t,s) = i(t) - P_L i(t_L) - P_R i(t_R)$

$i(t) = 2 * P(1|t) P(2|t)$

$= 2 * \dfrac{P(1,t)}{P(t)} \dfrac{P(2,t)}{P(t)}$

A. $i(t) = 2 * 0.5 * 0.5 = 0.5$

$P_L i(t_L) = \dfrac{400}{800} * 2 * 0.75 * 0.25 = 0.19$

$P_R i(t_R) = \dfrac{400}{800} * 2 * 0.25 * 0.75 = 0.19$

$\Delta(t,s) = 0.5 - 0.19 - 0.19 = 0.12$

B. $i(t) = 2 * 0.5 * 0.5 = 0.5$

$P_L i(t_L) = \dfrac{6}{8} * 2 * \dfrac{2}{3} * \dfrac{1}{3} = 0.33$

$P_R i(t_R) = \dfrac{2}{8} * 2 * 0 * \dfrac{1}{1} = 0$

$\Delta(t,s) = 0.5 - 0.33 - 0 = 0.17$

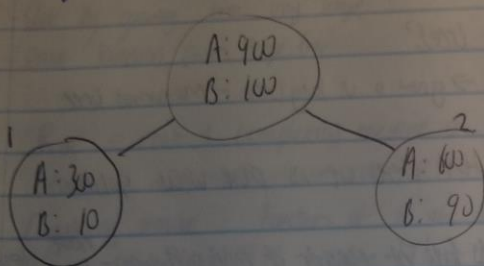A has lower impurity/improvement value for the split

Class Assignment Rule

- Assign a class to every terminal node
- Calculate $p(i|t)$ for each class $i$
- Let $j^*$ be max of $p(i|t)$
- Assign the node to class $j^*$
- Can assign a class or output a probability of being assigned to the node.

Priors
- usually use equal prior   0.5, 0.5



A: 900
B: 100

1
A: 30
B: 10

2
A: 600
B: 90

$p(i|t)$, $p(i,t)$, $p(t)$   using equal prior and data as prior

Equal prior

$$p(A,1) = \frac{\pi_i \, N_i(t)}{N_i}$$

$$0.5 \times \frac{300}{900} = .17$$   $$P(A,2) = 0.5 \times \frac{600}{900} = .33$$

$$P(B,1) = 0.5 \times \frac{10}{100} = 0.05$$   $$P(B,2) = 0.5 \times \frac{90}{100} = 0.45$$

$$P(1) = .17 + 0.05 = .22$$   $$P(2) = 0.33 + 0.45 = .78$$

$$P(A|1) = \frac{P(A,1)}{P(1)} \quad \frac{.17}{.22} = .77 \quad P(B|1) = \frac{.005}{.22} = .42$$

$$P(A|2) = P(B|2) \cdot \frac{.45}{.71} = .58 \quad P(B|1) = \frac{.05}{.22} = .23$$

Assign all cases in node 1 to A, all in node 2 to B

Data as priors

$$\pi_A = \frac{900}{1000} = .9 \qquad \pi_B = \frac{100}{1000} = .1$$

$$P(A,1) = \frac{\pi_i \, N_i(t)}{N_i} = 0.9 \times \frac{300}{900} = .3 \qquad P(A,2) = 0.9 \frac{600}{900} = .6$$

$$P(B,1) = .1 \left(\frac{10}{100}\right) = .01 \qquad P(B,2) = .1 \left(\frac{90}{100}\right) = 0.09$$

$$P(1) = .3 + .01 = .310 \qquad P(2) = .6 + .09 = .69$$

$P(A|1) = \frac{3}{31} = .97$                    $P(A|2) = \frac{.06}{.69} = .87$

$P(B|1) = \frac{.9}{.31} = 0.03$              $P(B|2) = \frac{.09}{.69} = .13$

Assign both nodes to A

When to Stop growing trees?

Bottom up – classical CART → grow a v big tree – maximal tree
  – Prune branches

Top down – Stop growing when there are no more ideal splits

Classical CART approach – need to look at concept of misclassification – node level / tree level

Confusion matrix → misclassification rate

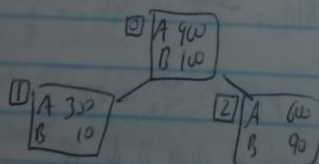|  | predicted | |
| --- | --- | --- |
|  | Y | N |
| observed  Y |  |  |
| N |  |  |

## Node Misclassification

Define $r(t)$ p of misclassification at node t $= 1 - \max(p(i|t))$

$R(t) = r(t) * p(t)$

For the tree we can write:

$R(\tilde{T}) = \sum_{\tilde{T}} r(t) p(t)$  where $\tilde{T}$ is all terminal nodes

Using tree:                                         and data w/ priors

```
          0 | A 900 |
            | B 100 |
   1 | A 300 |        2 | A  60 |
     | B  10 |          | B  90 |
```

Node 0: $r(t) = 1 - 0.9 = 0.1$

$R(t_0) = r(t_0) * p(t_0) = 0.1 * 1 = 0.1$

N1: $r(t) = 1 - .97 = .03$

$R(t_1) = r(t_1) * p(t_1) = 0.3 * \frac{310}{1000} = 0.0093$

N(2): $r(t) = 1 - \frac{600}{690} = .13$

$R(t_2) = r(t_2) * p(t_2) = .13 * \frac{690}{1000} = 0.0897$

$R(T_t) = 0.0093 + 0.0897 = 0.99$

09/04/16   Classification   +   R Tree

## Cost   Complexity   Pruning

-2 Stage     1: Develop a model sequence for evaluation

2: Choose final tree

- Start by growing tree very large
- Prune branch(s) from large tree

$R_\alpha$ = cost + complexity measure of the Tree T

Define cost as misclassification rate R(T)

Complexity measure : function of # of terminal nodes

$R_\alpha = R(T) + \alpha^* |T|$   where $|T|$ = # of terminal nodes.

$\alpha$ = penalty placed on complexity

- For a single node : t: $R(t)_\alpha + \alpha$
- For a subtree $t_t$   $R(T_t)_\alpha = R(T_t) + \alpha |T_t|$
- When $\alpha$ increases both $R(t)_\alpha$ and $R(T_t)_\alpha$ increase   but
  $R(T_t)_\alpha$ increases faster.
- Value of $R(t)_\alpha = R(T_t)_\alpha$
- Price paid for complexity - adding on to tree

$$\alpha = \frac{R(t) - R(T_t)}{|T_t| - 1}$$

- Bigger value for $\alpha$ implies it is a better branch
- Smaller → weaker

## Some example tree

For 2 node tree    $R(T_t)_\alpha = R(T_t) + 2\alpha$

$$\alpha = \frac{R(t_0) - R(T_T)}{|T_t| - 1} = \frac{.1 - 0.099}{2 - 1} = .001$$

- calculate α at each node
- Prune tree with lowest α - weakest link
- Recalculate again

In R   α is called complexity parameter
(calculated same way but is divided by (te r(s) - misclassification of the root
- keep deleting branches and re-calculating

misclassified

# nodes

- Can set a stopping value for CP
- Choose tree with min misclassification or cross validation

- Use training and test data ⇒ see where they cross over ⇒ choose the
  or value line.
- Gives a rate and S.E when you use multiple test sets

11/04/16    DA

CART vs Logistic Regression
- CART excels in detection of local structure
- Each half of tree is analysed separately
- Discovery of patterns becomes progressively more local
- Information from different nodes is not pooled or combined
- The fit at one node is never adjusted to take into account fit at another node
- Good at inference
- Automatic separation of relevant from irrelevant predictors - only used are
  which define "good" splits  - Screens data
- Does not require transform like log etc
- Automatic interaction detection        - variables don't need to be deleted in advance
- Impervious to outliers.
- Has methods for dealing with missing values
- Requires only moderate supervision by analyst.
- First time model is often as good as a neural net by an expert
- non parametric - doesn't require specification of any functional form

Disadvantage of CART
- May have unstable decision tree - insignificant modification of learning
  sample such as eliminating observations could lead to radical
  changes in the decision tree
- CART splits only by one variable - all splits are perpendicular to axis. ie.
  when splitting all data into boxes, all boxes are rectangles - if data
  has more complex structure, CART may not catch correct structure of data

## Logistic Regression

- Provides a smooth continuous predicted probability of class membership
- Effective capture of global features of the data
- Main effect model reflect show probability regard to predictor x over entire range of Y
- Some flexibility allowed with transformation, polynomial and interaction
- Provides standard errors of coefficients
- Independent variables don't have to be normally distributed
- Handles non-linear effect
- You derive model by selection
- No homogeneity of variance assumption

### Disadvantage

- Need larger sample of data
- have to identify correct independent variables
- limited out variables
- Independent observation(s) required
- Over fitting model
- Numeric approximation - does not always converge towards an optimal solution
- Does not handle missing values of continuous variable
- sensitive to extreme data

- Requires expert

- CART bad at detecting linear structure, recognizes it but can't represent it effectively
- With many variables, linear structure may not be obvious from CART analysis
- Can produce a very large tree in an attempt to represent very simple relationships
- Logistic Regression - good for linear relationships
- LR many non linear structures can still be reasonably approximated with a model
- Even incorrectly specified logistic regression can perform well

Combine CART and Logistic Regression ?
- No information left in terminal nodes to support further analysis
- In a well developed CART tree, no other model should be supportable in the nodes
- Run a shallow tree
- Assign every case a terminal node
- Terminal node assignment reparameterised by categorical variable with as many levels as terminal nodes
- Feed this categorical variable in the form of terminal node dummies to LR model

CART only with $j$ terminal nodes

$$y = \beta_0 + \beta_1 N_1 + \ldots + \beta_{j-1} N_{j-1}$$

$N, j-1$ dummy variables

CART - Logistic Regression Hybrid model

$$y = \beta_0 + \beta_1 N_1 + \ldots + \beta_{j-1} N_{j-1} + \alpha_1 Z_1 + \ldots \alpha_r Z_r$$

$Z_i$ extra variables

= CART node Dummies + Hybrid variables

## Surrogates

- A variable with possibly equivalent information
- A surrogate is a splitter that splits in a fashion similar to the primary splitter
- Reveals structure of the info in the variable at a particular node in the tree
- If the primary is expensive or difficult to gather - use surrogate instead
- Use surrogate split if data is missing
- Cases with data for both split
- For any node $t$, primary splitter $s$ sends $t_L$ could to left and $t_R$ to right
- For any other split $s^*$ of the node $t$ into $t_L^*$ and $t_R^*$
- $N_j(LL)$ number of cases in $t$ that both $s$ and $s^*$ send left for class $j$.
- $N_j(RR)$ # cases in $t$ that both $s$ and $s^*$ send right for class $j$.

- Surrogate are node dependant - calculated at the local level
- Surrogate versus competitors
- Useful for examining what node is trying to do
- Can choose how many you calculate at each nd.

## Variable Importance

- Assess the relative importance of the variables in a tree
- Measured by impurity improvement
- looks at primary splitter for each node and all the surrogate splits listed on rpart for every node
- Can control the number of surrogates.
- Calculated over tree
- Calculated as 0 %

## Calculation of $p(s, s^*)$

For this example $j = 2$

←from confusion matrix

$NN_{1L} = 500$     $NN_{2L} = 100$     $NN_{.LL} = 600$

$NN^*_{1L} = 450$     $NN^*_{2L} = 150$     $NN_2(LL) = 100$

$NN_{1R} = 100$     $NN_{2R} = 300$     $NN_1(RR) = 50$

$NN^*_{1R} = 150$     $NN^*_{2R} = 250$     $NN_2(RR) = 250$

## Calculating Similarity

- Split $S$ and Surrogate split $S^*_j$, for node $t$.
- $j$: no. of classes and $\pi_i$ is prior probability of class $j$
- Probability $S$ correctly predicts $S^* = p(s, s^*)$

$$p(s, s^*) = P_{RR}(s, s^*) + P_{LL}(s, s^*)$$

$$P_{LL}(s, s^*) = \frac{p(t_L, t^*_L)}{p(t)} = \sum_j \frac{\pi_j N_j(LL)}{N_j} / p(t)$$

e.g for $P_{RR}(s, s^*)$

11/04/16   DA

$$p(t) = \sum_j p(j,t) = \sum_j \frac{\pi_i N_j(t)}{N_j}$$

Using data of priors we can show   $p(t) = \frac{\sum_j N_j(t)}{\sum_j N_j}$

$\pi_1 = \frac{N_1}{N_1 + N_2} = 0.6$   $\pi_2 = 0.4$

- In this case, since we are at the root node, $p(t) = 1$

- $P_{LL}(S^x, S) = 0.6 \left(\frac{400}{600}\right) + 0.4 \left(\frac{100}{400}\right) = 0.5$

- $P_{RR}(S^x, S) = 0.6\left(\frac{20}{600}\right) + 0.4\left(\frac{20}{400}\right) = 0.3$

$p(S, S^x) = P_{LL}(S, S^x) + P_{RR}(S, S^x) = 0.5 + 0.3 = 0.8$

Error $= 1 - 0.8 = 0.2$

11/04/16     DA

Derived Variables
- Define new variables to make data more useful and informative
- Creative part of the process
- Improve quality of data
- Well chosen derived variables enhance the ability of models to be understood and interpreted
- Change over time

- Allows analysts incorporate human insights and background knowledge into modelling process
- Important skill to come up with "right" variables - what works in one model may not work in another
- Depends on country, setting / context
- Replace categorical variable with quantitive when many categories
- Use prior values of target variable or censed data
- Be careful with time

Phase
- What information to use? # apps? # phones call? more? model?
- Sometimes variables are too hard to use

Single Variable Transform
- Standardised numeric variables
- Centered variables - can increase interpretability
- Rescale ⇒ allow comparison between variables
- Turn number values into percentiles works with any distribution
- Useful when you are interested in relative position than absolute value
- Turn counts into rates

- Replace categorical variables with numeric - do not do so arbitrarily
- Create indicator variables, work well with few categories
- Capture important information
- Place - latitude and longitude
- Use previous data from target data - make sure it is going to be possible in the future
- Bin numerical data - equal intervals or quantiles
      ↳ may not be good with some techniques

Combination of Variables
- BMI
- Price earnings ratio
- Highly correlated variables - draw scatterplot before correlation calculation
- Nearly synonymous variables
- When two variables are equal, most of the time few places where the disagree may be informative
- Depends on DM technique used.      Multicolinery bad for regression

Correlated Variables
- Get rid of one - the one with less variance
- Talk to people with background knowledge as they may have more info or will understand the make!
- Try to derive a variable with high variance which is independent of its associate variables
- Ratio of the variables

Netflix example
- Ratings on scale 1-5, 480k users, 17770 movies
- Individual ratings
- total number of ratings
- number in first month

4/16   DA
   Derived Variables
   - Proportion in first month
   - Ratings per month
   - Proportion of 1 ratings etc
   - Average Rating
   - Average Rating for recent month
   - Ratio of 1 to 5 rating.
   - Date from release to rating
   - Comparison of rated to population
   - Population average - rule average per move

   - Reduce variables : PCA, FA, cluster membership

   KAPPA   STATISTIC
   - Cohen's  Kappa Statistic
   - Agreement between raters originally
   - Look at agreement taking into account to accuracy that would be
   generated by chance
      $$Kappa = \frac{Cohered\ accuracy\ -\ accuracy\ expected\ by\ chance}{(1 - accuracy\ expected\ by\ chance)} \left(\frac{O-E}{1-E}\right)$$

   - Values range from -1 to 1
   - 0 value means no agreement
   - Measures relative improvement over random predictor.

Predicted

|        |     | Yes | No  |     |
|--------|-----|-----|-----|-----|
| Actual | Yes | 100 | 100 | 200 |
|        | No  | 300 | 500 | 800 |
|        |     | 400 | 600 |     |

$$O = \frac{100 + 500}{1000} = 0.6 \quad \text{observed accuracy}$$

$$E = \left[\frac{400}{1000} \times \frac{200}{1000}\right] + \left[\frac{600}{1000} \times \frac{800}{1000}\right] = 0.56$$
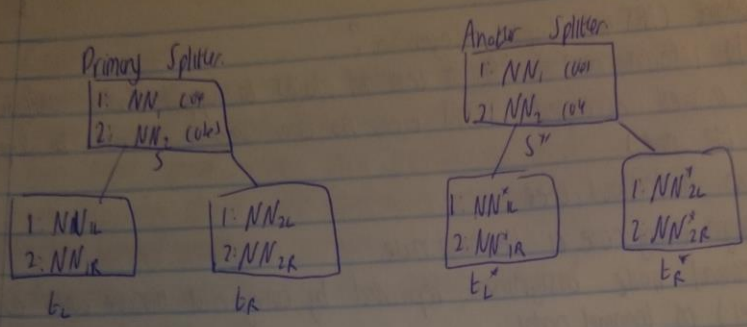
$$Kappa = \frac{0.6 - 0.56}{1 - 0.56} = 0.09$$

0 - 0.2 insufficient
0.21 - 0.4 schwach          even with low amount of data.
.41 - .6 suffia
.61 - 8 good
.81 - 1 exceln

Primary Splitter.

```
┌─────────────┐
│ 1: NN₁ (60)  │
│ 2: NN₂ (40)  │
└──────┬──────┘
       S
```

$$NN_1 \ (60) \qquad NN_2 \ (40)$$

```
┌──────────┐        ┌──────────┐
│ 1: NN₁ₗ  │        │ 1: NN₂ₗ  │
│ 2: NN₁ᵣ  │        │ 2: NN₂ᵣ  │
└──────────┘        └──────────┘
    tₗ                  tᵣ
```

Another Splitter.

```
┌─────────────┐
│ 1: NN₁ (60)  │
│ 2: NN₂ (40)  │
└──────┬──────┘
       S*
```

```
┌────────────┐      ┌────────────┐
│ 1: NN*₁ₗ   │      │ 1: NN*₂ₗ   │
│ 2: NN*₁ᵣ   │      │ 2: NN*₂ᵣ   │
└────────────┘      └────────────┘
    t*ₗ                 t*ᵣ
```
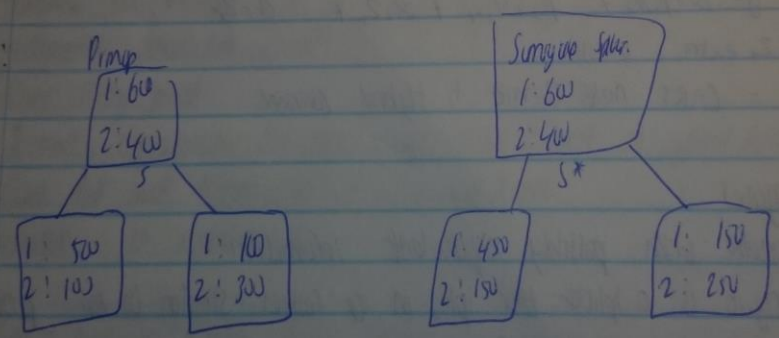
Perfect Splitter.

$$NN_{1L} = NN^*_{1L} \qquad NN_{2L} = NN^*_{2L} \qquad NN_{1R} = NN^*_{1R} \qquad \text{or} \quad NN_{2R} = NN^*_{2R}$$

- Same code - just not similar name.
- Need to calculate probability that both splitters S and S* send cases in the same direction
- Call this p(S, S*)

EXAMPLE:

Primary

```
┌──────┐
│ 1: 60 │
│ 2: 40 │
└───┬──┘
    S
```

```
┌──────┐      ┌──────┐
│ 1: 500│      │ 1: 100│
│ 2: 100│      │ 2: 300│
└──────┘      └──────┘
```

Surrogate Splitter.

```
┌──────┐
│ 1: 60 │
│ 2: 40 │
└───┬──┘
    S*
```

```
┌──────┐      ┌──────┐
│ 1: 450│      │ 1: 150│
│ 2: 150│      │ 2: 250│
└──────┘      └──────┘
```

Class 1

Surrog

|       | Primary L | Primary R |     |
|-------|-----------|-----------|-----|
| L     | 400       | 50        | 450 |
| R     | 100       | 50        | 150 |
|       | 500       | 100       |     |

Class 2

Surr

|       | Primary L | Primary R |     |
|-------|-----------|-----------|-----|
| L     | 100       | 50        | 150 |
| R     | 0         | 250       | 250 |
|       | 100       | 300       |     |

How good a predictor is $S^*$?

- How well can another split $S^*$ mimic the primary splitter $S$?
- How high is $p(S, S^*)$?
- Comparison or default rule
- For node $t$ suppose that $S$ sends cases with prob $P_L$ and $P_R$ right (reluse)
- News case $\Rightarrow$ predict $t_L$ if $p_i = \max(p_L, p_R)$ else $t_r$
- $r(t)$ (misclassification rate) $= \min(p_L, p_R)$

Same example
- Using the data for our priors we can compute these directly from the tree

$P_L = 0.6$ , $P_R = 0.4$

- Default rule $\Rightarrow$ send everything left, error $= 1 - 0.6 = 0.4$
- Compare this to surrogate rule

Surrogate rule:
- $p(S, S^*) = 0.8$ $\Rightarrow$ called agreement in R
- /error rate $= 0.2$

$$\frac{\overset{class\ 1}{LL + RR} \qquad \overset{class\ 2}{+LL + RR}}{n\ (1000)\ hed}$$

$\hookrightarrow$ = the tube var :

Association:
$$= \frac{\min(p_L, p_R) - (1 - p(S^*, S))}{\min(p_L, p_r)}$$

$$= \frac{\text{default mismatch} - \text{surrogate mismatch}}{\text{default mismatch}}$$

$P_L = 0.6$  $P_r = 0.4$   $\Rightarrow \dfrac{.4(1 - 0.8)}{.4} = 0.5$   50% reduction

Relative reduction in error driven by using $S'$ to predict $S$ instead of $\max(p_L, p_R)$ called adv in R