20/04/16  DA - RuleFit

- Developed by Jerry Friedman
- Combines trees and regression
- Used for descriptive or predictive purpose

- Turn a tree into a series of rules i.e. if (age < 70): left else right
- A series of indicator variables
- Geez call them $r_m(x)$
- Each terminal node defines a region
- Regression of terminal nodes.



$$F(x) = c_0 + \sum_{m=1} c_m r_m(x) + \sum_j b_j x_j \longrightarrow \text{ordinary variables in dataset as linear combination}$$

no upper limit → choose     value of $c_m$'s      add to capture linear effect

- $r_m(x)$ output from trees - indicator variables
- Need to determine $c_i$ and $b_j$'s
- Build many trees of a certain depth
- Use these to create new indicator variables
- Use intermediate nodes also

- $P_m$ are the split definitions for rule $r_m(x)$
- Can include rules for all non-terminal nodes aswell
- Counting all nodes, a J-terminal node tree generates $2 \cdot (J-1)$ rules

Combine the linear part and rules part
- Weight each of the terms
- Use some penalty function to reduce number of terms
- Lasso function
- Can use general loss function
- Post processing phase
- Could add non linear terms like $x^2$ or $x^3$

- Can be a large number of terms in relation function
- May type of penalty can be used
  - Sum of absolute value of coefficients - lasso
  - Sum of square of coefficients - ridge regression
  - Or both - elastic net
    - Brug grey

- $p \leq n$ number of continuous variables to be included as linear terms
- because of authors
- M is the number of trees
- $K = \sum_{m=1}^{M} 2 \times (J_m - 1)$ total number of trees
- $J_m$ is number of terminal nodes for tree$_m$

tree Size
- Important to consider
- 2 terminal nodes only consider one factor → no interactions
- To capture interaction need larger size trees
- For 2-way interaction → need 4 terminal nodes
- Use trees of varying size

- Number of terminal nodes is a random variable
- $t_m = 2 + fl(\gamma)$
- Whe $\gamma$ is drawn from an exponential dist with $Pr(\gamma) = \dfrac{\exp(-\gamma)(\bar{I}-2)}{(\bar{I}-2)}$
- $fl(\gamma)$ is the largest integer less than or equal to $\gamma$
- $(\bar{I}-2)$ is the average number of terminal nodes for trees in ensemble $(\bar{I} \geq 2)$

- $\bar{I} = 2$ the entire ensemble will have just 2 terminal nodes
- For $\bar{I} = 3$    $t_m = 2 + fl(\gamma)$
- When $\gamma$ drawn from exponential distribution with probability
    $Pr(\gamma) = e^{((-\gamma)(\bar{I}-2))}/(\bar{I}-2) = e^{-\gamma}$
- Exponential Distribution with mean of 1

### Linear Terms
- Windsorize the linear terms
- Replace the top and lower percentiles with the next higher or lower value
- Typical value for percentiles is 2.5%
- Can reduce outliers
  i.e. assign the top 2.5% of observed value all to be equal to the 97.5th percentile value (similar for the lower 2.5%)

### Spurious Interaction
- Outliers may cause problems
- Many Spurious interactions may occur especially if variables are highly correlated
- Have an incentive for fewer variables entering path
- Chose Split with maximum improvement $z_i$ splitting on $x_i$
- Adjust this to $k_i z_i$ where $k_i = 1$ if $x_i$ had not been used in this branch before
- Otherwise $k_i = k (z1)$
- Discourage highly correlated variables from appearing in the same rule

Spurious Relationship - maths relationship in which two events or variables have no direct causal connection, yet it may be wrongly identified that they do

When A is present, B is observed (A causes B)
When B is present, A is observed (B causes A)
  OR
When C is present, Both A and B are observed (C causes both A and B)

In the last case there is a spurious relationship between A and B in a regression model where A is regressed on B but C is actually the true cause of A.

## Importance

- Importance of any predictor in a linear model is the absolute value of corresponding standardised predictor.
- For rules this is $R_{lk} = |\hat{a}_k|^* \sqrt{S_k(1-S_k)}$
- where $S_k = \frac{1}{N} \sum_{i=1}^{N} r_k(x_i)$ — proportion of 1s if we have 0,1 variable → the proportion
- For linear terms this $l_j = |\hat{B}_j|^* \, std\,(l_j(x_j))$
- where $std\,(l_j(x_j))$ is the standard deviation of $(l_j(x_j))$ over data
- Relative importance here — nearly controlling for other variables in $\hat{f}$

## Local Measure of Importance

- Local measure of importance for each point $x$ as the absolute change in predictor when term is removed from the ensemble.
- For the tree part $R_{lk}(x) = |\hat{a}_k| \times |(r_k(x) - S_k)|$
- For the linear predictor $l_j = |\hat{B}_j|^* \, |l_j(x_i) - \bar{l}_j|$
- Where $\bar{l}_j$ is the mean of $l_j(x_i)$ over training set
- $S_k$ — 'the support of that rule'.

- The two can be combined:
  $J_i(x) = l_i(x) + \sum_{x_i \in R_k} R_{lk}(x)/m_k$
- $l_i(x)$ is the importance of the linear predictor
- Second term sums over the importance of the rules divided by the total number $m$ of input variables in the rule
- Can also be averaged over any subset of the input space

## Interactions

- Looking for variables which are involved in an interaction
- Compare this to what you would expect for no interactions present
- Reference distribution is computed using a bootstrap method
- For each variable identified, look to see which variable or variables interacts with
- Can look for 2-way or 3-way interaction — depends on size of tree — if you only have a stump will get no interaction

- Chart - Show whether they interact or not.
- Red is the interaction - what you would expect to get
- Don't add up to 1 → don't have to
- Plot a single variable and show interaction compared to other values
- Haven't got this information from ensembles or trees.
  ↳ In ensemble, you would have to know there was an interaction

## Partial Dependencies
- Can look at how response function changes with regard to important variables
- Can look at two variables at a time to understand interactions

- Plot $X_1$ against $X_5$ for when $X_1$ is $= -1$ and $X_1 = 1$
- Show $-1, 0, 1$ on X-axis for $X_5$.
- When there is zero height, this highlights interaction effect.

- When $X_1 = 1$, $X_5$ doesn't enter equation → no effect → has no bearing on $y$

- If pattern is roughly the same between (both) → no interaction

## Rules
- Support: proportion of cases in this when $X_1 = 1$ and $Y_2 = 1$ (vote)
- coef: coefficient in your equation $\beta$ or $\hat{a}$
    - $X_1$ not in lead 1 → so no -1 in the model here
- Be careful with factors
- Categories are numbered 1, 2, 3
- Factor coded as $-1, 0, 1$ rules will contain reference to 1, 2 and 3

- Use default as much as possible
- Maybe change tree size
- Make sure results make sense
- Be careful in comparing important values to random forest

- ensemble method which combined the prediction of a large number of simple models
- Predictive power of these rules is combined via regularized regression
- To clone rules - large number of CART trees grown
- Along with these rules, each predictor variable is also added to ensemble - to account for any linear dimension within the data that trees are poor at approximating

- Cross validation error - indicated how well a model will generalize to an independent data set
  REllc Fit use a k-fold cross val tech where original data is split into k subsets
  • A single subset is kept to validate a model built from remaining k-1 subset
  • Repeated with each subset taking a turn at validation, so that the final model may be a combination of the k-models.
  • By taking an expected error rate from these panels, one can estimate how well a model will generalize to an independent data set.