# DA - MARS

- form of regression analysis
- Non parametric method
- Can be used for classification
- Makes no assumptions about the underlying functional relationships between the dependent and independent variables
- Fits piecewise linear regression
- Use of separate regression slopes in distinct intervals of the independent variable space
- Also allows for interaction between variables

## Goal

- Model dependence of a response variable $y$ on one or more predictor variables
- Describe system that generated data by $y = f(x_1, \ldots x_n) + \varepsilon$ (underlying ground)
- We construct $\hat{f}(x_1, \ldots, x_n)$ which is an estimate of this function.

## How?

Construct a form of Basis function $\beta_i$ in a certain manner

$$\hat{f} = \sum_{m=1}^{m} a_m \beta_m(x)$$

- An extension to stepwise linear regression
- Modification of CART to improve it's performance in regression setting

## What are basis functions?

- Piecewise linear basis functions, sometimes called hinge functions

$$(x-t)_+ \quad \text{and} \quad (t-x)_+$$

- The $+$ means the positive part

$$(x-t)_+ = \begin{cases} x-t & \text{if } x \geq t \\ 0 & \text{otherwise} \end{cases}$$

$$t-x = \begin{cases} t-x & \text{if } x < t \\ 0 & \text{otherwise} \end{cases}$$

- Basis function can only be used once
- Interaction has to be between 2 separate variables
- If we have N cases and p variables with all cases having different values for each variable, there are 2Np possible basis functions

- Each knot has 2 basis function ≤ ≥
- Predict the knots across interval

One of 3 forms • A constant
• A hinge function $(x-t)_+$ or $(t-x)_+$
• Product of hinge function

Graph of hinge function - called a reflective pair
• (constant) term as part of the penonimal
• Take both sides as a variable for fit step
• Build model an prive backcncost

## Reflected Pair of Function

- For each input $x_i$, calculate reflected pair with knots at each individual value of $x_i$
- Collection of Basis Function C
  $$C = [(x_j - t)_+, (t - x_j)_+]$$
  $$t \in [x_{1j}, x_{2j}, \cdots, x_{nj}] \text{ individual unique value for each variable}$$
  $$j = 1, 2, \cdots p \text{ variables}$$

- Stepwise linear regression
- Use the reflected pair or product of reflected pairs
  $$f(x) = \beta_0 + \sum_{m=1}^{M} \beta_m h_m(x)$$
- $h_m(x)$ is a function in C - pair of reflected function
- Or a product of 2 or more such function
- $\beta_m$ estimated by minimizing the residual sum of square

## Procedure

- Start with constant function $h_0(x) = 1$  (like fitting mean of y to data)
- Define set of terms in model as $\phi$
- All product of a function $h_m$ in the model set $\phi$ with one of reflected pairs in C
- We add this to the model $\phi$ a term of the form:

$$\hat{\beta}_{m+1}\, h_i(x)\cdot(x_j - t)_+ + \hat{\beta}_{m+2}\, h_i(x)\cdot(t - x_j)_+ \qquad h_i \in \phi$$

- Calculate weights
- We add the term that produce the large decrease in training error - SSE

### Example

- First term is a constant $h_0$
- Second term
  → Now consider a function of the form $\beta_1(x_j - t)_+ + \beta_2(t - x_j)_+$
- For example, suppose this is $\beta_1(x_7 - t)_+ + \beta_2(t - x_7)_+$ for some value of $t$
- Multiplication by a constant does not change things

- Next stage we add $h_m(x)\cdot(x_j - t)_+$ and $h_m(x)\cdot(t - x_j)_+$
- We have 3 choices for $h_m$:
  - $h_0(x) = 1$ i.e. constant.
  - $h_1(x) = (x_7 - t)_+$
  - $h_2(x) = (t - x_7)_+$
- $h_1, h_2$ or interaction terms, can only have interaction with a variable that is already there
- At the end of this process we have overly large model
- We set a number limit of # of terms in the model

### Hierarchical Structure

- There is a hierarchical structure.
- Two way interactions are included if main effects are there
- A 4-way product only included if one of it's 3 way components is in the model
- Have facility to restrict order of interaction and particular terms in interaction
- each input can appear at most once in a product

### Forward Pass
Add terms in pairs until:
  • Reach maximum number of terms
  • Adding a term changes $R^2$ by less than 0.001

- Required an $R^2$ of 0.999 or more
- GRJQ < -10
- No new term increase $R^2$
- Default to nk is min $(200, max (20, 2 \cdot ncol (x)) ) + 1$

– GRSq –EARTHS estimate of the generalisation performance of the model

## Backward Pass

– Need to prune the number of basis function
– Can specify maximum number of terms here – nprune
– Assume we have nb basis function
  For each subset 1... nk Find best subset in terms of lowest RSS
– Then look at each subset and calculate GCV to find lowest value
– Gives us a term of basis function to use
– Calculate the coefficients, residual and fitted value using lm (no standard errors)

## Comments

– Typically overfit
– Go backward and prune and remove terms which cause smallest decrease in RSS
– Estimate model for each size $\lambda$ – # terms in model
– Use cross validation to estimate $\lambda$

G CV formula 
$$GCV(\lambda) = \frac{\sum_{j=1}^{N} (y_i - \hat{P}_\lambda (x_i))^2}{(1 - M (\lambda)/N)^2}$$

– $m (\lambda)$ is the effective # of params in model
– r linearly independent basis function
– K knots
– $M(\lambda) = r + JK$
– Substitute for cross validation

## GRSq

$GRSq = 1 - \frac{GCV}{GCV_{null}}$

– When GCV null is the measure of an intercept only model

DA - MARS

- The GCV and GRSq are measure of the generalisation ability of the model, ie how well the model would predict using data not in the training set.
- like adjusted $R^2$

OUTPUT PLOTS

1- Model Selection
- Plot of $R^2$ and GRSq
- $R^2$ a normalised form of the RSS
- GRSq - Meas of generalisation ability of model
- Should not run bigger -- else should be an increased penalty being applied to the GCV as # of model parameters increase
- GCV - Generalised cross validation : trade off goodness of fit against model complexity

2- Residuals vs Fitted
- Shows residuals of each value of predicted response ie (remainder of observed -expected value).
- Should be as close to zero as possible
- Constant variance of residuals not as important as in linear model
- Highlights outliers

3 - Cumulative Distribution
- Cumulative sum of absolute value of residuals
- Ideally, start at 0 and shot sharply up to one
- Can calculate mean etc.

4 - Quantile Quantile
- Compares the distribution of residuals to the assumed normal distribution
- Want all points to fall on a straight line
- Look for outliers

### Variable Importance

- look relationship between var importance in model compared
- Variance of variable importance is high
- Run different datasets → bootstrap → may get different orders
- Highly correlated variables - one variable is chosen over the other
- In interactions, each variable gets credit for entire term
- Uses 3 different criteria:
  - # times variable appears in subset relative to predicted subset of variables size
  - Decrease in RSS for each subset relative to predicted subset
  - Decrease in GCV for each " "
- Normalized to largest dense is 100
- ie variable with "highest importance" will have large number for #subset and number will be low for GCV and RSS

### Comment

- Can calculate CI and PI
- Can develop a variance model
- Assumes errors are independent
- Ability to operate locally
- Regression surface built up parsimoniously using non-zero components any where they are needed
- Can use GLM method of to build model that has dataset combined to determine final weights

- Can be computationally slow to fit to model and for more complex than a tree
- Has advantage of being able to be used as binary of continuous outcome
- Non parametric regression procedure - model no assumption about underlying functional relationship between dependent and independent variable
- Low FP and identifies interactions
- Typically overran
- Automatically model non-linearity and interactions

DA - MARS

- More flexible than linear regression models
- Simple to understand
- Often requires little or no data preparation - effect of outliers is contained
- Automatic variable selection
- Tend to have a good bias-variance trade off - models are flexible enough to model non-linearity and variable interaction (low bias) yet the constrained form or MARS basis function prevents too much flexibility (low variance)
- Suitable for handling large datasets
- Cross validation and related techniques must be used for validating the model
- Doesn't give as good a fit as boosted trees but can be built quicker and are more interpretable

- Can enter (impred - variables) their be less important and shouldn't be added in as considered for a hinge function - reduces complexity, fit is the low
- Cross validation - partition data into hold subsets, repeatedly build model on all but one of these subsets, mean performance in the left out data

Variable importance
- Asks to relative importance of the variables in a tree
- Measured by impurity improvement
- Look at primary splitter for each node and all the surrogate splits used in place for every node
- Can control the number of surrogates
- Calculated over tree

- For each variable we add up the improvement scores generated by variable in primary splitters.
- Also go through each node this variable used as a splitter and add in their improvement score
- Gets a raw imp input for
- Variables that dont appear get a zero
- Rescale the results so best variable is 100 and work down
- Top competitor gets credit but the second best splitter in a node gets zero credit for being second best.

- 95% CI for accuracy $\pm 1.96 \sqrt{\frac{P(1-P)}{N}}$
- No information rate: how well we do if predict everything as a 'yes'
- P-value: how better your did than what is
- McNemar P-Value - paired chi-square test. Suggests no difference (luchy or GS) in example probably of 1
- Tells if FP and FN are to same, if they are the same they prove zero otherwise concluded they are not the same
$$\frac{(B-C)^2}{B+C}$$

16/04/16    DA   -   Multivariate Adaptive Regression Splines   MARS

### What are MARS?
- Form of regression analysis
- Non parametric method
- Can be used for classification
- Makes no assumptions about the underlying functional relationships between the dependent and independent variables
- Fit piecewise linear regressions
- Use of seperate regression slopes in distinct intervals of the independent variable space
- Also allow for interaction between variables
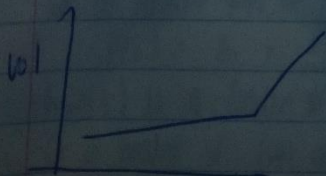- Divides space - "very step wise regressions"

### Goal
- Model the dependance of a response variable y on one or more predictor variables.
- Set of data    $y_i, x_{1i}, ..., x_{ni}$
- Describe system that generated data by $y = f(x_1, ..., x_n) + \varepsilon$   (underlying function)
- We construct $\hat{f}(x_1, ..., x_n)$ which is an estimate of this function

### How?
- Construct a series of Basis functions $\beta_i$ in a certain manner
- $\hat{f} = \sum_{m=1}^{M} a_m \beta_m(x)$
- An extension to stepwise linear regression
- Modification of CART to improve its performance in regression setting
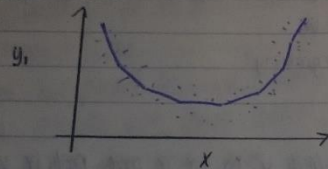
### Cherry Tree Example
- Measurement of girth, height, volume of timber of 31 trees
- Want to predict volume based on girth of tree
- Can fit a SLR of girth to column

Vol

- Given value for intercept, $h(13.8 - Girth)$ use for when $G < 13.8$ and $h(G - 13.8)$ when $G > 13.8$
- GCV - like cross validation
- RSS - like $R^2$
- GRSq

- Result is value of $h()$ function w/ the basis vector

Quadratic example



$y_i \sim X$

Return 8 parameters → intercept and 7 $h()$ function

- Can also use interactions in model to predict

Form of Models
- MARS build model in form of $\hat{f}(x) = \sum_{i=1}^{h} c_i \beta_i(x)$
- $c_i$ are coefficients
- $\beta_i(x)$ are called basis function

What are basis function?
- Piecewise linear basis function sometimes called hinge function
  $(x-t)_+$ and $(t-x)_+$
- The $+$ means the positive part $(x-t)_+ = \begin{cases} x-t & \text{if } x > t; \\ 0 & \text{otherwise} \end{cases}$

$(t-x) = \begin{cases} t-x & \text{if } x < t \\ 0 & \text{otherwise} \end{cases}$

- Sometimes called hinge function

16/04/16   DA MARS

- Basis function can only be used once
- Interaction has to be between 2 separate variables

- If we have N (obs) and p variables with all rows having distinct values for each variable, there are $2Np$ possible basis functions
    - each knot has 2 basis functions $<$ and $>$
    - doesn't take into account interactions.

Basis function in MARS?
One of 3 forms:
• a constant
• a hinge function $(x-t)_+$ or $(t-x)_+$
• Product of hinge functions

Graph of hinge function - called a reflective pair.
- consider them a) pairs at the beginning
- takes both sides @ a variable for first step
- Build model then prune backwards

Reflected Pair of functions
- For each input $x_i$ calculate reflected pairs with knots at each individual value of $x$.
- Collection of Basis functions C
  $C = [ (x_j - t)_+, (t - x_j)_+ ]$
- $t \in [x_{1j}, x_{2j}, \ldots x_{Nj}]$ individual unique value for each variable
- $j = 1, 2, \ldots p$ variables

- Stepwise linear regression
- Use the reflective pairs or product of reflective pairs
  $$f(x) = \beta_0 + \sum_{m=1}^{M} \beta_m h_m(x)$$
- $h_m(x)$ is a function in C - pair of reflective functions
- Or a product of two or more such functions.

$\beta_m$ estimated by minimizing the residual sum of squares

### Procedure
- Start with constant function $h_0(x) = 1$ (like fitting mean of y to data)
- Define set of terms in model as $\Phi$
- All products of a function $h_m$ in the model set $\Phi$ with one of refined pairs in $C$
- We add this to the model $\Phi$ a term of the form

$$\beta_{m+1} h_1(x) \cdot (x_j - t)_+ + \beta_{m+2} h_1(x) \cdot (t - x_j)_+ \qquad h_1 \in \Phi$$

- Calculate weights
- We add the term that produces the largest decrease in the training error $-SSE$

### Example
- First term is a constant $h_0$
- Second term:
  → We now consider a function of the form $\beta_1 (x_j - t)_+ + \beta_2 (t - x_j)_+$
- For example, suppose this is $\beta_1 (x_7 - t)_+ + \beta_1 (t - x_7)_+$ for some value of $t$.
- Multiplication by a constant does not change things

- The next stage we add $h_m(x) \cdot (x_j - t)_+$ and $h_m(x) \cdot (t - x_j)_+$
- We have three choices for $h_n$:
  - $h_0(x) = 1$ i.e. constant
  - $h_1(x) = (x_7 - t)_+$
  - $h_2(x) = (t - x_7)_+$
- $h_1, h_2$ are interaction terms, can only have interaction with a variable that is already there
- At end of this process we have a very large model
- We set a limit on the number of terms in model

### Hierarchical Structure
- There is a hierarchical structure
- Two way interaction are included if main effects are there
- A 4-way product only included if one of its 3-way components is in model
- Have facility to restrict order of interaction and particular terms in interaction

DA - MARS
- each input can appear at most, one in a product.

Forward Pass
- Adds terms in pairs until you
  - reach maximum number of terms
  - Adding a term changes $R^2$ by less than 0.001
  - Reached an $R^2$ of 0.999 or more
  - GRSQ $< -10$
  - No new term increases $R^2$
  - Default for $n_n$ is $\min(200, \max(20, 2 \times ncol(x))) + 1$

GRSQ - earth's estimate of the generalisation performance of the model.

Backward Pass
- Need to prune the number of basis functions.
- Can specify maximum number of terms here - nprune.
- Assume that we have $n_b$ basis functions.
- For each Subset $1 \ldots n_b$ find best subset in terms of lowest RSS.
- Then look at each subset and calculate GCV to find lowest value
- Gives us a series of basis functions to use
- Calculate the coefficients, residuals and fitted values using lm. (no standard errors)

Comments
- Typically overfit)
- Go backward and reduce and remove terms which cause smallest decrease in RSS.
- Estimate model for each size $\lambda$  - # terms in model
- Use cross validation to estimate $\lambda$

General Cross Validation formula
$$GCV(\lambda) = \frac{\sum_{i=1}^{N} (y_i - \hat{F}_\lambda(x_i))^2}{(1 - M(\lambda)/N)^2}$$

- $M(\lambda)$ is the effective # of parameters in model
  - $r$ linearly independent basis functions
  - $k$ knots
  - $M(\lambda) = r + 3k$
- Substitute for cross validation

GRSq

$$GRSq = 1 - \frac{GCV}{GCV_{null}}$$

- Where $GCV_{null}$ is the GCV of an intercept only model
- The GCV and GRSq are measures of the generalization ability of the model ie how well the model would predict using data not in training set.
- like adjusted $R^2$

Output plots

1st - Model Selection
- Plot of $R^2$ and GRSq
- R-Square - a normalized form of the RSS
- GRSq - measure of generalization ability of the model
- Should not run together - there should be an increased penalty being applied to the GCV as # of model parameters increases
- GCV - Generalized Cross Validation: trades off goodness of fit against model complexity

2 - Residuals vs Fitted
- Show residual for each value of predicted response (variance of observed - expected value)
- Should be close to 0 as possible
- Constant variance of residual not as important as in linear model

3 - Cumulative Distribution
- Cumulative sum of absolute value of residuals
- Ideally Start at 0 and Shoot straight up to one
- Can calculate mean proportion 50% residual value, and where 95% of values are predicted within

4 - Quantile, Quantile
- Compares the distribution of residuals to a normal distribution
- Want all points to fall on a straight line
- look for outliers

Variable Importance
- lose relationship between var importance in model compared to data
- Variance of variable important is high
- Run different datasets → bootstrap → may give different answers
- Highly correlated variables - one variable chosen over the other
- In each interactions each variable gets credit for entire term.

Comments
- Can calculate prediction intervals and CI.
- Develop a coarse model
- Assume errors are independent
- Ability to operate locally
- Regression surface built up parsimoniously using non-zero components only where they are needed
- Can use GLM method after the basis functions have been continued to determine final weights

- Can be computationally slow to fit the model and for more complex than a tree
- Has advantage of being able to be used on binary or continuous outcome
- non-parametric regression procedure, makes no assumptions about underlying functional relationship between dependent and independent variables
- will fit and identifies interactions
- Typically overfit

## DA - Combining Classifiers

Two classifier A and B on same test
- $n_{00}$ # cases misclassified by A and B
- $n_{01}$ # misclassified by A but not B
- $n_{10}$ # misclassified by B but not A
- $n_{11}$ # classified correctly by A and B

### McNemar Test
- Info only in diagonal) $n_{01}$ but $n_{10}$
- No difference between classifier would expect $\frac{n_{01}+n_{10}}{2}$ in each cell

$$\chi_1^2 = \frac{(O-E)^2}{E} = \frac{\left(n_{01} - \frac{n_{01}+n_{10}}{2}\right)^2}{\frac{n_{01}+n_{10}}{2}} + \frac{\left(n_{10} - \frac{n_{01}+n_{10}}{2}\right)^2}{\frac{n_{01}+n_{10}}{2}}$$

Reduces to $\frac{\left[|n_{01}-n_{10}|-1\right]^2}{n_{01}+n_{10}}$  $\chi^2$ with $df=1$

### Combining classifiers
- Simple new red if both red
- Otherwise blue
- Ensemble → instead of building one model, build many

- 25 independent classifier each with error rate of 0.35
- Ensemble predicts using majority voting
- Assume classifiers are independent of each other.
- Errors are uncorrelated
- Ensemble makes a wrong prediction if more than half of classifiers predict incorrectly

- Ensemble Error rate
$$e_{ensemble} = \sum_{i=13}^{25} \binom{25}{i} e^i (1-e)^{25-i} = 0.06 \quad \text{instead of } 0.35 \text{ of individual}$$

- Combine to reduce overall error
- 10 classifier 0.4 error rate = overall error of 0.02

### Bias Variance Decomposition

- look at ensembles from the viewpoint of bias and variance

$$\mathbb{E}\,[\,f-t\,]^2 = (\mathbb{E}[f]-t)^2 + \mathbb{E}\,[\,f-\mathbb{E}[f]^2\,]$$

$$\text{MSE} \quad = \quad \text{Bias}^2 \quad + \text{Variance}$$

(expectation of estimated value - true value)²

- For any single estimator $f$ where $t$ = true value
- MSE = Mean Square Error
- Collection of $f$'s $f_1, \dots f_m$     (M models)
- Train each seperately and take average of result $\bar{f} = \sum_{i=1}^{m} f_i$

- Treating the ensemble as a single estimator we can define the variance and bias

$$\mathbb{E}\,[\,\bar{f}-t\,]^2 = (\mathbb{E}[\bar{f}]-t)^2 + \mathbb{E}\,[\,\bar{f} - \mathbb{E}[\bar{f}]\,]^2$$

$$= \text{Bias}^2 \quad + \text{Variance}$$

- Bias stays the same, variance can be reduced

### Ensemble

- Using a committee of models
- When each model is constructed independently, variance of the committee reduce by factor of $\frac{1}{m}$. m=most
- Tighter prediction
- No unique decomposition of MSE

- Bagging
- Random Forest
- Boosting

"the collective knowledge of a diverse and independent body of people typically exceeds the knowledge of any one individual and can be harnessed by voting"