

DA - MARS

Multivariate Adaptive Regression Splines

- Form of regression analysis
- Non parametric - makes no assumptions about the underlying functional relationship between the dependent and independent variables
- Fits piecewise linear regressions.
- Use of separate regression slopes in different intervals of the independent variable space
- Allows for interaction between variables
- Construct a series of basis functions β_i in a certain manner
$$\hat{f} = \sum_{i=1}^M a_i \beta_i(x)$$
- An extension of stepwise linear regression
- Modification of CART to improve its performance in regression setting
- MARS build model in the form of $\hat{f}(x) = \sum_{i=1}^M c_i \beta_i(x)$
- c_i are coefficients
- $\beta_i(x)$ are called basis functions

Basis function

- Piecewise linear basis functions sometimes called hinge functions.
- $(x-t)_+$ and $(t-x)_+$
- The + means the positive part.

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

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

- If we have N cuts and p variables with all cuts having distinct values for each variable, there are 2^{pN} possible basis functions

- Can take 1 of 3 forms - a basis function -
 - a constant
 - a hinge function $(x-t)_+$ or $(t-x)_+$
 - Product of hinge functions

Refined Pairs

- Refined pairs at a particular knot value
- For each input x_i , calculate refined pairs with knots at each individual value of x_i
- Collection of basis functions C
- $C = \{(x_i-t)_+, (t-x)_+\}$
- $t \in \{x_{1j}, x_{2j}, \dots, x_{nj}\}$ individual unique values for each variable
- $j = 1, 2, \dots, p$ variables
- Stepwise linear regression
- Use the refined pairs or product of refined pairs

$$f(x) = \beta_0 + \sum_{m=1}^M \beta_m h_m(x)$$
- $h_m(x)$ is a function in C - pair of reference functions.
- Or a product of two or more such functions.
- β_m estimated by minimizing the residual sum-of-squares

Procedure

- We start with constant function $h_0(x) = 1$
- Define set of terms in model as Φ
- All products of a function h_m in the model set Φ with one of the refined pairs in C
- We add this to the model as a term of the form

$$\beta_{m+1} h_m(x) \cdot (x_i-t)_+ + \beta_{m+2} h_m(x) \cdot (t-x)_+ \quad h_m \in \Phi$$
- Calculate weights
- We add the term that produces the largest decrease in the training error - residual sum of squares

28/04/16

DA - MARS

3

Example:

- First term is a constant h_0
- Second term
- We now consider a function of the form $p_1(x_1 - t) + p_2(t - x_1) +$
- For example, Suppose en turns out to be $p_1(x_7 - t) + p_2(t - x_7)$ for some value of t
- Multiplication by a constant does not change things.

- Next stage we add $h_m(x) \cdot (x_1 - t) +$ and $h_m(x) \cdot (t - x_1) +$
- We have 3 choices for h_m :

$$h_0(x) = 1 \quad \text{i.e. constant.}$$

$$h_1(x) = (x_7 - t) +$$

$$h_2(x) = (t - x_7) +$$

- The last two are interactive terms.
- At end of process we have a very large model
- Set a limit on the number of terms in model

Hierarchical Structure

- Two way interactions included if main effects are there.
- A 4-way product only included if one of its three way components in the model already
- Can restrict interaction
- Each input can appear at most once in a product.

Forward Pass

- Add terms in pairs until you:
 - Reach max # 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 increased R^2

- Default for n_n is $\min(20, \max(20, 2^{\text{ndf}(x)})) + 1$

Backward Pass

- Need to give the number of Basis Functions.
- Can specify max # of terms here - n_{prior}
- Assume that we have n_b basis functions
- For each subset $1 \dots n_b$ find best subset in terms of lowest RSS
- Then look at each of these subsets and calculate GCV to find lowest value
- Then give us set of basis functions to use
- Calculate the coefficients, residuals and fitted values by lm
- Typically overfits.
- Go backwards and reduce and remove the terms which cause the smallest decrease in the residual sum of squares
- Estimate model for each size λ
- Use cross validation to estimate λ

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 number of parameters in model
- r linearly independent basis functions
- k knots
- $M(\lambda) = r + 3k$

$$GRS_g = 1 - \frac{GCV}{GCV_{\text{null}}}$$

- where GCV_{null} is the GCV of an intercept only model
- The GCV and GRS_g are measures of the generalization ability of the model i.e. how well the model would predict using data not in the training set.

28/04/16

DA-MAR

5

OUTPUT PLOTS

1- Model Selection

- Plot of R^2 and GRS_g
- R^2 - a normalized version of the RSS
- GRS_g - measure of generalization ability of the model
- Should not run together - should be on increased penalty being applied to the GCV as # of model parameters increase
- GCV - Generalized Cross Validation: trade off between fit & model complexity

2- Residual vs Fitted

- Shows residuals of each value of predicted response: (remainder of observed - expected value)
- Should be as close to zero as possible
- Constant variance of residuals not as important as in linear model
- Highlight outliers

3- Cumulative Distribution

- Cumulative sum of absolute values of residuals
- Ideally, start at zero and shoot straight up to one
- Can calculate mean proportion 50% residual value at which 95% of values are contained within

4 - Quantile Quantile

- Compare 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 to data
- Variance of variable importance is high
- Run different datasets \rightarrow bootstrap \rightarrow may give different answers
- Highly correlated variables - one variable chosen over the other

- 6
- In interaction, each variable gets credit for entire term
 - Use 3 different criteria:
 - # times variable appears in subset of varied sizes
 - Decrease in RSS for each subset relative to previous subsets.
 - Decrease in GCV for each subset relative to previous subsets

Normalized so that the largest decrease is 100

i.e. variable with "higher importance" will have large number for number and number close to 100 for GCV and RSS

Comments

- Can calculate prediction and (I).
- Develop a variance 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 matrix has been calculated 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 assumption about underlying functional relationship between dependent and independent variables
- Used for and identical means
- Typically useful
- Automatically models non-linear and interaction

28/04/16

DA-MARS

7.

- 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 \rightarrow model is flexible enough to model non-linearity and variable interactions (low bias) yet the constrained form of MARS basis functions 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 linear prod - variables that are less important and suitable for added in as a contribution for a hinge sum - reduces complexity, fit is less sensitive
- Cross validation - partition data into n fold subsets, repeatedly build model on all but one of these subsets, measure performance on the left out data