

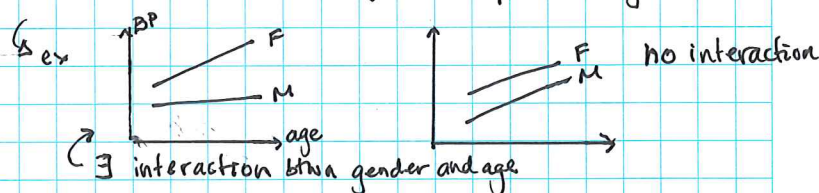
- Recall Generalized Additive Models (GAMs): (can use GAMs for regression and for classification)

$$y_i = \beta_0 + \sum_{j=1}^p f_j(x_{ij}) + \epsilon_i$$

• non-linearity makes GAM generalizable

- GAMs do not allow interactions, i.e.: Where X_j and X_k interact in their "effects" (non-causal) on Y if the "effect" of X_j on Y depends only on the value of X_k

→ simplifies the model.



- (RCode11)

- line 152: $ns(\cdot)$ = natural spline

- education is an ordinal variable in the model and has dummy variables (ex. High School, College, ...)
- Plots:
 - $ns^{(yr)}$ vs yr: trend is from inflation and ~~starting~~
 - $ns^{(age)}$ vs age: increased salaries over time and retirement (pensions) starting

- line 162: backfitting for GAM

- line 177:
 - we go with model m2 because $Pr(>F)$ it has the smallest
 - a good model has: simple structure, (non linearity), low $Pr(>F)$

Backfitting (1985 Breiman & Friedman)

(comparable to ccf Gauss - Seidel method for solving systems of linear equations)

- Let $S_j: (u_1, \dots, u_n)^T \rightarrow (\hat{u}_1, \dots, \hat{u}_n)$ denote a smoother of u_1, \dots, u_n against X_{1j}, \dots, X_{nj}



Procedure:

- Initialize: mean $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n Y_i$ (est.) ; set all predictor functions $\hat{g}_j = 0, j=1, \dots, p$
- Do until convergence:

- $\forall j^{th} = 1, \dots, p$ predictor: (a) fit $\hat{g}_j = S_j \left(\underbrace{Y - \hat{\mu}}_{\text{center } Y} - \underbrace{\sum_{k \neq j} \hat{g}_k}_{\text{sum of vectors in } \mathbb{R}^n} \right)$. For $j=1, \sum_{k \neq 1} \hat{g}_k = 0$.

- (b) Normalize: $\hat{g}_j \leftarrow \hat{g}_j - \frac{1}{n} \sum_{i=1}^n \hat{g}_j(x_{ij})$

where $\hat{g}_j = (\hat{g}_j(x_{1j}), \dots, \hat{g}_j(x_{nj}))^T, (\hat{g}_j^{ic}(x_{ij}))$. \hat{g}_j is the contribution of the j^{th} variable

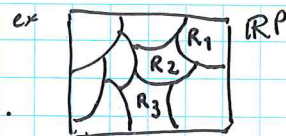
$$Y = (Y_1, \dots, Y_n)^T$$

ie: we fix all variables $k \neq j$ and smooth on the residuals.

- (c) Measure convergence by, eg. ~~max~~ $\frac{\|\hat{g}_{j,new} - \hat{g}_{j,old}\|_2}{\|\hat{g}_{j,old}\|_2}$ (if we have) the relative change of \hat{g}_j vector.

if $\max_j \left\{ \frac{\|\hat{g}_{j,new} - \hat{g}_{j,old}\|_2}{\|\hat{g}_{j,old}\|_2} \right\} \leq \text{tolerance}$, then we have reached convergence.

⇒ other side pg. 2.

• Trees (ISUR chapter 8). CART.

• CART: Classification And Regression Trees (1984 Breiman).

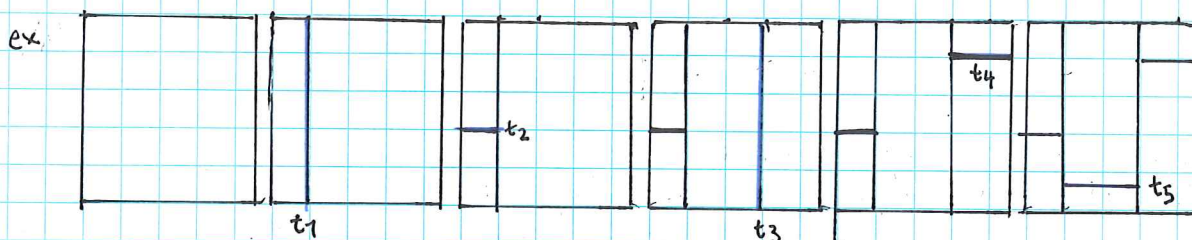
• The underlying model: $g_{\text{tree}}(x) = \sum_{r=1}^M \beta_r \mathbf{1}_{\{x \in R_r\}}$, where $x \in \mathbb{R}^p$. g_{tree} is piecewise constant.

• $P = \{R_1, \dots, R_M\}$ is a partition of \mathbb{R}^p (ie union of disjoint sets $\bigcup_{r=1}^M R_r = \mathbb{R}^p$)
 P covers whole space.

• If the partition is given, then things are easy. (

$$\hat{\beta}_j = \frac{\sum_{i=1}^n y_i \mathbf{1}_{\{x_i \in R_j\}}}{\sum_{i=1}^n \mathbf{1}_{\{x_i \in R_j\}}} \equiv \text{average } y \text{ value among observations in } R_j \text{ (subspace of } \mathbb{R}^p)$$

* How do we find a good partition? (See Fig 8.1, 2 (ISUR) in Slides)



• Fig. 8.3 (lower right): shows that the model is piecewise constant function

• If we grow the tree too deeply, then we overfit.

Estimation (of the best partition).

1. Start with $M=1$. Then $R_1 = \mathbb{R}^p$, $P = \{R_1\}$.

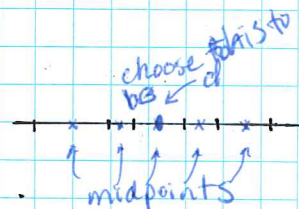
2. Refine the partition $R_1 = R_{\text{left}} \cup R_{\text{right}}$, (disjoint sets)

where $R_{\text{left}} = \mathbb{R} \times \mathbb{R} \times \dots \times \mathbb{R} \times (-\infty, d] \times \mathbb{R} \times \dots \times \mathbb{R}$

$R_{\text{right}} = \mathbb{R} \times \mathbb{R} \times \dots \times \mathbb{R} \times (d, \infty) \times \mathbb{R} \times \dots \times \mathbb{R}$

at the j th axis.

This involves choosing an axis j (ie a predictor X_j) and a split point $d \in \{\text{midpoints between } X_j\text{'s}\}$.



• Choose them such that the $(-\log\text{-likelihood})$ is reduced the most; in certain cases, this is ~~just~~ the same as reducing the RSS the most.

3. Refine the current partition by refining one of the partition cells.

4. Iterate until we have a large tree.

5. Then, prune the tree to get a simpler model using cross-validation

$$\text{cost complexity of pruning} = \min_{T \in T_0} \{ \text{err}(T) + \alpha |T| \}$$

full tree

we want low error

penalty term for # of nodes
ie. # of parameters
so that model is not unnecessarily complex.

(RCode13)

- line 19: "high" is a binary variable.

// End of Lecture //

Trees

$$y_i = \sum_{r=1}^M \beta_r \mathbb{1}_{\{x_i \in R_r\}} + \varepsilon_i$$

- Difficult part is getting partition.
- We limit ourselves to R_r being boxes.
- Find partition using greedy approach. \rightarrow recursive binary splitting
 \hookrightarrow Pick the partition that reduces $(-\log\text{-likelihood})$ the most

Pruning: $\min_{T \leq T_0} \text{err}(T) + \alpha |T|$
 $\uparrow -\log\text{like}(T)$: misclassification rate

ex. 20 observations: 16 yes, 4 no.

• Put them in one leaf node: $n_y \log(\hat{p}_y) + n_n \log(\hat{p}_n)$
 $= 16 \log(\frac{16}{20}) + 4 \log(\frac{4}{20})$

• Split them over two leaf nodes $\rightarrow -\log\text{likelihood}: 10.008$

$\begin{array}{ll} 9Y & \rightarrow 9 \log(1) + 0 \log(0) \\ 7Y, 4N & \rightarrow 7 \log(\frac{7}{11}) + 4 \log(\frac{4}{11}) \end{array} \} \rightarrow -\log\text{like}: 7.23$

(Rcode 13)

- line 55: "dev" \equiv deviance \equiv # of misclassifications.
"size" \equiv # of nodes in tree after each pruning step.
(remaining)
- line 71: "prune.{#}" \rightarrow best tree with \geq {#} nodes

Bagging: "Bootstrap Aggregating"

Procedure: Consider a base procedure, ex. a tree, that gives us some estimated function \hat{g} : ex. $\hat{g}: \mathbb{R}^p \rightarrow \mathbb{R}$ (prediction)
 $\mathbb{R}^p \rightarrow [0,1]$ (classification)

Algorithm: (1) Generate a bootstrap sample $(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$ and compute the corresponding estimator $\hat{g}^*(\cdot)$ by running our base procedure on the bootstrap sample.

(2) Repeat this B times, yielding $\hat{g}^{*1}(\cdot), \dots, \hat{g}^{*B}(\cdot)$
(store these tree objects s.t. we can predict at any point.)

(3) Aggregate the bootstrap estimators by averaging all the bootstrap trees.

• The bagged estimator $\hat{g}_{\text{bag}}(\cdot) = \frac{1}{B} \sum_{b=1}^B \hat{g}^{*b}(\cdot)$ is an approximation of $E^* \hat{g}^*(\cdot)$, which is exact if $B = \infty$.

$\hat{g}_{\text{bag}}(\cdot)$: aggregation
reduces variance but doesn't incur much bias.

individual trees have low bias, high variance
bag has low bias, low variance

// End of Lecture //

- Ref last week's series, problem 1: ^(correction) hint: justify that $y_i^{(-i)}$ coincides with \hat{z}_i , where \hat{z} is the fitted value obtained from regression of z on x_1, \dots, x_p , where z is defined as:

$$z_j = y_j \quad \forall j \neq i$$

$$z_i = y_i^{(-i)}$$
- $H_0: y = \beta_0 + \varepsilon$ (global null)
 $H_1: y = \beta_0 + \sum_j \beta_j x_j + \varepsilon$
- $T_1: H_0': \beta_1 = 0$ in the model $y = \beta_0 + \sum_j \beta_j x_j + \varepsilon$
- If H_0 is true, then $\rightarrow H_0'$
- $p_1 = P_0(1 \times 1 \rightarrow |T_1|)$
 $p_1 = P_0(x^2 > T_1^2) = 1 - F_0(T_1^2)$
 where F_0 is the CDF of T_1^2 under the null
- Goal: Show that for $\alpha \in [0, 1]$, (Prob) $P(p_1 \leq \alpha | H_0) = \alpha \quad \forall \alpha \in [0, 1]$
 \uparrow sample from null under
- $p_1 \sim \text{Unif}(0, 1) \Rightarrow 1 - p_1 \sim \text{Unif}(0, 1)$

$$P(1 - p_1 \leq \alpha | H_0) = P(F_0(T_1^2) \leq \alpha | H_0) = P(T_1^2 \leq F_0^{-1}(\alpha) | H_0)$$

$\uparrow (T_1^2 \sim F_0 \text{ under null})$
 \uparrow property of monotonicity of F_0 and F_0^{-1}

"What is the probability that a random variable will be \leq quantile of distribution from which we are asking the problem?"

(T-test \Rightarrow continuous distribution)

So, $P(T_1^2 \leq F_0^{-1}(\alpha) | H_0) = \alpha$ under continuity of F_0 .

- Assuming under the null, there is an equal chance of getting a very low or ~~and~~ getting a very high pval.

ie. Under the null, the p-value is (completely) random, so ~~it makes no sense to use it to describe model~~
 it for anything really...

// End of Section //