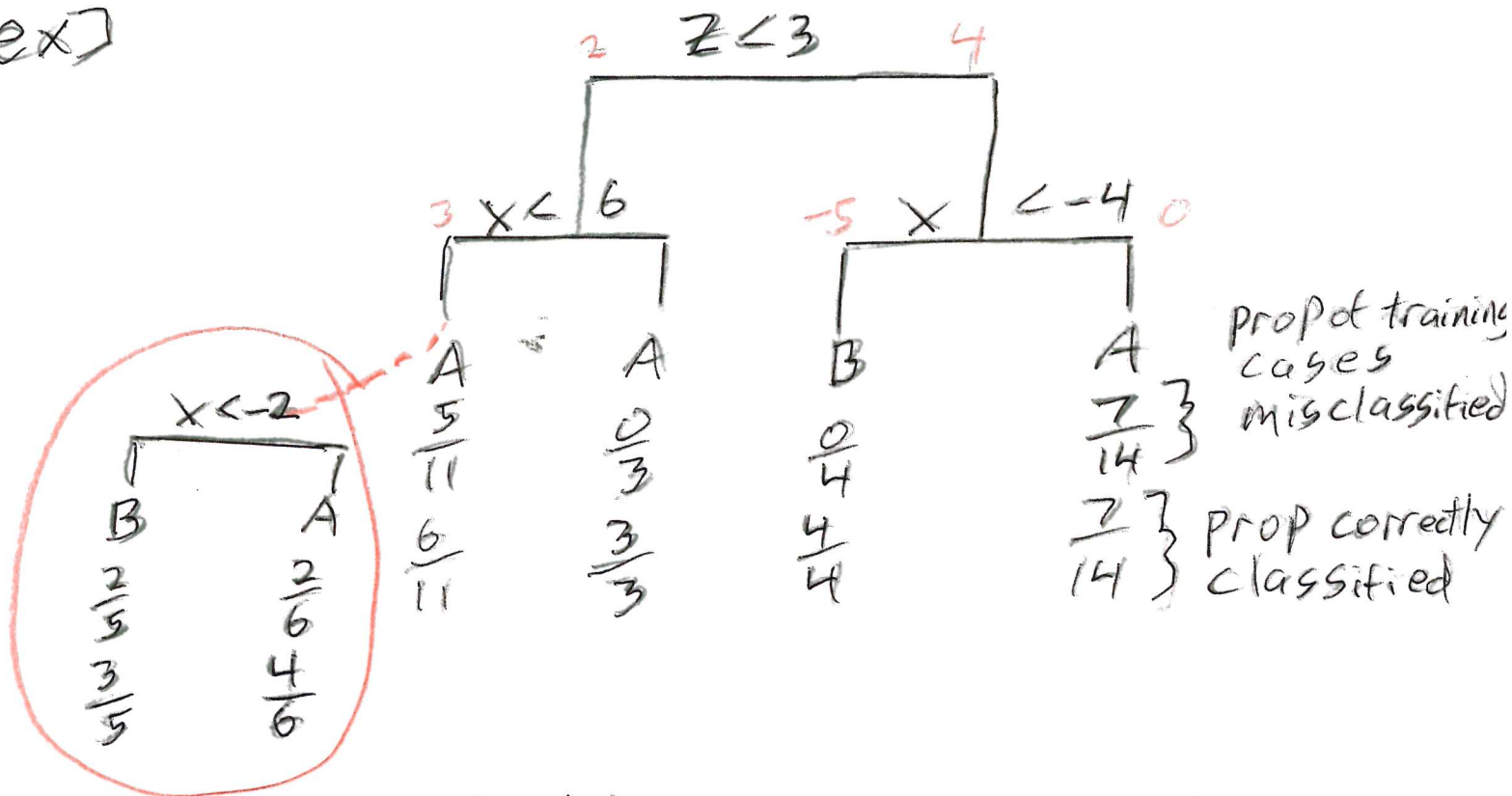


A single vertical line at $z=3$ is HD91
 the 1st partition.

The horizontal line at $x=-4$ is the 2nd partition.

∪ $x=-6$ 3rd ∪ ∪

ex]



add one more partition

- 6] Trees
- i) give prediction rules that can be rapidly and repeatedly evaluated,
 - ii) are useful for screening predictors (variable selection, interactions = products of variables like $x_1 x_2, x_3 x_5 x_7$)
 - iii) can be used to summarize large multivariate data sets

iv) regression trees can be used 91.5
to assess the adequacy of linear models.

7) If Y is a ^{categorical} factor with levels
 $1, \dots, G$, classification rules are of the
form "if $x_1 \leq 2.3$ and $x_3 \in \{A, B\}$, then Y
is most likely to be in level 5."

If Y is numerical, regression rules
are of the form "if $x_2 \leq 1.7$ and
 $x_9 \in \{C, D, F\}$ and $x_5 \leq 3.5$, then $\hat{Y} = 4.75$."

written $x_9 = C, D, F$

8) Trees can be easier to interpret ^{than MLR} when
some predictors are numerical and some
categorical, Trees are invariant to
monotone (increasing or decreasing)
transformations of the predictors x_i .
Trees can handle complex unknown
interactions.

9) Regression trees handle missing values better

than MLR, and can beat MLR if there is nonadditive behavior ($y = m(x) + e$), HD 92

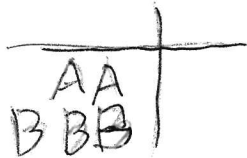
10) In growing a tree, the binary partitioning algorithm recursively splits the data in each node until either the node is homogeneous (roughly 0 training data misclassifications for a classification tree, $y \approx \text{constant}$ for a regression tree) or the node contains too few observations (default ≤ 5).

11) The deviance is a measure of node homogeneity, and deviance = 0 for a perfectly homogeneous node.

12) Often use the mean of the region for \hat{y} , reg tree

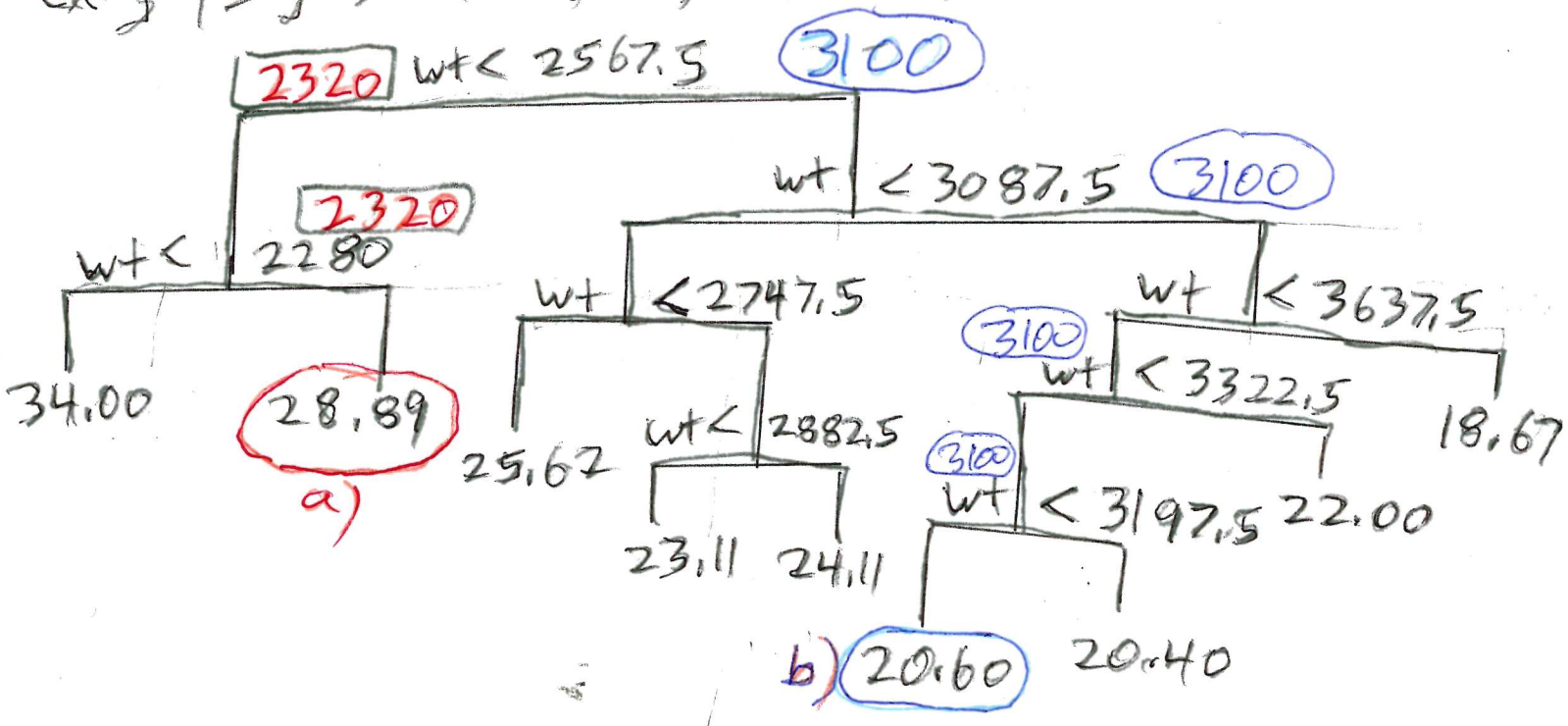
mode

\hat{y} , classification tree



mode = B

ex } $Y = \text{gas mileage}$, $x = \text{weight of car}$



a) Predict mileage if $wt = 2320$

$\hat{y} = 28.89$

see red: 2320 left
2320 right

b) predict mileage if $wt = 3100$

$\hat{y} = 20.60$

see blue R R L L
right left

13} a) Divide the predictor space

(= set of possible values for x_1, \dots, x_p) into training data

J distinct and nonoverlapping regions

R_1, \dots, R_J

b) For every observation that falls in

region R_j , make the same prediction HD 93

$$\hat{y}_{R_j} = \begin{cases} \bar{y}_j = \text{sample mean of } y \text{ in } R_j & \text{reg tree} \\ \text{mode}_j = \text{mode} & \text{class. tree} \end{cases}$$

14) The R_j are high dimensional ^{boxes} rectangles.

Choose R_j so $RSS \stackrel{\text{reg}}{\downarrow} = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$

or $RSS \stackrel{\text{class.}}{\downarrow} = \sum_{j=1}^J \sum_{i \in R_j} \mathbb{I}(y_i \neq \hat{y}_{R_j})$

is small,

15) Let $\{ \underline{x} \mid x_j < s \}$ be the region in the predictor space such that $x_j < s$ where

$\underline{x} = (x_1, \dots, x_p)^T$. Define 2 regions

$$R_1(j, s) = \{ \underline{x} \mid x_j < s \}, R_2(j, s) = \{ \underline{x} \mid x_j \geq s \}$$

and seek "cutpoint" s and j to minimize

$$\sum_{i: \underline{x}_i \in R_1(j, s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: \underline{x}_i \in R_2(j, s)} (y_i - \hat{y}_{R_2})^2 \quad \text{for reg trees.}$$

This can be done "quickly" if p is small (could use order statistics),

Then repeat the process looking for the best predictor and best cutpoint in order to split the data further so as to minimize the RSS within each resulting region. Only split one of the regions. So now there are regions R_1, R_2 , and R_3 . Continue this process until a stopping criterion is reached such as no region contains less than 5 obs's (and stop splitting if the region is homogeneous, eg all Y in region belong to class k for classtree).

16] Classification trees are usually inferior to LDA and SVM. Bagging, random forests, or boosting makes trees more competitive.

17] Trees use regions R_1, \dots, R_J . If J is too large, the tree overfits (not good for test data). One

strategy is to grow a large tree T_0 with J_0 regions. Then "prune" it to get a subtree T_α with J_α regions.

18) cost complexity pruning = weakest link pruning:

Let $T \subseteq T_0$, $\alpha \geq 0$, and $|T| = \#$ terminal nodes of tree T . Each terminal node corresponds to a region (hyperrectangle) R_i . Let R_m be the region corresponding to the m th terminal node and \hat{Y}_{R_m} be the predicted response for R_m .

For each value $\alpha > 0$, there corresponds a subtree $T \subseteq T_0$ such that

(*) $\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{Y}_{R_m})^2 + \alpha |T|$ is as small as possible. (Use $I(y_i \neq \hat{Y}_{R_m})$ for classification.) Note that $\alpha = 0$ has $T = T_0$ and (*) = $RSS(T_0)$ = training data

RSS for T_α .

19) Much like lasso, as α increases, there is a sequence of nested subtrees

(**) $T_0 \supseteq T_{\alpha_1} \supseteq T_{\alpha_2} \supseteq \dots \supseteq T_{\alpha_m}$. Then branches get "pruned" from T_0 in a nested and predictable fashion.

20) a) Build tree T_0 , stopping when each node has ≤ 5 obs's (or is homogeneous)
further splitting would cause fewer than 5 obs's

b) Do 18) to obtain subsequence (**).

c) use k -fold CV to choose $\alpha = \alpha_j^*$.

For each $i \in \{1, \dots, k\}$

i) repeat steps a) and b) on all but the i th fold.

ii) Evaluate the mean square prediction error on data left out of fold i as a function of α .

Average the results for each value of α and pick α_j^* to minimize the average error.

d) Use tree $T_{\alpha_j^*}$ from b).

For a regression tree,

$$MSE_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ji} - \hat{y}_j(i))^2 \quad \text{for data}$$

y_{ji} in the left out fold i .

$$CV_{(K)} = \frac{1}{K} \sum_{i=1}^K MSE_i$$

For a classification tree, use

$$MSE_i = \frac{1}{n_i} \sum_{j=1}^{n_i} I[y_{ji} \neq \hat{y}_j(i)] =$$

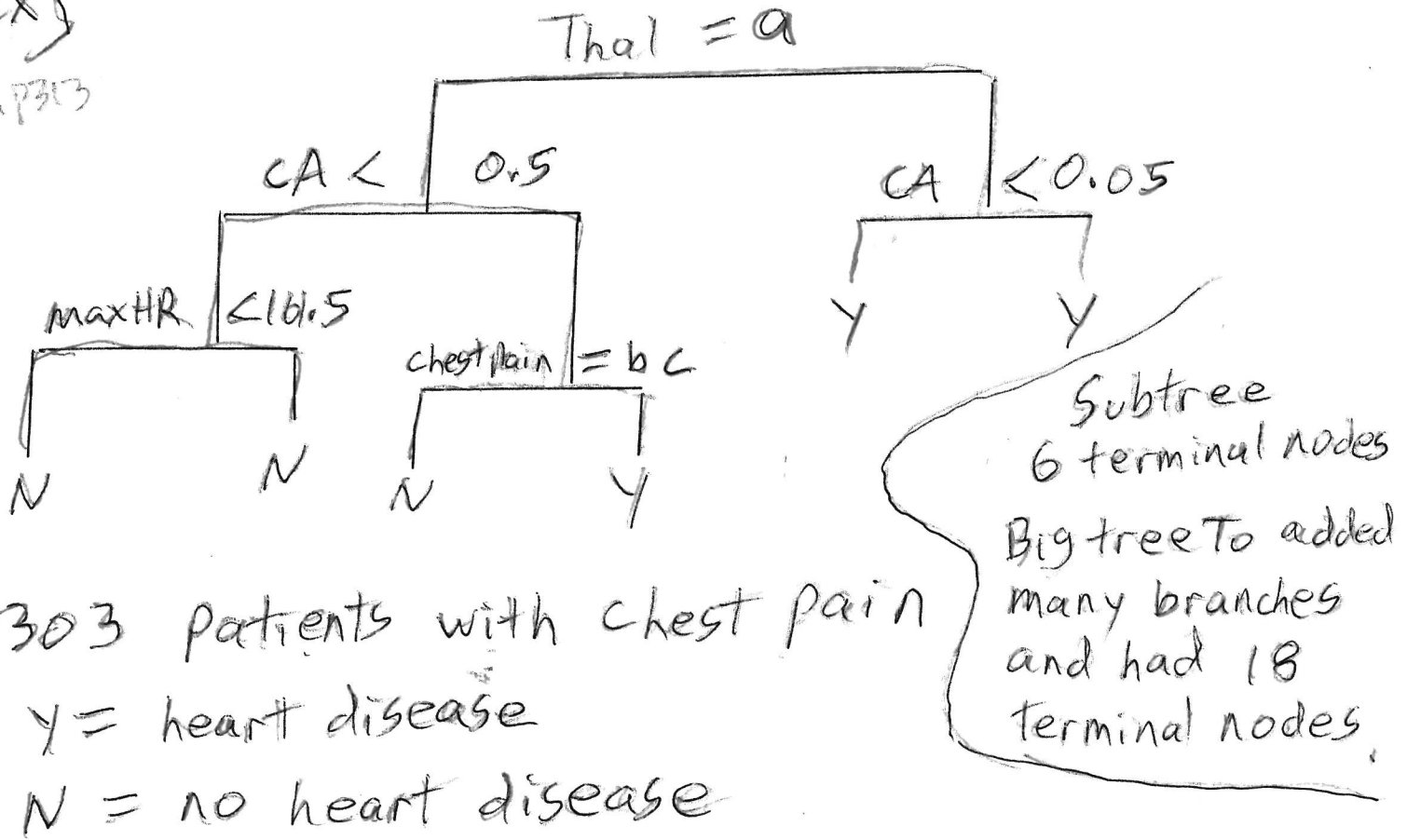
proportion misclassified in the i th fold.

If $n_i = \frac{n}{K}$, then

$$CV_{(K)} = \left[\frac{1}{n} \sum_{j=1}^n (y_{ji} - \hat{y}_j(i))^2 \right], \text{ reg tree}$$

$$\left[\frac{1}{n} \sum_{j=1}^n I[y_{ji} \neq \hat{y}_j(i)] \right], \text{ class tree}$$

ex] J.P.313



Thal = Thallium stress test: a = normal

chestpain a = typical aginal, b = atypical aginal,

c = non-aginal pain, d = asymptomatic

maxHR maximum heart rate?

CA blood calcium?

21] For a tree T_α , $Y_i = m(\underline{x}_i) + \sigma_i e_i$

and $\hat{y} = \hat{m}(\underline{x}_i) = \sum_{m=1}^{J_\alpha} c_m I(\underline{x}_i \in R_m)$ where

T_α uses regions R_1, \dots, R_{J_α} and $c_m = \hat{y}_m$ for $\underline{x}_i \in R_m$.

22) Trees can handle categorical variables (qualitative, factors) without creating indicators = dummy variables. HD 96

23) Bagging was used with the bootstrap:

compute T_1^*, \dots, T_B^* with the bootstrap

and the sample mean $\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^*$

is the bagging estimator.

24) For a regression tree, draw a sample of size n with replacement

from x_1, \dots, x_m . Fit the tree and

find $\hat{f}_1(x)$. Repeat B times to get

$\hat{f}_1^*(x), \dots, \hat{f}_B^*(x)$. Then the bagging

estimator $\hat{f}_{\text{bag}}^*(x) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i^*(x)$.

The trees are not pruned, so terminate

when each terminal node has 5 or

fewer obs's or is homogeneous.

if further splitting would result in fewer than 5 obs's

25] For classification, draw a sample of size n_j from each group with replacement. Let

$$\hat{f}_i^*(x) = j_i(x) \in \{1, \dots, G\} \text{ where}$$

Y takes on G levels (groups) $1, \dots, G$.

Compute $\hat{f}_1^*(x), \dots, \hat{f}_G^*(x)$ and let

$$m_k = \# j_i(x) = k \quad \text{for } k = 1, \dots, G.$$

Take $\hat{f}_{\text{bag}}^*(x) = d$ where $m_d =$

$$\max\{m_1, \dots, m_G\}.$$

26] As before, the prob that x_j is not in the bootstrap dataset $\rightarrow e^{-1} \approx .3679$
 $\approx \frac{1}{3}$ as $n \rightarrow \infty$

$$\frac{1 - \frac{1}{n}}{1} \dots \frac{1 - \frac{1}{n}}{n}$$

$$\left(1 - \frac{1}{n}\right)^n \rightarrow e^{-1}$$

prob x_j is not the k th obs
 $= 1 - \frac{1}{n}$ since each obs
 has $\frac{1}{n}$ chance of being
 selected for the k th position

27] For each bootstrap data set b , HD97

let x_{i1}, \dots, x_{ik_b} be the k_b obs's not used in b . These are out of bag (OOB)

obs's. Predict \hat{y} for each OOB obs.

Doing this for all B bootstrap data sets produces about $\frac{B}{3}$ predictions for each

x_i . Let $\hat{y}_{io} = \begin{cases} \text{ave } \hat{y}_i & \text{reg tree} \\ \text{mode level} & \text{class} \end{cases}$

$$\text{OOB MSE} = \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{io})^2 & \text{reg} \\ \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_{io}) & \text{class} \end{cases}$$

The OOB MSE is "virtually equivalent" to the leave one out CV estimate for sufficiently large B .

28] Bagging typically gives better accuracy than a single tree.

29) For classification trees, let

97.5

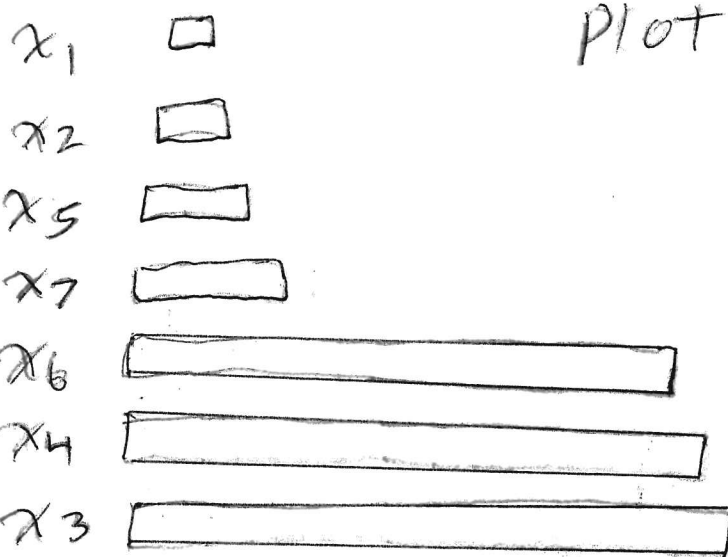
\hat{P}_{mk} = proportion of training obs's in R_m from the k th class.

$$\text{Gini's index} = \sum_{k=1}^G \hat{P}_{mk} (1 - \hat{P}_{mk})$$

is small if all \hat{P}_{mk} are close to 0 or 1.

30) For bagging with B trees, a measure of variable importance can be computed for each variable using splits for each variable.

variable importance plot



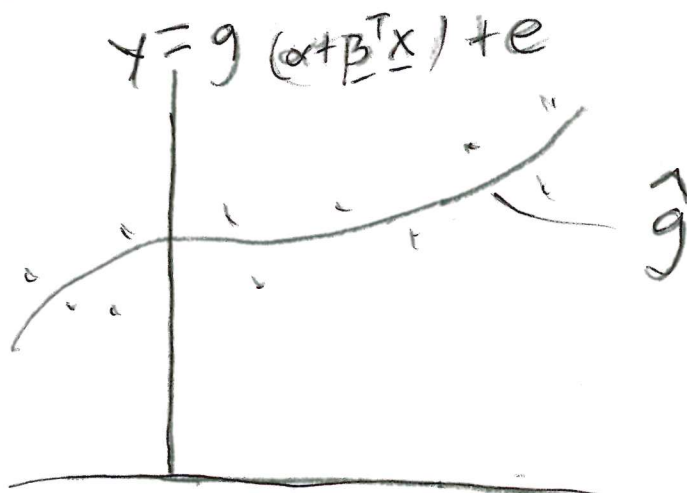
x_3, x_4, x_6 are the most important predictors for the tree

31) If $Y = \alpha + \sum_{j=1}^p S_j(x_j) + e$ or $Y = m(x) + e$

with $m(x) = g(\alpha + \beta^T x)$, then slicing the

$$\text{ESP } \hat{\alpha} + \sum_{j=1}^p \hat{S}_j(x_j) \text{ or } \hat{\alpha} + \hat{\beta}^T x$$

is more effective than partitioning HD 98
 the predictor space \underline{x} with hyperboxes R_m
 (curse of dimensionality (\mathbb{R}^p)).



If $y = g(w)$ and you
 plot w vs y , then you
 "see g !"

So if $y = g(SP) + e$, you
 "see g " up to noise of
 $ESP \approx SP \approx w$ and e .

$$ESP = \hat{\alpha} + \hat{\beta}^T \underline{x}$$

32] **MS project** If $y = m(\underline{x}) + e$, we can
 make prediction intervals for y_f with the
 regression tree using $\hat{y} = ESP = \hat{m}(\underline{x})$ and
 $r = y - \hat{y}$ as before.

Random Forests

33] For random forests, the bootstrap is
 used, but each time a split is considered,
 a random sample of $m = \lceil \sqrt{p} \rceil \approx \sqrt{p}$
 predictors is chosen as split candidates.
Much faster than examining all p variables.

estimator

Random forests produces bootstrap (98.5 trees that are less correlated than bagged trees (that use $m=p$) and the random forests estimator tends to have better test error and OOB error than the bagging estimator.

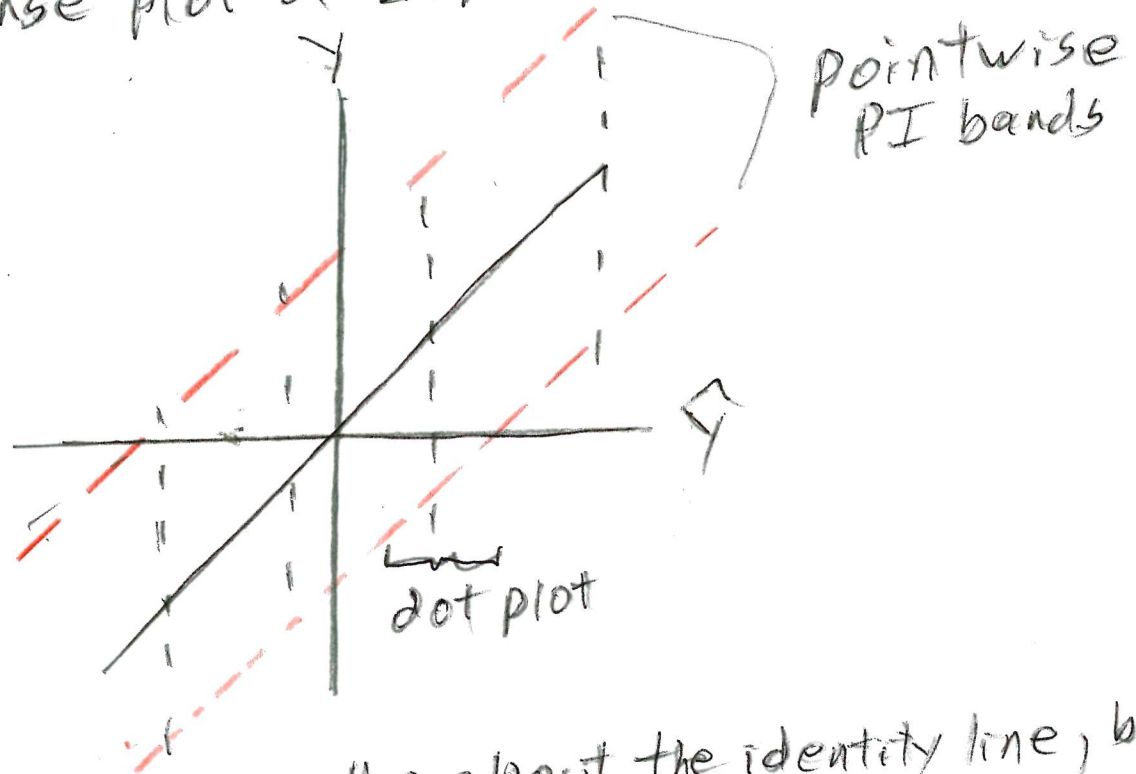
34] B around a few hundred seems to work.

35] If there is a single strong predictor, bagged trees tend to use that predictor in the 1st split. For random forests, the strong predictor is not considered for $\frac{p-m}{p}$ splits, on average. Trees from random forests also tend to be less correlated than bagged trees if there are many correlated predictors.

36] For classification, the null classifier has $\hat{y} = d$ where d is the dominant (mode) class. So if $k\%$ of obs's belong to the dominant class, the test error = $\frac{100-k}{100} \leq 1 - \frac{1}{G}$ where there are G groups, since $k \geq 100 \frac{1}{G}$.
(assumes training data dominant class = that of test data)

classifiers that do not beat the null classifier are very bad. HD 99

37] Since a reg tree uses J_α regions, the response plot of $ESP = \hat{y} = \hat{m}(x)$ vs y looks like



The plotted points scatter about the identity line, but there is a dot plot of n_m cases with $\hat{y} = \hat{y}_{Rm}$ for each of the J_α regions.

(One way Anova models have a similar dot plot, but each dot plot crosses the identity line at $\hat{y}_i = \bar{y}_{i0}$.)

38] A dot plot of z_1, \dots, z_m consists of an axis and m points corresponding to the values of z .

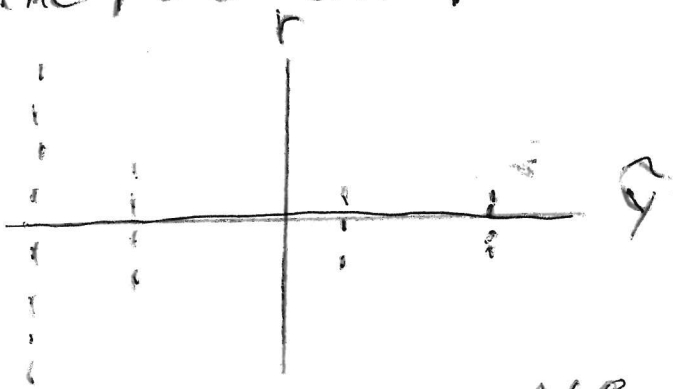
ex] $z_1=3, z_2=-1, z_3=0,$
 $z_4=-4$



99.5

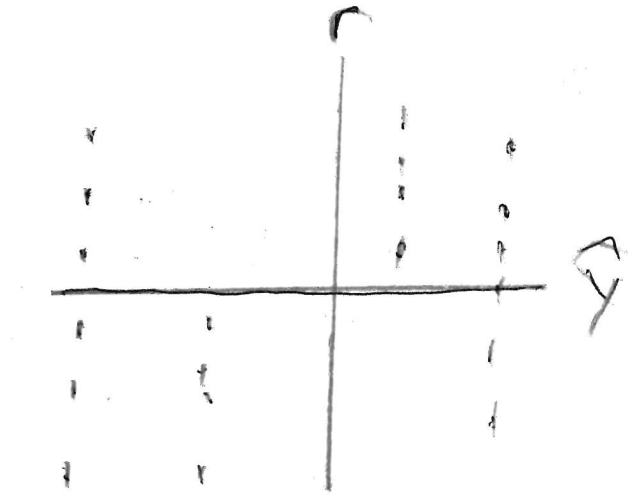
dot plot

39] The residual plot consists of dot plots that scatter about the $r=0$ line.



non constant variance

$$y_i = m(x_i) + \sigma_i e_i$$



evidence that the fit could be better

5582.3

40] The third technique for improving trees is boosting. Like bagging,

boosting can be applied to many statistical methods, not just trees.

41] Boosting for regression trees

i). set $\hat{f}(x) \equiv 0$ and $r_i = y_i, i=1, \dots, N$
 \mathcal{R} training data