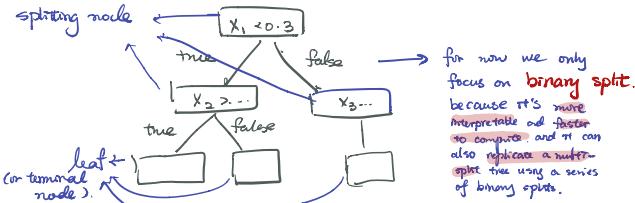


07 Trees



(1) Regression Tree

$$\hat{Y} = \hat{f}(x) = \sum_{m=1}^M C_m I(x \in R_m) \rightarrow x \text{ is categorized in } R_1, R_2, \dots, R_M.$$

where the choice of C_m depends on loss function: $\sum_{m=1}^M \sum_{i: x_i \in R_m} L(Y_i, C_m) = \sum_{i=1}^n \sum_{m=1}^M L(Y_i, C_m) I(x_i \in R_m)$.

e.g. square loss: $\sum_{m=1}^M \sum_{i: x_i \in R_m} (Y_i - C_m)^2 = \sum_{i=1}^n \sum_{m=1}^M (Y_i - C_m)^2 I(x_i \in R_m)$

$$\begin{aligned} \text{minimizing it} &\Rightarrow \sum_{m=1}^M \sum_{i=1}^n 2(Y_i - C_m) I(x_i \in R_m) = 0 \Rightarrow \sum_{i=1}^n \sum_{m=1}^M Y_i I(x_i \in R_m) = \sum_{i=1}^n \sum_{m=1}^M C_m I(x_i \in R_m) \Rightarrow \sum_{i=1}^n Y_i I(x_i \in R_m) = \sum_{m=1}^M C_m \sum_{i=1}^n I(x_i \in R_m) \\ &\Rightarrow C_m = \frac{\sum_{i=1}^n Y_i I(x_i \in R_m)}{\sum_{i=1}^n I(x_i \in R_m)} = \frac{1}{N_m} \sum_{i=1}^n Y_i I(x_i \in R_m) \text{ for } m=1, 2, \dots, M \Rightarrow \text{sample mean in each region.} \end{aligned}$$

absolute loss: $\sum_{m=1}^M \sum_{i: x_i \in R_m} |Y_i - C_m| = \sum_{m=1}^M \sum_{i=1}^n |Y_i - C_m| I(x_i \in R_m)$.

minimizing it give $C_m = \text{sample median in each region}$

steps:

(1) Deciding on splits by using a greedy approach.

i.e. at each splitting node we choose to use the best predictor and threshold to give the smallest loss.

let R_k be the subset of predictors data represent by the k^{th} split, then in a binary split,

$$R_{k1}(j, s) = \{x_i : x_{ij} < s\} \text{ and } R_{k2}(j, s) = \{x_i : x_{ij} > s\}$$

\Rightarrow using square loss.

$$\sum_{i: x_i \in R_{k1}(j, s)} (Y_i - C_1)^2 + \sum_{i: x_i \in R_{k2}(j, s)} (Y_i - C_2)^2$$

Minimising it,

$$\min_{j, s} \left\{ \min_{C_1} \sum_{i: x_i \in R_{k1}(j, s)} (Y_i - C_1)^2 + \min_{C_2} \sum_{i: x_i \in R_{k2}(j, s)} (Y_i - C_2)^2 \right\}$$

As sample mean minimises it

$$\min_{j, s} \left\{ \sum_{i: x_i \in R_{k1}(j, s)} (Y_i - \bar{Y}_{R_{k1}})^2 + \sum_{i: x_i \in R_{k2}(j, s)} (Y_i - \bar{Y}_{R_{k2}})^2 \right\}$$

i.e. we need try a few predictor (x_j) and a few threshold values (s) to find the best predictor and threshold to do the split at each split.

(2) Then decide on tree size.

let T_0 be the full tree (i.e. only a few observations in each leaf). Then if $T \subset T_0$, T is a subtree.

We decide on the tree size by pruning the full tree T_0 that minimises:

$$\text{cost complexity: } C_d(T) = \frac{\sum_{m=1}^{17} N_m Q_m(T)}{d! T!} \text{ for a given } d. \quad \Rightarrow \begin{cases} d=0, T \rightarrow T_0 \\ d \rightarrow \infty, T \rightarrow \text{root tree (no splits).} \end{cases}$$

\downarrow number of leaves of tree T .
 \downarrow penalty on the size of tree.
 \downarrow $\frac{1}{N_m} \sum_{i: x_i \in R_m} (Y_i - \bar{Y}_m)^2$: average square loss.

$$\Rightarrow C_d(T) = Q(T) + d! T!$$

how to choose the tree that minimises $C_d(T)$:

(3) How to prune an optimal tree for a given d :

(1) Grow a full tree T_0 , compute $C_d(T_0) = Q(T_0)$ and set $d_0 = 0$

(2) Prune off the weakest terminal split from T_0 to get T_{d_1} :

$$T_{d_1} = \arg \min_{T \subset T_0, |T_0 - T|=1} Q(T) - Q(T_0). \quad \text{i.e. remove the split with the least loss reduction}$$

and set $d_1 = Q(T_{d_1}) - Q(T_0) \Rightarrow$ improvement of error.

(3) Prune off the weakest terminal split from T_{d_1} to get T_{d_2} :

$$T_{d_2} = \arg \min_{T \subset T_{d_1}, |T_{d_1} - T|=1} Q(T) - Q(T_{d_1})$$

and set $d_2 = Q(T_{d_2}) - Q(T_{d_1})$

(4) Repeat until the root tree is left

\Rightarrow This gives a sequence of trees $\{T_0, T_{d_1}, T_{d_2}, \dots\}$ with increasing $0 \leq d_1 \leq d_2 \leq \dots$

\Rightarrow then there must be an optimal tree T_{d_k} in the sequence $\{T_0, T_{d_1}, T_{d_2}, \dots\}$ with d_k being the smallest one that is $\geq d$.

proof: $C_d(T_{d_k}) = Q(T_{d_k}) + d! T_{d_k}! = Q(T_0) + d_1 + d_2 + \dots + d_k + d!(T_0 - k)! \quad \text{since } Q(T_0) - Q(T_{d_1}) = d_1, Q(T_{d_1}) - Q(T_{d_2}) = d_2, \dots$

as $k \uparrow$, $+d_k$ and $-d \Rightarrow C_d(T_{d_k})$ will start to increase once $d_k > d$ so we stop when $d_k \geq d$

\Rightarrow it's equivalent to growing the strongest terminal node each time, which gives a sequence of optimal trees (T_0, T_1, \dots) with decreasing d \Rightarrow so we choose T_d where d_k is the first one that is $\leq d$.

note: tree size controls the bias-variance tradeoff

- overly large tree \rightarrow complex/overfitted model (low bias and high variance)
- too small tree \rightarrow over-simplify/unfitted model (high bias and low variance).

(2) Classification Trees

$$\hat{Y} = \hat{f}(X) = \sum_{m=1}^M k_m I(X \in R_m) \rightarrow X \text{ is categorized in } R_1, R_2, \dots, R_M.$$

where Y has K levels/classes; k_m is the most commonly occurring class in R_m .

define $\hat{P}_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(Y_i = k)$, $k=1, 2, \dots, K \Rightarrow$ proportion of observation in R_m that belong to level k .

$\Rightarrow k(m)$ can be written as $\arg \max_k \hat{P}_{mk}$. i.e. the most common class in region m .

classification tree uses different loss function to prune/grow trees: e.g.

Misclassification error:

$$Q_m = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i \neq k(m)) = 1 - \max_k \hat{P}_{mk}. \text{ i.e. proportion of observation that are not in the most common class}$$

$$Q_m = 1 - \max(\hat{P}_m, 1 - \hat{P}_m) \text{ for 2 classes situations}$$

Gini index:

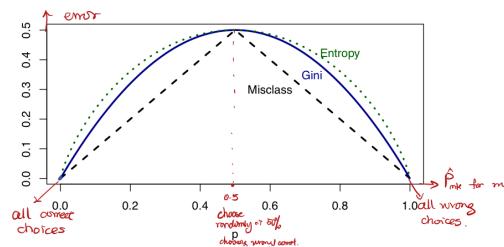
$$Q_m = \sum_{k \neq k'} \hat{P}_{mk} \hat{P}_{mk'} = \sum_{k=1}^K \hat{P}_{mk} (1 - \hat{P}_{mk}) \rightarrow$$
 if \hat{P}_{mk} is close to 1 or 0 then Gini $\rightarrow 0$.

$$Q_m = 2\hat{P}_m(1 - \hat{P}_m) \text{ for 2 classes situations}$$

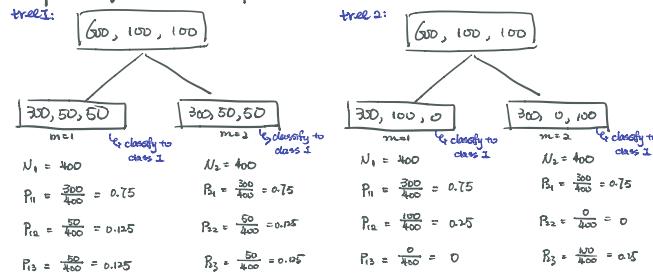
Cross-entropy (binomial deviance):

$$Q_m = - \sum_{k=1}^K \hat{P}_{mk} \log \hat{P}_{mk} \rightarrow$$
 numerically similar to Gini.

$$Q_m = -\hat{P}_m \log \hat{P}_m - (1 - \hat{P}_m) \log (1 - \hat{P}_m)$$



example: given a sample with 600, 100, 100 in 3 classes:



total Misclassification rate:

$$\text{Tree 1: } 400(1-0.75) + 400(1-0.75) = 200$$

$$\text{Tree 2: } 400(1-0.75) + 400(1-0.75) = 200$$

total Gini loss:

$$\text{Tree 1: } 400[(0.75)(0.75) + (0.125)(0.125) + (0.125)(0.125)] \approx 325 \quad \text{Tree 2: } 400[(0.75)(0.75) + (0.25)(0.25) + (0)(0)] \approx 300$$

note: some types misclassification are worse than other types (e.g. false positive or false negative).

\Rightarrow we can adjust loss expression Q_m to upweight some types of losses relative to others.

* we can also use trees for feature selection.

* nice properties.

Speed: Running through the possible splits is fairly fast, and there are efficient algorithms to do this. Trees scale well for large n , compared to other approaches.

* Problems:

In a word - **ACCURACY**.

It turns out that even for large dataset, decision trees give good but not great accuracy. The main reasons are

Interpretability: It is very easy to see what variables are driving the model, to make interpretations and to explain to a less numerate person.

Handles all types of variables well: We do not have to massage ordinal and categorical variables into continuous variables, as with some other approaches. All types of variables can be handled in the single framework.

Treatment of missing variable observations: For a given split, it is (often) possible to find splits using other variables that replicate the main split fairly well. These are referred to as "surrogate splits". If the observation for the main split variable is missing, we can use a surrogate split to assign the appropriate direction to an observation, greatly improving the ability to make predictions when observations are missing.

Robustness: Some modelling approaches will give extreme predictions for new response observations that lie outside the original data domain. By contrast trees give piecewise constant predictions, so will be very stable near boundaries. Similarly the fits are relatively insensitive to outliers (of the response variable).

In sensitivity to predictor transformations: A monotonic transformation of a continuous predictor variable will not affect the tree in any way (since all we do is separate high and low values), so we need not worry about scaling or transforming our data.

Nonparametric: Besides some smoothness, the model does not assume any function form for the underlying model structure. In particular, given enough data a decision tree can approximate the true model structure to arbitrary accuracy.

Curse of dimensionality: Since splits use just one variable at a time and ignore the rest, trees are relatively insensitive to extra dimensionality. Predictions will be almost as good, even if you add extra nonsense variables.

(3) **Bagging** { - many bootstrap samples.
- grow but don't prune (high variance on each tree)
- ↓ variance due to average.
- but trees are correlated

methods:

grow deep, unpruned trees → ↗ regression tree: average the prediction.
→ classification tree: majority vote from all trees.

steps:

- (1) Generate bootstrap sample $\mathcal{B}_b^* = \{(X_{b1}, Y_{b1}), \dots, (X_{bn}, Y_{bn})\}$, for $b=1, 2, \dots, B$
- (2) Use the same prediction method on these bootstrap dataset to give $\hat{f}_1(x), \dots, \hat{f}_B(x)$.
- (3) The bagged model prediction is:

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$
 or majority vote.

(4) **Random Forest** { - many bootstrap samples
- grow but don't prune (high variance on each tree)
- ↓ variance due to average.
- de-correlate trees due to feature bagging.
- feature bagging: at each split consider only m randomly selected predictors → more randomness → double bagging.

steps:

- (1) Generate bootstrap sample $\mathcal{B}_b^* = \{(X_{b1}, Y_{b1}), \dots, (X_{bn}, Y_{bn})\}$, for $b=1, 2, \dots, B$
- (2) For each bootstrap sample \mathcal{B}_b^* , grow a random forest tree T_b by repeating:
 - 2.1 Select m predictors randomly
 - 2.2 Pick the best predictor and split-threshold among the m predictors
 - 2.3 Split the node into two nodes
 - 2.4 Repeat until minimum node size n_{\min} is reached
- (3) Output the ensemble trees $\{\mathcal{T}_b\}_1^B$
- (4) The bagged model prediction is:

$$\hat{f}_{\text{RF}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$$
 for regression tree.

$$\hat{C}_{\text{RF}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$$
 for classification tree
 $m = p/3$
 $m = \sqrt{p}$

(5) **Boosting** { - one original tree
- trees grow sequentially
- each tree is fitted on the residuals
- improve the prediction for the observations that the model is poor at estimating.
- slow-learner: each tree is small.

- ① ► Piecewise constant structure is probably not true in "reality".
- ② ► It struggles to fit main effects.
- ③ ► Model complexity can grow only slowly with sample size (see below). \Rightarrow adding a split would require twice the amount of data.
- ④ ► The greedy procedure does not guarantee optimal splits (see below). \Rightarrow
- ⑤ ► Tree stability can be an issue (see below). \Rightarrow cause a different split every time there's a little change in data.
- ⑥ ► Categorical variables with many levels can distort trees (see below). \Rightarrow the tree stops growing before it has more meaningful splits

Note that the first two points suggest that trees will struggle when linear models will do well. Indeed, the linear model is almost the "chief enemy" of a decision tree.

- we can overcome some of these problems with the following methods:

there are many algorithms but the main difference is the method of training points weighting.
Adaboost:

Methods:

- Produce a collection of classifier $g_m(x_i) \in \{-1, 1\}$ for $m=1, 2, \dots, M$. These are each produced from the training data, re-weighting each time.
- Then the error is :

$$\bar{err} = \frac{1}{n} \sum_i I\{y_i \neq g_m(x_i)\} \text{ for classifier } m.$$

- Then combine these to obtain a weighted majority vote :

$$G(x) = \text{sign}\left(\sum_{m=1}^M d_m g_m(x)\right) \xrightarrow{\text{should be larger for models that are more accurate}}$$



Algorithm:

1. Initialize weights $w_i = 1/n$ for all i
2. for $m=1, \dots, M$
 - a. Fit classifier $g_m(x)$ to the training data, with weights w_i
 - b. Compute weighted error,

$$err_m = \frac{\sum_i w_i I\{y_i \neq g_m(x_i)\}}{\sum_i w_i}$$

- c. Set $d_m = \log((1-err_m)/err_m)$ $\Rightarrow err_m \uparrow \Rightarrow \log(\frac{1-err_m}{err_m}) \uparrow \Rightarrow d_m \uparrow \Rightarrow$ less/more weight for poor/good models
- d. Update w_i with $w_i \exp(d_m I\{y_i \neq g_m(x_i)\})$ for each observation. \Rightarrow increase the weight of observations that the model misclassified

3. Output $G(x) = \text{sign}\left[\sum_m d_m g_m(x)\right]$.

Note:

- Adaboost doesn't specify what classifier to use at each step. Rather it reweights the observations each time \Rightarrow so we can use different models at each stage, such as decision trees, simple spline models etc.
- Adaboost can take a class of weak classifiers and produce an ensemble model that predicts accurately

Forward Stagewise Modelling

Methods:

Suppose $f(x) = \sum_{m=1}^M \beta_m b_m(x; \gamma_m)$, where β_m and γ_m are determined to minimize :

$$\sum_i L(Y_i, \sum_{m=1}^M \beta_m b_m(x_i; \gamma_m))$$

forward stagewise adds a basic function one at a time.

Algorithm:

1. Initialise $f_0(x) = 0$
2. For $m=1, \dots, M$,
 - a. Compute $(\beta_m, \gamma_m) = \underset{\beta, \gamma}{\operatorname{arg\min}} \sum_i L(Y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma))$.
 - b. Set $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$.
3. $f(x) = \sum_{m=1}^M \beta_m f_m(x)$

Note:

- If γ represent the choice of the number of splits and level for each terminal node, Adaboost is equivalent to (i.e. $\gamma_m = 2\beta_m$) forward stagewise modelling with exponential loss function:

$$L(Y_i, f(x_i)) = \exp(-Y_i f(x_i)) = \begin{cases} \exp(-1) & \text{if misclassified} \\ \exp(-1) & \text{if correct.} \end{cases}$$

- So we can :

- build stagewise models for both classification and regression.
- choose any loss function we like
- Not be restricted to classifiers taking values in $\{-1, 1\}$.
- modify the basis on which new terms are added to the model (e.g. shrinkage).



For regression:

$$\{Y_i - f(x_i)\}^2 \text{ (square loss),}$$

$$|Y_i - f(x_i)| \text{ (absolute loss),}$$

$$h(Y_i - f(x_i)) \text{, where } h(u) = \begin{cases} u^2 & \text{if } |u| \leq \delta \\ 2\delta(|u| - \delta/2) & \text{if } |u| > \delta \end{cases} \text{ (Huber loss)}$$

For classification:

$$I[Y_i \neq \text{sign}(f(x_i))] \text{ (misclassification rate).}$$

$$-\left[I(Y_i = 1) \log \left(\frac{e^{f(x_i)}}{1 + e^{f(x_i)}} \right) + I(Y_i = -1) \log \left(\frac{1}{1 + e^{f(x_i)}} \right) \right] \text{ (binomial loss)}$$

$$\exp(-Y_i f(x_i)) \text{ (exponential loss)}$$

$$\{Y_i - \frac{e^{f(x_i)}}{1 + e^{f(x_i)}}\}^2 \text{ (squared loss)}$$

boosting with forward stagewise.

1. Initialise $f_0(x) = 0$

2. For $m=1, \dots, M$,

2.1 Compute $(\beta_m, r_m) = \arg \min_{\beta, r} \frac{1}{n} \sum_{i=1}^n L(Y_i, f_{m-1}(x_i) + \beta b_m(x_i; r))$.

2.2 Set $f_m(x) = f_{m-1}(x) + \beta_m b_m(x; r_m)$. \rightarrow can set new basic function each time.

3. $f(x) = \sum_{m=1}^M \beta_m f_m(x)$