

Random Forest

- uses bagging on decorrelated decision trees for regression or classification.
- decision trees are low bias + high variance models and can thus benefit from bagging (moreover, they are quick to train)

- an ensemble method

bagging

- ① bootstrap a sample of the original ^{dataset} size from the data
- ② train a model, M_1
- ③ repeat for N models
- ④ regression: avg. results ; classification: take majority vote

$$\begin{aligned}\text{Var}[F(x)] &= \text{Var}\left[\frac{1}{N} \sum_{i=1}^N f_i(x)\right] = \frac{1}{N^2} \left\{ \sum_i \overbrace{\text{Var}[f_i(x)]}^{\equiv \sigma^2} + 2 \sum_{i < j} \overbrace{\text{corr}(f_i(x), f_j(x))}^{\equiv \rho} \sqrt{\text{Var}[f_i(x)]} \sqrt{\text{Var}[f_j(x)]} \right\} \\ &= \frac{1}{N^2} \left\{ N\sigma^2 + 2\binom{N}{2} \rho \sigma^2 \right\} = \frac{\sigma^2}{N} + \frac{N-1}{N} \rho \sigma^2\end{aligned}$$

⇒ increasing N decreases variance to a limiting value of $\approx \rho \sigma^2$ ^(0.5) better than the variance of a single model σ^2

bias remains unchanged:

$$\begin{aligned}\text{Bias}[f(x)] &= \text{Bias}\left[\frac{1}{N} \sum_i f_i(x)\right] = E\left[\frac{1}{N} \sum_i f_i(x)\right] - f_{\text{act}}(x) \\ &= E[f_1(x)] - f_{\text{act}}(x) \\ &= \text{Bias}[f_1(x)]\end{aligned}$$

building a decision tree

- algo is a greedy process where we intend to reduce squared loss (for regression) and 0-1 loss for misclassification
 - recursively partition the feature space w/ binary splits on a certain feature, at a certain value
 - the feature, and feature value are found by greedily minimizing squared loss ($\hat{y}_R = \bar{y}_R$) for classification and gini index / cross-entropy for classification
 - gini index or cross entropy tend to result in more "pure" nodes rather than the misclassification rate, which would correspond to minimizing 0-1 loss

- stop building tree at a certain stopping criterion (eg. # data points in leaves)

RF Algo

- repeat N times
- ① Bootstrap sample
 - ② build decision tree where you only consider a randomly chosen subset^{of features} at each node along which to split. This reduces correlation amongst the trees (introduces randomness, & avoids all trees splitting on very predictive features), thus reducing variance.
 - ③ Average preds in the end (regression), take majority vote (classification)

OOB estimate

- each bootstrap contains $\approx \frac{2}{3}|D|$, so that for each data point, z_i , it is not used to train $\approx \frac{1}{3}$ of the trees
- let $y_{\text{OOB}}^{(i)}$ = mean prediction from $x^{(i)}$ over all trees not trained w/ z_i .
- $MSE_{\text{OOB}} = \sum_{i=1}^n (y^{(i)} - y_{\text{OOB}}^{(i)})^2$

← m_{OOB}

and technically be $< n$ if there are data points used in training every tree

Variable importance

tree t

$$FI_j^t = \frac{\sum_{s: j \text{ is the splitting feature}} (\Delta RSS)_s}{RSS_j^t - RSS_o^t} \Rightarrow FI_j = \langle FI_j^t \rangle_t$$

Feature j

- Feature importance is roughly the decrease in the RSS for Feature j for tree t , divided by total decrease in RSS (then averaged over all trees)
- Similarly, use $\Delta(\text{quality measure})$ for categoricals

Pros

- typically great performance
- easily parallelizable (just train each tree on a different core)
- difficult to over-fit
- sometimes slow

Cons

- not a regularized model (can lead to issues w/ high dims)
 - can fix w/ feature selection
- difficult to interpret