

Introduction to Machine Learning

Random Forest Basics



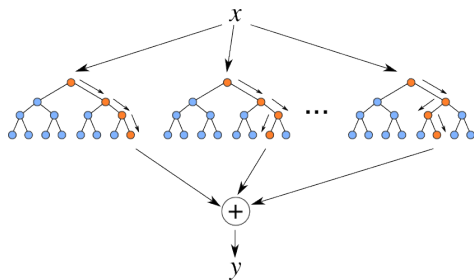
Learning goals

- Know how random forests are defined by extending the idea of bagging
- Understand general idea to decorrelate trees
- Understand effects of hyperparameters
- RFs and overfitting

RANDOM FORESTS

► Breiman 2001

- RFs use **bagging with CARTs as BLs**
- **Random feature sampling** decorrelates the base learners
- **Fully expanded trees** further increase variability of trees



INTUITION BEHIND DECORRELATION

- Since bootstrap samples are similar, models $\hat{b}^{[m]}$ are correlated, affecting the variance of an ensemble \hat{f}
- We would like variance to go down linearly with ensemble size, but because of correlation we cannot really expect that
- Assuming $\text{Var}(\hat{b}^{[m]}) = \sigma^2$, $\text{Corr}(\hat{b}^{[m]}, \hat{b}^{[j]}) = \rho$, semi-formal analysis, without proper analysis of prediction error:

$$\begin{aligned}\text{Var}(\hat{f}) &= \text{Var}\left(\frac{1}{M} \sum_{m=1}^M \hat{b}^{[m]}\right) = \frac{1}{M^2} \left(\sum_{m=1}^M \text{Var}(\hat{b}^{[m]}) + 2 \sum_{m < j} \text{Cov}(\hat{b}^{[m]}, \hat{b}^{[j]}) \right) \\ &= \frac{1}{M^2} \left(M\sigma^2 + 2 \frac{M(M-1)}{2} \rho\sigma^2 \right) = (1 - \rho) \frac{\sigma^2}{M} + \rho\sigma^2\end{aligned}$$

- Ensemble variance is “convex-combo of linear-reduction and no-reduction, controlled by ρ ”
- Maybe we can decorrelate trees, to reduce ensemble variance?
And get less prediction error?



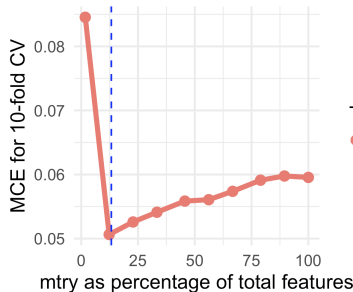
RANDOM FEATURE SAMPLING

RFs decorrelate trees with a simple randomization:

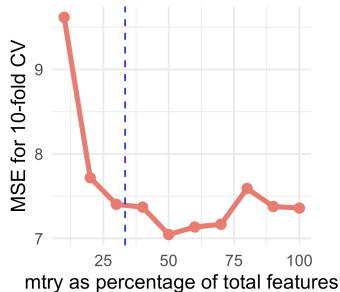
- For each node of tree, randomly draw $m_{\text{try}} \leq p$ features (m_{try} = name in some implementations)
- Only consider these features for finding the best split
- Careful: Our previous analysis was simplified! The more we decorrelate by this, the more random the trees become!
This also has negative effects!



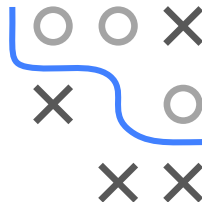
EFFECT OF FEATURE SAMPLING



Task
● spam



Task
● mtcars



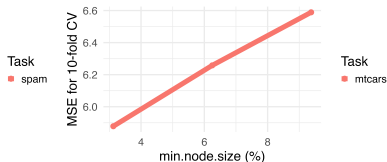
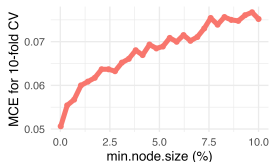
- Optimal `mtry` typically larger for regression than for classification
- Good defaults exist, but still most relevant tuning param
- Rule of thumb:
 - Classification: $mtry = \lfloor \sqrt{p} \rfloor$
 - Regression: $mtry = \lfloor p/3 \rfloor$

TREE SIZE

In addition to `mtry`, RFs have two other important HPs:

- Min. nr. of obs. in each decision tree node

Default (ranger): `min.node.size = 5` ▶ Breiman 2001

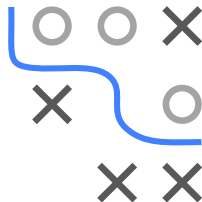


- Depth of each tree

Default (ranger): `maxDepth = ∞`

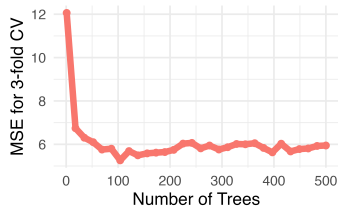
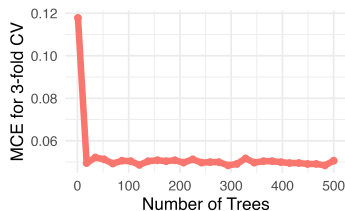
- There are more alternative HPs to control depth of tree: minimal risk reduction, size of terminal nodes, etc.

RF usually use fully expanded trees, without aggressive early stopping or pruning, to further **increase variability of each tree**. ▶ Louppe 2015

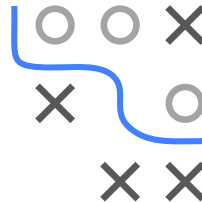
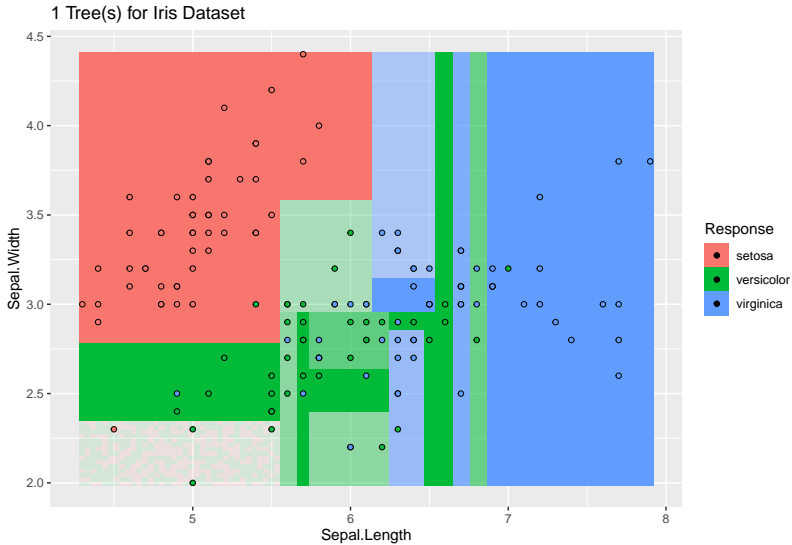


ENSEMBLE SIZE

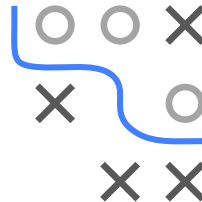
- RFs usually better if ensemble is large ▶ Breiman 2001
- But: Increases computational costs, and diminishing returns
- 100 or 500 is a sensible default
- Can also inspect the OOB error (see later)



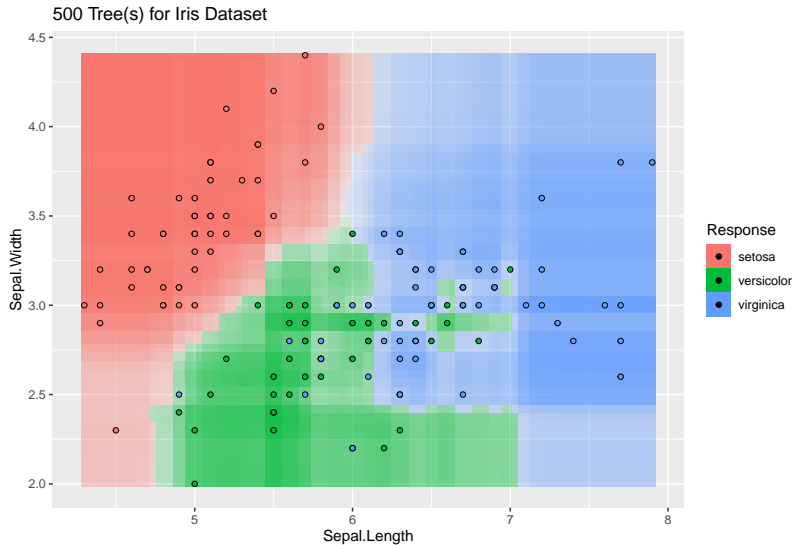
EFFECT OF ENSEMBLE SIZE



This plot displays the result of a 10-tree ensemble model for the Iris dataset. The x-axis represents 'Sepal.Length' (ranging from approximately 4.5 to 8.0) and the y-axis represents 'Sepal.Width' (ranging from 2.0 to 4.5). The background is divided into three distinct color-coded regions: red (top-left), green (bottom-left to center), and blue (center-right to top-right). These regions represent the predicted class for each point in the feature space based on the ensemble of 10 trees. Numerous individual data points are overlaid as small circles, showing their distribution across the plot. The decision boundary between the red and green regions is particularly complex and irregular, while the boundary between the green and blue regions is more linear.



EFFECT OF ENSEMBLE SIZE



► Probst and Boulesteix 2018

-



©

RF IN PRACTICE

Benchmarking bagged ensembles with 100 BLs each on spam
versus RF ($\text{ntrees} = 100$, $\text{mtry} = \sqrt{p}$, $\text{minnode} = 1$),
we see how well RF performs!

