

Motivation

- One of the successes of modern statistics and machine learning is that we have a wide range of models that work well on high-dimensional data. Out of the box, these models can build predictive rules from a large number of features – no manual curation needed. However, unless the model performs explicit variable selection, this automation can leave the model’s developers and users at a loss for which features were actually relevant.
- Variable importance measures attempt to answer this question. Most approaches rely on the intuition that a feature is important if “removing it” causes model performance to deteriorate. Of course, the devil is in the details, and there are many ways to formalize this intuition.
- This intuition might remind you of attribution using Shapley values, but there are two key (important?) differences. First, Shapley values are focused on local explanations, while most (though not all) variable importance attempts to characterize importance globally. Second, Shapley values focus on the estimated model, while variable importance focuses on the estimated model’s performance.

Variable Importance in Random Forests

- The phrase “variable importance” is closely associated with random forests, and for good reason – the random forest paper included not just a method for estimating the model but also a method for estimating variable importance. The important variables were reported throughout the paper’s illustrations, and its particular implementation has been influential.
- Let’s briefly review the original proposal before exploring some of its issues (and the follow up work that addresses them). Random forests work by averaging the predictions from many different [decision trees](#). To make the decision tree predictions less correlated with one another, each is trained on a different subsample of both training data and features.
- To compute the importance of feature d within tree t , we compare the out-of-sample prediction accuracy when using the true test data and a version of the test data with feature d randomly permuted:

$$VI_{t,d} = \frac{1}{n} \left[\sum_{i=1}^n J(y_i, x_i^{(-d)}) - J(y_i, x_i^{(-)}) \right]$$
$$x_i^{(-)} = (x_{i1}, \dots, x_{i(d-1)}, \dots, x_{in})$$

- The importance of that feature in the RF averages importances across all trees: Despite shuffling our initial intuition so well, we have to be extremely careful with how we interpret these importances. The absurdly simple table below comes from Elton’s Prediction, Estimation, and Attribution [paper](#).

Table B: Results of test set scores for gradient-boosted random forest predictions, measuring top prediction errors in Figure 7.

Feature	0	1	2	3	4	5	6	7	8	9	100	300
Mean	1	0	3	1	2	2	2	0				

There are 6332 total features in the original random forest. Replacing the 348 “most important” of them leads to no perceptible decrease in test set performance!

- The trouble is that real data have correlated features. This correlation wreaks havoc on the usual intuition about variable importance. For this reason, a running theme in the literature since the original RF paper is: We need to be precise about what exactly we are estimating.
- Or, in the case of Strobl et al. (2008), can we be precise about what we are testing? The permutation step in the VI calculation introduces a test of conditional independence. To see this, first note that if the response were independent of the features, then

$$P(y, X) = P(y)P(X)$$

If we define a test statistic that is large when this factorization doesn’t hold, then we will be able to test

$$H_0: y \perp\!\!\!\perp X$$

- For example, our test statistic could be R^2 . If R^2 is no higher on the real data than it is on datasets with the Y permuted, then we don’t have any evidence that X is associated with Y . For random forest variable importance, we don’t permute the response – we permute individual features. If we had

$$P(y, X_d, X_{-d}) = P(X_d)P(y, X_{-d})$$

then the two sides of the VI definition would be comparable. Only when this independence is violated will a variable be flagged as independent. So RF’s variable importance measure is implicitly testing

$$H_0: X_d \perp\!\!\!\perp y, X_{-d}$$

The issue here is that if $X[d]$ is correlated with other features, then that will also lead to large variable importance.

- Strobl et al. (2008) suggest to instead develop a permutation scheme that tests

$$H_0: X_d \perp\!\!\!\perp y \mid X_{-d}$$

Feature d has to be associated with y even after controlling for the remaining features. This is accomplished by separately permuting $X[d]$ within neighborhoods

Y	X_d	Z
y_1	$x_{1,d}$	$z_1 = a$
\vdots	\vdots	\vdots
y_2	$x_{2,d}$	$z_2 = a$
\vdots	\vdots	\vdots
y_3	$x_{3,d}$	$z_3 = b$
\vdots	\vdots	\vdots
y_4	$x_{4,d}$	$z_4 = b$
\vdots	\vdots	\vdots
y_5	$x_{5,d}$	$z_5 = b$
\vdots	\vdots	\vdots
y_6	$x_{6,d}$	$z_6 = b$
\vdots	\vdots	\vdots
y_7	$x_{7,d}$	$z_7 = a$
\vdots	\vdots	\vdots
y_8	$x_{8,d}$	$z_8 = a$
\vdots	\vdots	\vdots
y_9	$x_{9,d}$	$z_9 = a$
\vdots	\vdots	\vdots
y_{10}	$x_{10,d}$	$z_{10} = a$

If X is continuous, then we can’t require exact matches. Instead, the proposal defines a grid using the decision tree’s original partition. Within each of those grid elements, we can permute the d th feature



A further relaxation is to ignore grid dimensions for features $X[d]$ that are known not to be very correlated with $X[d]$. This is a bit heuristic, but seems accurate enough in practice, and it does help with ensuring there are more samples to permute.

MDI

- We actually didn’t need to use permutations to have a notion of variable importance in trees. If we consider how trees are built, we’ll see that there’s already a notion of importance built into the construction process. Specifically, when a decision tree is grown, it repeatedly computes the “decrease in impurity” associated with defining a new split at an existing tree node t :

$$\Delta i(t) = \frac{1}{N(t)} \left[\sum_{x \in L} (y_x - \bar{y}_L)^2 - \sum_{x \in R} (y_x - \bar{y}_R)^2 \right]$$
$$= \frac{1}{N(t)} \left[\sum_{x \in L} y_x^2 - 2 \bar{y}_L \sum_{x \in L} y_x + N(t) \bar{y}_L^2 - \sum_{x \in R} y_x^2 + 2 \bar{y}_R \sum_{x \in R} y_x - N(t) \bar{y}_R^2 \right]$$

L and R are the left and right descendant subtrees obtained by allowing a split along d . $N(t)$ is the number of samples that descend from the current tree node.

- Q. What tree does the drawing above correspond to? How would you modify the drawing to reflect a second split in the left node?

- The mean decrease in impurity (MDI) for feature d is the weighted mean of these decreases for all splits that involve feature d .

$$MDI_d = \frac{1}{N} \sum_{t \in \mathcal{T}} N(t) \Delta i(t)$$

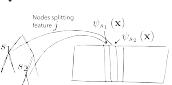
- From the drawings above, you might already notice that decision trees can be represented as linear regressions onto a basis that’s derived from the initial tree estimate. It’s not particularly useful for learning the tree (because the basis depends on the tree), but it’s useful for understanding properties of the model, including variable importances.

$$f(x) = \frac{1}{\sqrt{N(t)}} \left[N(t) \mathbb{1}\{x \in L\} - N(t) \mathbb{1}\{x \in R\} \right]$$



- The fitted values from this regression are exactly the same as the decision tree fits.

$$\hat{g} = [\psi_1(x) \mid \psi_2(x) \mid \dots \mid \psi_d(x)]^T \hat{\beta}$$



Further, consider replacing all basis elements that don’t involve feature d with their averages:

$$\hat{g} = [\bar{\psi}_1(x) \mid \dots \mid \bar{\psi}_{d-1}(x) \mid \psi_d(x) \mid \dots \mid \bar{\psi}_D(x)]^T \hat{\beta}$$

The R^2 of this model is exactly the MDI for feature d !

- The essential insight of Agarwal et al. (2023) is that we don’t have to restrict ourselves to linear regression and R^2 . We can use any high-dimensional regression model and any measure of model performance. They call these new measures of feature d ’s importance MDIs.

- For example, using ridge regression (still with R^2) leads to more stable variable rankings than ordinary MDI. Further, if we construct regressions like these across many trees in an ensemble, we can define a new type of random forest (they call it RF 2) with analogous feature importance measures.

Importance, Retaining, and Approximations

- How might we generalize these ideas beyond trees? Suppose that the data is drawn according to $\pi = P$ and that we want to analyze $\mathcal{Y}(x) = E[Y \mid X = x]$. Also, consider these variants that either ignore or replace feature d :

$$x_{-d} \sim \mathcal{P}_{-d}$$
$$x_{-d} = (x_1, \dots, x_{d-1}, p_d, x_{d+1}, \dots, x_D) \sim \mathcal{P}^{(d)}$$
$$x_{-d}(x^{(d)}) = E[y \mid x_{-d}]$$

(Remove)

(Replace)

(Ignore: replace on left and remove from right)

These are defined in Gao et al. (2023). We will also use the notation $\mathcal{P}^{(d)}[d]$ to denote the distribution of $\mathcal{P}^{(d)}[d]$, though this is not defined in that paper.

- Then, if V is a measure of model performance, we can define the importance of feature d as:

$$VI_d = V(\hat{f}, \hat{\mathcal{P}}) - V(\hat{f}_{-d}, \hat{\mathcal{P}}_{-d})$$

In other words, does the model performance deteriorate when we ignore feature d ?

- In reality, we know none of the distributions and conditional distributions needed to compute this variable importance measure. But at least now we can be precise about what to estimate. There are two common ways to approach the problem, and a third interesting compromise introduced by Gao et al. (2023).

- The dropout approach first estimates the performance of a model $V(\text{full})$ trained using the entire dataset. Then, it gets new predictions by querying the function on versions of the data with the d th coordinate replaced by its mean. The two performances are compared. In other words:

$$\widehat{VI}_d = V(\hat{f}, \hat{\mathcal{P}}) - V(\hat{f}_{-d}, \hat{\mathcal{P}}^{(d)})$$

This approach is very fast, but it can be misleading (Prop. 3.2) (Like in the original variable importance definition for RFs, it can overstate the importance of a variable that is correlated with an important feature, even if conditionally it has no influence on the response).

- Intuitively, suppose features 1 and 2 are almost perfectly correlated, but that only feature 1 is actually used to generate the response. As long as $\text{Var}(\text{beta}[1]) + \text{Var}(\text{beta}[2])$ are close to the true $\text{beta}[1]$, then the full model will have good predictions, but the dropout version will perform poorly (especially if $\text{Var}(\text{beta}[2])$ is large). This performance gap leads to high variable importance. If, however, d has a chance to interact with feature 2 removed, then the estimated $\text{Var}(\text{beta}[1])$ will be close to the true $\text{beta}[1]$, and we will have no drop in performance.

- Alternatively, in the second step, we could have retained a new model using a version of the dataset that replaces coordinate d with its mean,

$$\hat{f}_{-d} = \arg \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i^{(-d)}))^2$$

and then evaluate the performance of this new model, which never had access to the real feature d . We can express these steps as:

$$\widehat{VI}_d^{\text{RT}} = V(\hat{f}, \hat{\mathcal{P}}) - V(\hat{f}_{-d}, \hat{\mathcal{P}}^{(-d)})$$

This is called the *retrofit* approach. It doesn’t suffer from the correlation issue of dropout, but it requires retraining the model d times, which can be intractable.



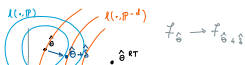
- The idea of Gao et al. (2023) is to make a one-step update in the direction that full retaining would have taken us. The update will be much cheaper than full retaining, but it will also address the basic issues with dropout. Specifically, if we have a parameterized model, then the model trained using all features corresponds to the estimator:

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N (y_i - \theta(x_i))^2$$

The lazy estimation approach is to find one-step updates:

$$\hat{\delta}_d = \arg \min_{\delta} \frac{1}{N} \sum_{i=1}^N (y_i - \theta(x_i^{(-d)}) - \delta \psi_d(x_i^{(-d)}) + \lambda \|\delta\|_1)^2$$

Notice that the first function evaluation in this expression computes the dropout prediction. Instead of using the dropout prediction directly, though, we use the corrected version.



Comparing the full and one-step models leads to a new estimate of variable importance:

$$VI_d^{\text{LA}} = V(\hat{f}, \hat{\mathcal{P}}) - V(\hat{f}_{\hat{\delta}_d}, \hat{\mathcal{P}}^{(-d)})$$

- This approach works as an especially nice interpretation in the case that f is a neural network. Estimating the one-step update is like fitting a neural tangent kernel (NTK) model, a kind of linearized neural network model. The only difference from NTK is that we use the full model’s parameter as our initialization, not a random initialization.

- The main result of Gao et al. (2023) is a central limit theorem for the estimator $VI^{\text{LA}}[\delta][d]$. By adapting concentration bounds from the NTK literature, they show that these estimators are asymptotically normal with mean $VI[d]$.

$$\sqrt{n} (\widehat{VI}_d^{\text{LA}} - VI_d) \xrightarrow{d} N(0, \sigma_d^2)$$

Moreover, a plug-in estimate of the associated covariance yields confidence intervals for the variable importance measure. Of course, the result requires conditions on the data and model class. They are surprisingly simple (only three), but we nonetheless defer these details to the reading.

- A side note: Lazy estimation has a more heuristic but general analog. Instead of solving the optimization above, we can imagine taking the full model and then “fine-tuning” it on the dropout-transformed data. Taking a few steps is much cheaper than training to convergence, but it can already improve on naive dropout.

Shapley-Based Importance

- We now have a nice way of estimating $VI[d]$. But why should we be targeting this quantity in the first place? It has intuitive appeal, but certainly doesn’t follow from any first principles. To address this, we can try to connect variable importance to Shapley values. In the same way that ordinary Shapley attribution averages function evaluations over arbitrary subsets of predictors, we will need a version of our conditional expectation that ignores more than just the d th feature at a time. Specifically, for a subset of features S , define:

$$x_S^{(-)} = (x_1, \dots, x_{d-1})$$
$$x_S(x^*) = E[y \mid x_S]$$

$\mathcal{P}[S]$ is a submodel that can only use information from features in S .

- With this notation, we can write the Shapley-based importance introduced by Williamson and Feng (2020):

$$SVI_{d,S} = \sum_{S \in \mathcal{B}(d)} \frac{1}{d} \frac{1}{\binom{p-1}{d-1}} \left[V(x_{S \cup d}, \mathcal{P}) - V(x_S, \mathcal{P}) \right]$$

Like the Shapley values we covered in week 2, this definition satisfies a number of nice properties. For example, it satisfies an additivity property:

$$\frac{\partial}{\partial x_d} SVI_{d,S} = V(x_S, \mathcal{P}) - V(x_{\emptyset}, \mathcal{P})$$

Individual feature importances gap between full and null models