

## Prediction Error Reduction Function as a Variable Importance Score

E. Fokoué<sup>1\*</sup>

<sup>1</sup>*School of Mathematical Sciences, Rochester Institute of Technology, 98 Lomb Memorial Drive, Rochester, NY 14623, USA.*

### *Author's contribution*

*The sole author designed, analyzed and interpreted and prepared the manuscript.*

### *Article Information*

DOI: 10.9734/BJMCS/2016/23945

#### *Editor(s):*

(1) Carlo Bianca, Laboratoire de Physique Thorique de la Matire Condense, Sorbonne Universits, France.

#### *Reviewers:*

(1) Er. Anuj Goel, Maharishi Markandeshwar University, India.

(2) S. K. Srivatsa, Prathyusha Engg College, Chennai, India.

Complete Peer review History: <http://sciencedomain.org/review-history/13781>

*Received: 29<sup>th</sup> December 2015*

*Accepted: 27<sup>th</sup> February 2016*

*Published: 21<sup>st</sup> March 2016*

**Original Research Article**

## Abstract

This paper introduces and develops a novel variable importance score function in the context of ensemble learning, and demonstrates its appeal empirically. Our proposed score function is simple and more straightforward than its counterpart proposed in the context of random forest, and by avoiding permutations, it is by design computationally more efficient than the random forest variable importance function. Just like the random forest variable importance function, our score handles both regression and classification seamlessly. One of the distinct advantage of our proposed score is the fact that it offers a natural cut off at zero, with all the positive scores indicating importance and significance, while the negative scores are deemed indications of insignificance. An extra advantage of our proposed score lies in the fact it works very well beyond ensemble of trees and can seamlessly be used with any base learners in the random subspace learning context. Our examples, both simulated and real, demonstrate that our proposed score does compete mostly favorably with the random forest score.

*\*Corresponding author: E-mail: [epfeqa@rit.edu](mailto:epfeqa@rit.edu);*

*Keywords:* High-dimensional; variable importance; random subspace learning; out-of-bag error; random forest; large  $p$  small  $n$ ; classification; regression; ensemble learning; base learner.

**2010 Mathematics Subject Classification:** 62H30, 62H25.

## 1 Introduction

Consider a data set  $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$  where  $\mathbf{x}_i$  is a  $p$ -dimensional vector of attributes of potentially different types observable on some input space denoted here by  $\mathcal{X}$ , and  $\mathbf{y}_i$  are the responses taken from  $\mathcal{Y}$ . We shall consider various scenarios, but mainly the regression scenario with  $\mathcal{Y} = \mathbb{R}$  and the classification scenario with  $\mathcal{Y} = \{1, 2, \dots, K\}$ . We consider the task of building the estimator  $\hat{f}(\cdot)$  of the true but unknown underlying  $f$ , and seek to build  $\hat{f}(\cdot)$  such that the true error (generalization error) is as small as possible. In this context, we shall use the average test error  $\text{AVTE}(\cdot)$ , as our measure of predictive performance, namely

$$\text{AVTE}(\hat{f}) = \frac{1}{R} \sum_{r=1}^R \left\{ \frac{1}{m} \sum_{j=1}^m \ell(\mathbf{y}_j^{(r)}, \hat{f}^{(r)}(\mathbf{x}_j^{(r)})) \right\}, \quad (1.1)$$

where  $(\mathbf{x}_j^{(r)}, \mathbf{y}_j^{(r)})$  is the  $j$ th observation from the test set at the  $r$ th random replication of the split of the data. Throughout this paper, we shall use the zero-one loss (1.2) for all our classification tasks.

$$\ell(\mathbf{y}_j^{(r)}, \hat{f}^{(r)}(\mathbf{x}_j^{(r)})) = 1_{\{\mathbf{y}_j^{(r)} \neq \hat{f}^{(r)}(\mathbf{x}_j^{(r)})\}} = \begin{cases} 1 & \text{if } \mathbf{y}_j^{(r)} \neq \hat{f}^{(r)}(\mathbf{x}_j^{(r)}) \\ 0 & \text{otherwise.} \end{cases} \quad (1.2)$$

For regression tasks, we shall use the squared error loss (1.2), namely

$$\ell(\mathbf{y}_j^{(r)}, \hat{f}^{(r)}(\mathbf{x}_j^{(r)})) = (\mathbf{y}_j^{(r)} - \hat{f}^{(r)}(\mathbf{x}_j^{(r)}))^2. \quad (1.3)$$

Besides seeking the optimal predictive estimator of  $f$ , we also seek to select the most important (useful) predictor variables as a byproduct of our overall learning scheme. Indeed, while accurate prediction is very important in and of itself, it's often desirable or even crucial in some cases, provide the added description of the importance of the variables involved in the prediction task. The statistical literature is filled with thousands of papers on variable selection and measurement of variable importance. [1] and [2] propose a measure of variable importance for random forest. [3] and [4] have written many interesting contributions to the estimation of variable importance specifically in the context of classification and regression trees and their random forest ensemble learning extension. [5] implements some of those measures of variable importance for tree-based learning in their wonderful R package *caret*. [6] give a nice summary of variable importance and [7] also touches on this subject. [8] proposed a bias-correction improvement to the above measures of variable importance. Unlike all these authors whose work on variable importance measures concentrated solely on tree-based models, we herein propose a measure that goes beyond tree-based methods.

## 2 Construction of the Prediction Error Reduction Function

### 2.1 Definitions and tools for defining the score function

We consider the common framework of a  $p$ -dimensional input space  $\mathcal{X}$  with typical input vector  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_p)^\top$ . We also consider building different models with different subsets of the  $p$  original

variables. Let  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)^\top$  denote the  $p$ -dimensional indicator such that

$$\gamma_j = \begin{cases} 1 & \text{if } \mathbf{x}_j \text{ is active in the current model indexed by } \boldsymbol{\gamma} \\ 0 & \text{otherwise.} \end{cases} \quad (2.1)$$

Assume that we are given an ensemble (collection or aggregation) of models, say

$$\mathcal{H} = \{\mathbf{g}(\cdot, \boldsymbol{\gamma}^{(1)}), \mathbf{g}(\cdot, \boldsymbol{\gamma}^{(2)}), \dots, \mathbf{g}(\cdot, \boldsymbol{\gamma}^{(B)})\} \quad (2.2)$$

where  $\mathbf{g}(\cdot, \boldsymbol{\gamma}^{(b)})$  denotes the function built with only those variables that are active in the  $b$ th model of the ensemble (aggregation), and  $\boldsymbol{\gamma}^{(b)} = (\gamma_1^{(b)}, \dots, \gamma_p^{(b)})$  with

$$\gamma_j^{(b)} = \begin{cases} 1 & \text{if } \mathbf{x}_j \text{ is active in the } b\text{-th model of the ensemble} \\ 0 & \text{otherwise.} \end{cases} \quad (2.3)$$

For instance, we may consider a homogeneous ensemble, i.e, an ensemble in which all the functions are of the same family, like the case where all the base learners are multiple linear regression (MLR) models differing by the variables upon which they are built. Consider a score function  $\text{score}(\mathbf{g}(\cdot, \boldsymbol{\gamma}^{(b)}))$  used to assess the performance of model indexed by the variables active in  $\boldsymbol{\gamma}^{(b)}$ . We propose a variable importance score in the form of a function that measures the importance of a variable  $\mathbf{x}_j$  in terms of the reduction in average score

$$\text{PERF}(\mathbf{x}_j) = \frac{1}{B} \sum_{b=1}^B \text{score}(\mathbf{g}(\cdot, \boldsymbol{\gamma}^{(b)})) - \frac{1}{B_j} \sum_{b=1}^B \gamma_j^{(b)} \text{score}(\mathbf{g}(\cdot, \boldsymbol{\gamma}^{(b)})) \quad (2.4)$$

where  $B_j$  is the number of models containing the variable  $\mathbf{x}_j$ , specifically  $B_j = \sum_{b=1}^B 1_{\{\gamma_j^{(b)}=1\}}$ . In words,

$$\text{PERF}(\mathbf{x}_j) = \text{Average score over all models} - \text{Average score over all models with } \mathbf{x}_j$$

## 2.2 Properties and benefits of the PERF score function

Intuitively,  $\text{PERF}(\mathbf{x}_j)$  somewhat measures the impact of variable  $\mathbf{x}_j$ . In a way similar to the approach used by sports writers to decide the Most Valuable Player (MVP) on a team or in a league,  $\text{PERF}(\mathbf{x}_j)$  looks at the overall performance of the whole ensemble and then for each variable  $\mathbf{x}_j$  computes the direction and magnitude of the change to that overall performance of the ensemble brought by its presence in models. In a sense, the variable with the highest PERF score is the Most Valuable Predictor (MVP) variable. In other words, if

$$j^* = \underset{j=1, \dots, p}{\text{argmax}} \{\text{PERF}(\mathbf{x}_j)\}, \quad \text{then } \mathbf{x}_{j^*} = \text{MVP}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_p).$$

From the above definition, it follows that *If a variable  $\mathbf{x}_j$  is important, then its presence in any model will cause that model to perform better in the sense of having a lower than common average error (score). The average score of all models containing an important variable  $\mathbf{x}_j$  should therefore be lower than the overall average score.* Essentially, (i)  $|\text{PERF}(\mathbf{x}_j)|$  measures the magnitude of the importance/impact. (ii)  $\text{sign}(\text{PERF}(\mathbf{x}_j))$  measures the direction of the impact. (iii) If  $\text{sign}(\text{PERF}(\mathbf{x}_j)) = +1$  and  $|\text{PERF}(\mathbf{x}_j)|$  is relatively large, then  $\mathbf{x}_j$  is an important variable. Some of the benefits of the PERF score include the following: (a) The PERF score is seamlessly applied to both large  $p$  small  $n$  and small  $p$  large  $n$  machine learning settings, whether it be a classification or a regression task. (b) All variables with  $\text{PERF}(\mathbf{x}_j) \leq 0$  are unimportant and can be discarded. (c) The  $\text{PERF}(\cdot)$  score can be used whenever an ensemble  $\mathcal{H}$  is available along with a suitable score function for each base learner. (d) This works with any base learner and can be adapted to parametric, nonparametric and semi-parametric models and one can imagine ensembles with any base learners as its atoms. (e) A great advantage over the traditional variable importance [1], [2] score functions is that the clear cut-off at zero, in the sense that all variables with  $\text{PERF}(\mathbf{x}_j) > 0$  are kept and all those variables with  $\text{PERF}(\mathbf{x}_j) \leq 0$  are thrown away.

### 2.3 Random subspace learning estimation of PERF

A natural implementation of  $\text{PERF}(\cdot)$  can be done using the ubiquitous bootstrap along with the random subspace learning scheme. The **Out-of-Bag** (oob) error in the bagging or random subspace learning context is a good (in fact excellent) candidate score function, especially when the goal is the selection of variables that lead to the lowest prediction error. The advantage of using oob as the score lies in the fact that the score is obtained as part of building the ensemble in the random subspace learning framework. Consider the training set  $\mathcal{D} = \{\mathbf{z}_i = (\mathbf{x}_i^\top, y_i)^\top, i = 1, \dots, n\}$ , where  $\mathbf{x}_i^\top = (x_{i1}, \dots, x_{ip})$  and  $y_i \in \mathcal{Y}$  are realizations of two random variables  $X$  and  $Y$  respectively. Let  $\mathbf{x}_{i,\pi_j} = (x_{i,1}, \dots, x_{i,\pi_j}, \dots, x_{i,d})$ . The permutation  $\pi_j$  acts on the  $|\mathcal{D}^{(b)}|$ -dimensional  $j$ th column of the out-of-bag data matrix. Essentially,  $\pi_j$  simply permutes the  $|\mathcal{D}^{(b)}|$  elements of the  $j$ th column of the out-of-bag data matrix.

**Algorithm 2.1.**

Choose a base learner  $\widehat{\mathbf{g}}(\cdot)$  ▷ e.g.: Trees, MLR  
 Choose an estimation method ▷ e.g.: Recursive Partitioning or OLS  
 Initialize all the  $\text{PERF}(\mathbf{x}_j)$  and  $\widehat{\text{VI}}(\mathbf{x}_j)$  at zero  
**for**  $b = 1$  to  $B$  **do**

Draw with replacement from  $\mathcal{D}$  a bootstrap sample  $\mathcal{D}^{(b)} = \{\mathbf{z}_1^{(b)}, \dots, \mathbf{z}_n^{(b)}\}$   
 Draw without replacement from  $\{1, \dots, p\}$  a subset  $\mathcal{V}^{(b)} = \{j_1^{(b)}, \dots, j_d^{(b)}\}$  of  $d$  variables.  
 Form the indicator vector  $\boldsymbol{\gamma}^{(b)} = (\gamma_j^{(b)}, \dots, \gamma_p^{(b)})$  with

$$\gamma_j^{(b)} = \begin{cases} 1 & \text{if } j \in \{j_1^{(b)}, \dots, j_d^{(b)}\} \\ 0 & \text{otherwise.} \end{cases}$$

Drop unselected variables from  $\mathcal{D}^{(b)}$  so that  $\mathcal{D}_{\text{sub}}^{(b)}$  is  $d$  dimensional  
 Build the  $b$ th base learner  $\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)})$  based on  $\mathcal{D}_{\text{sub}}^{(b)}$   
 Compute score of the  $b$ th base learner  $\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)})$  ▷ e.g. Out-of-bag error

$$\mathbf{s}^{(b)} = \text{score}(\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)})) = \frac{1}{|\mathcal{D}^{(b)}|} \sum_{\mathbf{z}_i \notin \mathcal{D}^{(b)}} \ell(y_i, \widehat{\mathbf{g}}(\mathbf{x}_i, \boldsymbol{\gamma}^{(b)}))$$

**for**  $j \in \mathcal{V}^{(b)}$  **do**

Generate the permutation of the  $j$ th column of  $\mathcal{D}^{(b)}$ , namely

$$\pi_j$$

Compute the permutation impacted score

$$\mathbf{s}_{\pi_j}^{(b)} = \text{score}_{\pi_j}(\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)})) = \frac{1}{|\mathcal{D}^{(b)}|} \sum_{\mathbf{z}_i \notin \mathcal{D}^{(b)}} \ell(y_i, \widehat{\mathbf{g}}(\mathbf{x}_{i,\pi_j}, \boldsymbol{\gamma}^{(b)}))$$

Compute the  $b$ th instance of the importance of  $\mathbf{x}_j$

$$\widehat{\text{VI}}^{(b)}(\mathbf{x}_j) = \mathbf{s}^{(b)} - \mathbf{s}_{\pi_j}^{(b)}$$

**end for**

**end for**

Use the ensemble  $\mathcal{H} = \{\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)}), b = 1, \dots, B\}$  to form the estimator

$$\widehat{\text{PERF}}(\mathbf{x}_j) = \frac{1}{B} \sum_{b=1}^B \text{score}(\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)})) - \frac{1}{B_j} \sum_{b=1}^B \gamma_j^{(b)} \text{score}(\widehat{\mathbf{g}}(\cdot, \boldsymbol{\gamma}^{(b)})) \quad (2.5)$$

$$\widehat{\text{VI}}(\mathbf{x}_j) = \frac{1}{B_j} \sum_{b=1}^B \gamma_j^{(b)} \widehat{\text{VI}}^{(b)}(\mathbf{x}_j) \quad (2.6)$$

### 3 Computational Demonstrations

We herein assess the goodness and usefulness of the PERF variable importance score by applying to both simulated data and real life data. We specifically use benchmark machine learning data sets like the the spam detection dataset and the pima indian diabetes dataset in the classification context and the attitude dataset in the regression context. Our first example features simulated data with different scenarios on the level of correlation among the variables, and the ratio  $n$  and  $p$ . In this particular example, the true function is

$$f(\mathbf{x}) = 1 + 2\mathbf{x}_3 + \mathbf{x}_7 + 3\mathbf{x}_9$$

with  $\mathbf{x} \sim \text{MVN}(\mathbf{1}_9, \Sigma_\rho)$  and  $\epsilon \sim \text{N}(0, 1)$ . The dataset in this example is simulated data with different scenarios on the level of correlation among the variables, and the ratio  $n$  and  $p$ . Specifically, we simulate data by defining  $\rho \in [0, 1)$ , then we generate our predictor variables using a multivariate normal distribution. Throughout this paper, the multivariate Gaussian density will be denoted by  $\phi_p(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$

$$\phi_p(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \tag{3.1}$$

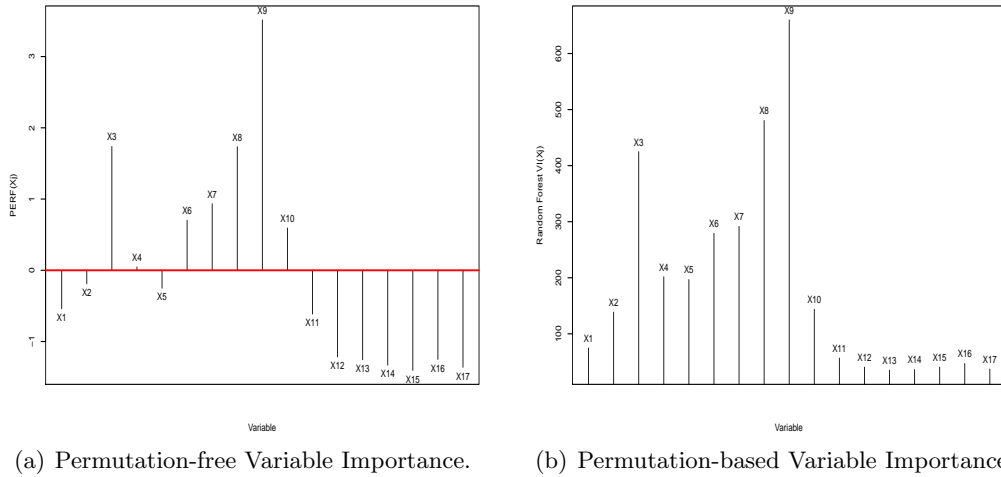
Furthermore, in order to study the effect of the correlation pattern, we simulate the data using a covariance matrix  $\Sigma$  parameterized by  $\tau$  and  $\rho$  and defined by  $\tau\Sigma$  where  $\Sigma = (\sigma_{ij})$  with  $\sigma_{ij} = \rho^{|i-j|}$ .

$$\Sigma = \Sigma(\tau, \rho) = \tau \begin{pmatrix} 1 & \rho & \dots & \rho^{p-2} & \rho^{p-1} \\ \rho & 1 & \rho & \dots & \rho^{p-2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \rho^{p-2} & \ddots & \rho & 1 & \rho \\ \rho^{p-1} & \rho^{p-2} & \dots & \rho & 1 \end{pmatrix}$$

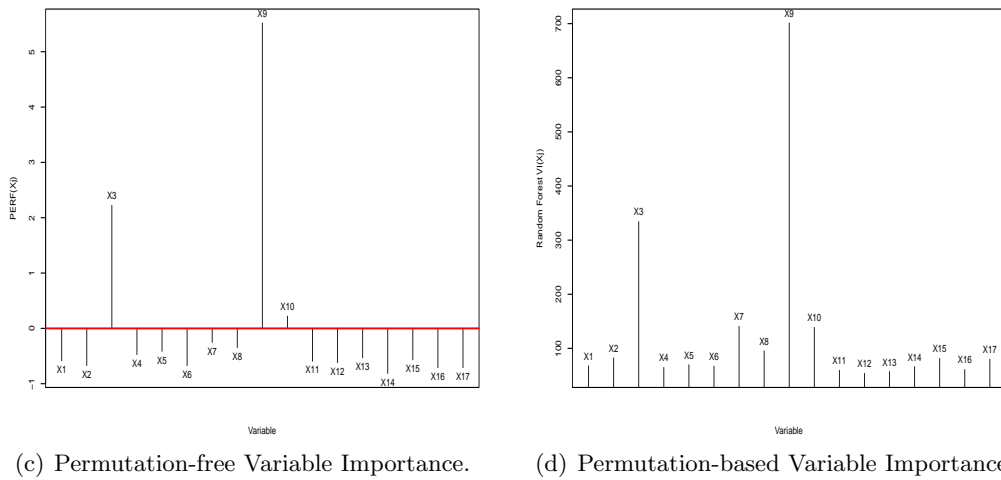
For simplicity however, we use the first  $\Sigma$  with  $\tau = 1$  throughout this paper. For the remaining parameters, we use  $\rho \in \{0, 0.25, 0.75\}$  and  $p \in \{17, 250\}$ , with the same  $n = 200$ . The plots depicting the comparisons between Random Forest and PERF variable importance scores are given in subsequent parts of this paper. In each plot, the dimensionalities of the data are appropriately indicated, and comments are provided as well.

### 4 Discussion and Conclusion

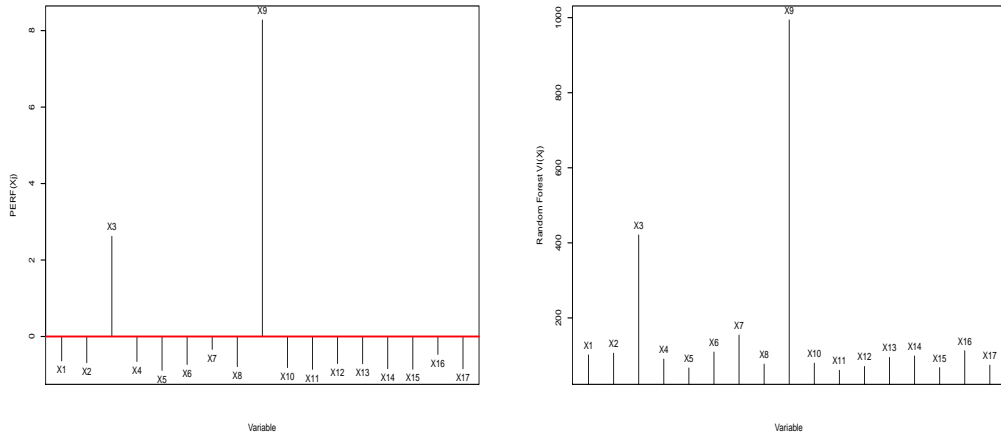
We have presented a variable importance score function in the context of ensemble learning. Our proposed score function is simple and more straightforward than its counterpart proposed in the context of random forest, and by avoiding permutations, it is by design computationally more efficient than the random forest variable importance function. Just like the random forest variable importance function, our score handles both regression and classification seamlessly. One of the distinct advantage of our proposed score is the fact that it offers a natural cut off at zero, with all the positive scores indicating importance and significance, while the negative scores are deemed indications of insignificance. An extra advantage of our proposed score lies in the fact it works very well beyond ensemble of trees and can seamlessly be used with any base learners in the random subspace learning context. Our examples, both simulated and real, demonstrated that our proposed score does compete mostly favorably with the random forest score. In our future work, we present and compare the corresponding average test errors of the single models made up of the most important variables. We also provide in our future work theoretical proofs of the connection between our score function and the significance of variables selected using existing criteria. It is also our plan to address the fact that sometimes the correlation structure among the predictor variables obscures the ability of our proposed score to correctly identify some significant variables.



**Fig. 1. Variable score for simulated data with high correlation among the variables in low dimension high sample size setting**

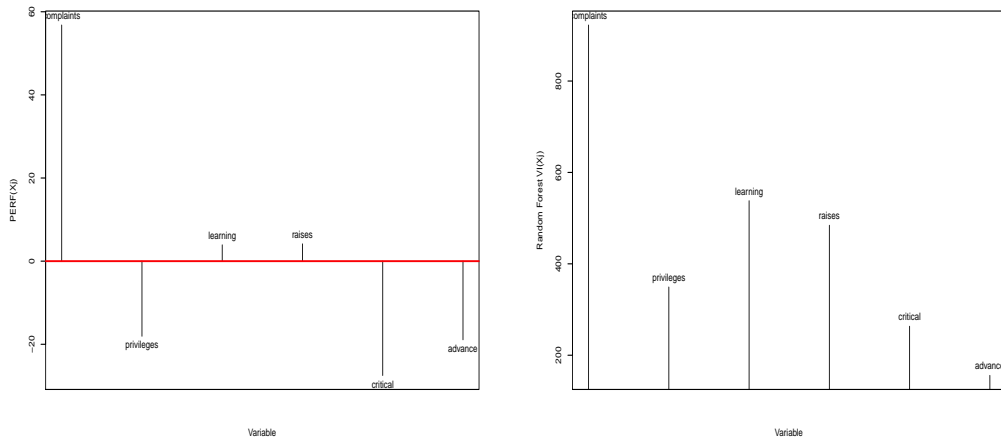


**Fig. 2. Variable Importance Scores for simulated data with mild correlation among the variables in low dimension high sample size setting**



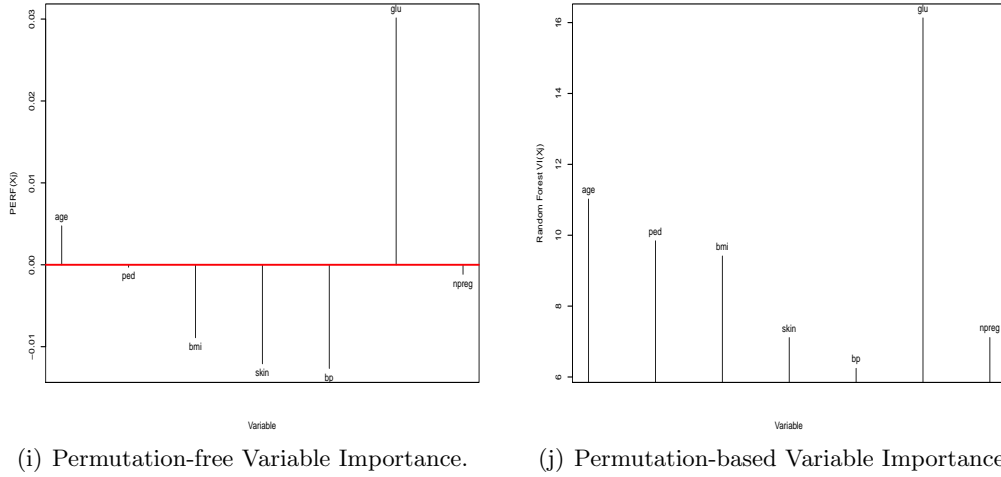
(e) Permutation-free Variable Importance. (f) Permutation-based Variable Importance.

**Fig. 3. Variable Importance Scores for simulated data with zero correlation among the variables in low dimension high sample size setting**

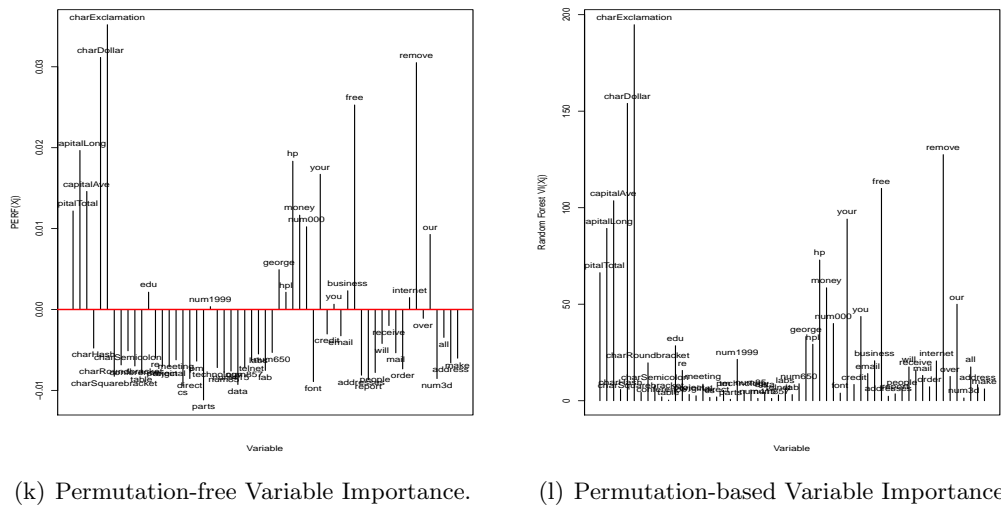


(g) Permutation-free Variable Importance. (h) Permutation-based Variable Importance.

**Fig. 4. Variable Importance Scores for the Attitude Data Set.  $n = 30$  and  $p = 6$**



**Fig. 5. Variable Importance Scores for the Spam Detection Dataset.  $n = 200$ ,  $p = 7$ , and  $K = 2$  classes**



**Fig. 6. Variable Importance Scores for the Spam Detection Dataset,  $n = 4601$ ,  $p = 57$  and  $K = 2$  classes**



## Acknowledgements

Ernest Fokoué wishes to express his heartfelt gratitude and infinite thanks to Our Lady of Perpetual Help for Her ever-present support and guidance, especially for the uninterrupted flow of inspiration received through Her most powerful intercession.

## Competing Interests

Author has declared that no competing interests exist.

## References

- [1] Breiman L. Random forests. *Machine Learning*. 2001;45:5-32.
- [2] Breiman L. Statistical modeling: The two cultures. *Statistical Science*. 2001;16(3):199-215.
- [3] Strobl C, Boulesteix AL, Zeileis A, Hothorn T. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*. 2007;8(25).
- [4] Strobl C, Boulesteix AL, Kneib T, Augustin T, Zeileis A. Conditional variable importance for random forests. *BMC Bioinformatics*. 2008;9(307):1-12.
- [5] Kuhn M. Building predictive models in r using the caret package. *Journal of Statistical Software*. 2008;28(5).
- [6] Strobl C, Hothorn T, Zeileis A. Party on. *R Journal*. 2009;1(2):14-17.
- [7] Shih YS. Regression trees with unbiased variable selection. *Statistica Sinical*. 2009;12:361-386.
- [8] Sandri M, Zuccolotto P. A bias correction algorithm for the gini variable importance measure in classification trees. *Journal of Computational and Graphical Statistics*. 2008;17(3):611-628.

---

©2016 Fokoué; This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**Peer-review history:**

The peer review history for this paper can be accessed here (Please copy paste the total link in your browser address bar)

<http://sciencedomain.org/review-history/13781>