

Testing Variable Importance in Random Forests

Carolin Strobl (LMU München) and Achim Zeileis (WU Wien)

lifestat 2008

The permutation
importance

The suggested test

Summary and
outlook

References

Introduction

random forests

The permutation
importance

The suggested test

Summary and
outlook

References

Introduction

random forests

- ▶ have become increasingly popular in, e.g., genetics and the neurosciences

The permutation
importance

The suggested test

Summary and
outlook

References

Introduction

random forests

- ▶ have become increasingly popular in, e.g., genetics and the neurosciences [imagine a long list of references here]

The permutation
importance

The suggested test

Summary and
outlook

References

Introduction

random forests

- ▶ have become increasingly popular in, e.g., genetics and the neurosciences [imagine a long list of references here]
- ▶ can deal with “small n large p”-problems, high-order interactions, correlated predictor variables

The permutation importance

The suggested test

Summary and outlook

References

Introduction

random forests

- ▶ have become increasingly popular in, e.g., genetics and the neurosciences [imagine a long list of references here]
- ▶ can deal with “small n large p”-problems, high-order interactions, correlated predictor variables
- ▶ are used not only for prediction, but also to assess variable importance and

The permutation importance

The suggested test

Summary and outlook

References

Introduction

random forests

- ▶ have become increasingly popular in, e.g., genetics and the neurosciences [imagine a long list of references here]
- ▶ can deal with “small n large p”-problems, high-order interactions, correlated predictor variables
- ▶ are used not only for prediction, but also to assess variable importance and
- ▶ on the official random forest website Breiman and Cutler (2008) even suggest a significance test for the variable importance...

The permutation importance

The suggested test

Summary and outlook

References

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

The permutation
importance

The suggested test

Summary and
outlook

References

Random forests

learn a random forest as a classification/regression model

to predict Y from X_1, \dots, X_p

The permutation
importance

The suggested test

Summary and
outlook

References

Random forests

learn a random forest as a classification/regression model

to predict Y from X_1, \dots, X_p

result: almost a black-box

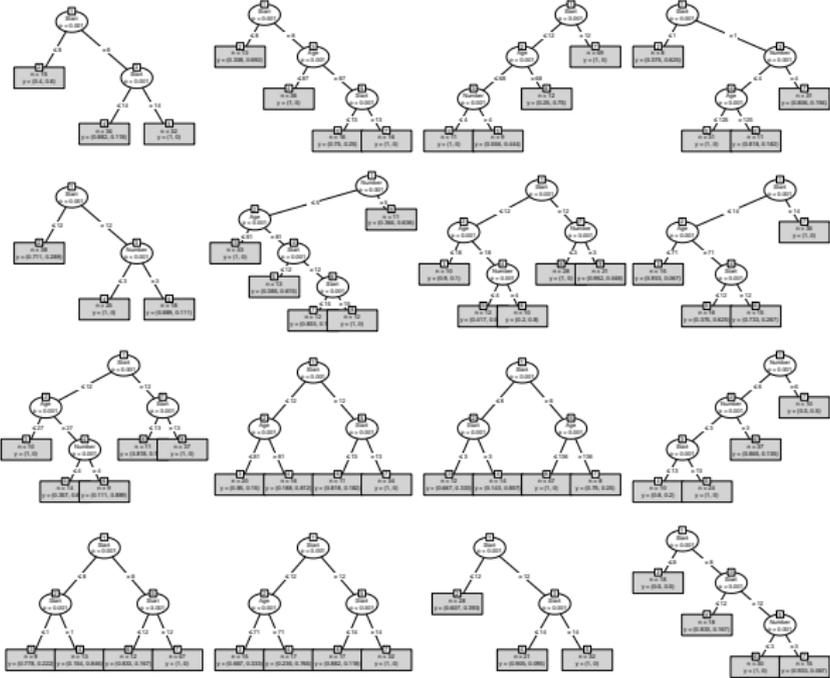
The permutation
importance

The suggested test

Summary and
outlook

References

(Small) random forest



The permutation importance

The suggested test

Summary and outlook

References

Measuring variable importance

The permutation importance

The suggested test

Summary and outlook

References

Measuring variable importance

- ▶ Gini importance

mean Gini gain produced by X_j over all trees

(can be severely biased due to estimation bias and multiple testing; Strobl et al., 2007)

The permutation
importance

The suggested test

Summary and
outlook

References

Measuring variable importance

- ▶ Gini importance

mean Gini gain produced by X_j over all trees

(can be severely biased due to estimation bias and multiple testing; Strobl et al., 2007)

- ▶ permutation importance

mean decrease in classification accuracy after permuting X_j over all trees

- ▶ informative variables produce a systematic decrease in accuracy when permuted
- ▶ uninformative variables produce a random decrease or increase in accuracy when permuted

(unbiased when subsampling is used; Strobl et al., 2007)

The permutation
importance

The suggested test

Summary and
outlook

References

The permutation importance

within each tree t

$$VI^{(t)}(\mathbf{x}_j) = \frac{\sum_{i \in \overline{\mathfrak{B}}^{(t)}} I(y_i = \hat{y}_i^{(t)})}{|\overline{\mathfrak{B}}^{(t)}|} - \frac{\sum_{i \in \overline{\mathfrak{B}}^{(t)}} I(y_i = \hat{y}_{i, \pi_j}^{(t)})}{|\overline{\mathfrak{B}}^{(t)}|}$$

$\hat{y}_i^{(t)} = f^{(t)}(\mathbf{x}_i)$ = predicted class before permuting

$\hat{y}_{i, \pi_j}^{(t)} = f^{(t)}(\mathbf{x}_{i, \pi_j})$ = predicted class after permuting X_j

$\mathbf{x}_{i, \pi_j} = (x_{i,1}, \dots, x_{i,j-1}, x_{\pi_j(i),j}, x_{i,j+1}, \dots, x_{i,p})$

Note: $VI^{(t)}(\mathbf{x}_j) = 0$ by definition, if X_j is not in tree t

The permutation importance

The suggested test

Summary and outlook

References

The permutation importance

over all trees:

1. raw importance

$$VI(\mathbf{x}_j) = \frac{\sum_{t=1}^{ntree} VI^{(t)}(\mathbf{x}_j)}{ntree}$$

The permutation importance

The suggested test

Summary and outlook

References

The permutation importance

over all trees:

2. scaled importance (z-score)

$$\frac{VI(\mathbf{x}_j)}{\frac{\hat{\sigma}}{\sqrt{ntree}}} = z_j$$

The permutation importance

The suggested test

Summary and outlook

References

- ▶ individual $VI^{(t)}(\mathbf{x}_j)$ have standard deviation σ and

The permutation
importance

The suggested test

Investigating the
statistical properties

Specifying the null
hypothesis

Summary and
outlook

References

- ▶ individual $VI^{(t)}(\mathbf{x}_j)$ have standard deviation σ and
- ▶ are computed from n_{tree} independent bootstrap samples (where n_{tree} is large)

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

- ▶ individual $VI^{(t)}(\mathbf{x}_j)$ have standard deviation σ and
- ▶ are computed from n_{tree} independent bootstrap samples (where n_{tree} is large)
- ▶ central limit theorem for the mean $VI(\mathbf{x}_j)$
 \Rightarrow normal with standard error $\sigma/\sqrt{n_{tree}}$

The permutation
importance

The suggested test

Investigating the
statistical properties

Specifying the null
hypothesis

Summary and
outlook

References

- ▶ individual $VI^{(t)}(\mathbf{x}_j)$ have standard deviation σ and
- ▶ are computed from n_{tree} independent bootstrap samples (where n_{tree} is large)
- ▶ central limit theorem for the mean $VI(\mathbf{x}_j)$
 \Rightarrow normal with standard error $\sigma/\sqrt{n_{tree}}$

under the null hypothesis of zero importance:

$$z_j \stackrel{as.}{\sim} N(0, 1)$$

The permutation
importance

The suggested test

Investigating the
statistical properties

Specifying the null
hypothesis

Summary and
outlook

References

The suggested test

if z_j exceeds the α -quantile of $N(0,1) \Rightarrow$ reject the null hypothesis of zero importance for variable X_j

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

Simulation study

- ▶ generate data sets of sample size $n = 100, 200$ and 500
- ▶ five predictor variables of which only X_1 is relevant with
- ▶ $y \sim \begin{cases} B(n, 0.5 - \rho) & \text{for } X_1 = 0 \\ B(n, 0.5 + \rho) & \text{for } X_1 = 1 \end{cases}$ with relevance
 $\rho = 0, 0.05, \dots, 0.5$
- ▶ fit random forests with $n_{tree} = 100, 200$ and 500
- ▶ for 1000 iterations in each combination:
count how many times the null hypothesis for X_1 was rejected

The permutation importance

The suggested test

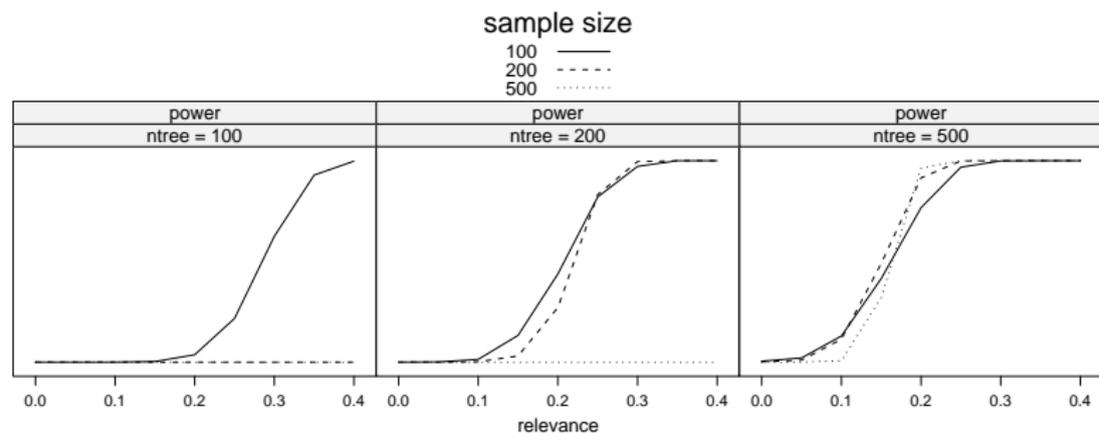
Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

The power



The permutation
importance

The suggested test

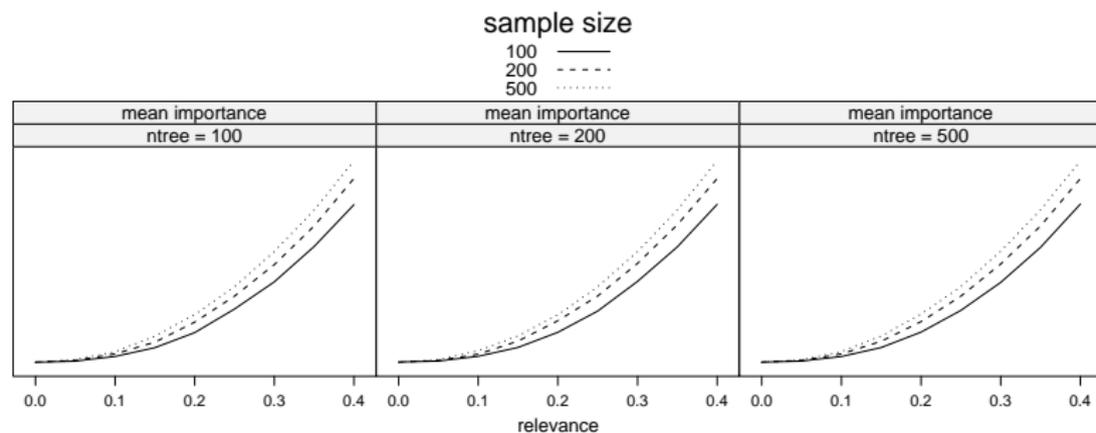
Investigating the
statistical properties

Specifying the null
hypothesis

Summary and
outlook

References

The average raw importance



The permutation
importance

The suggested test

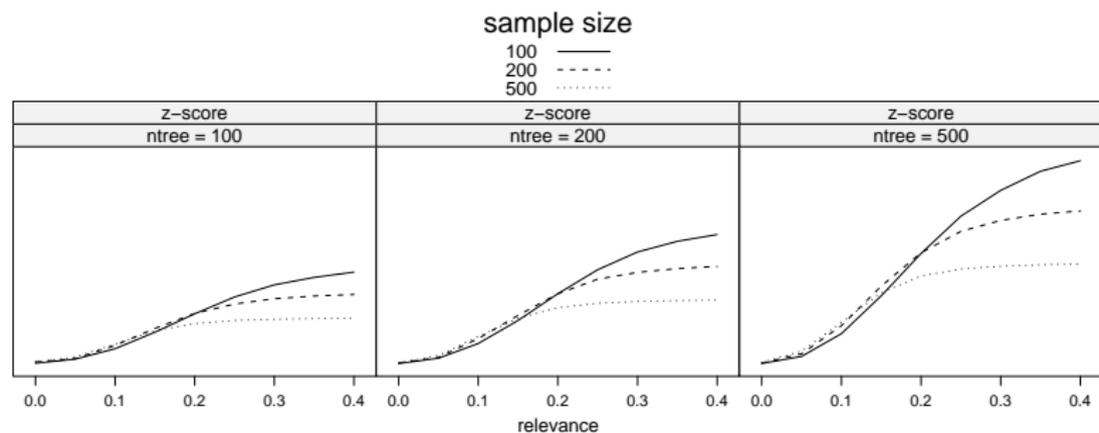
Investigating the
statistical properties

Specifying the null
hypothesis

Summary and
outlook

References

The average z-score



The permutation importance

The suggested test

Investigating the statistical properties

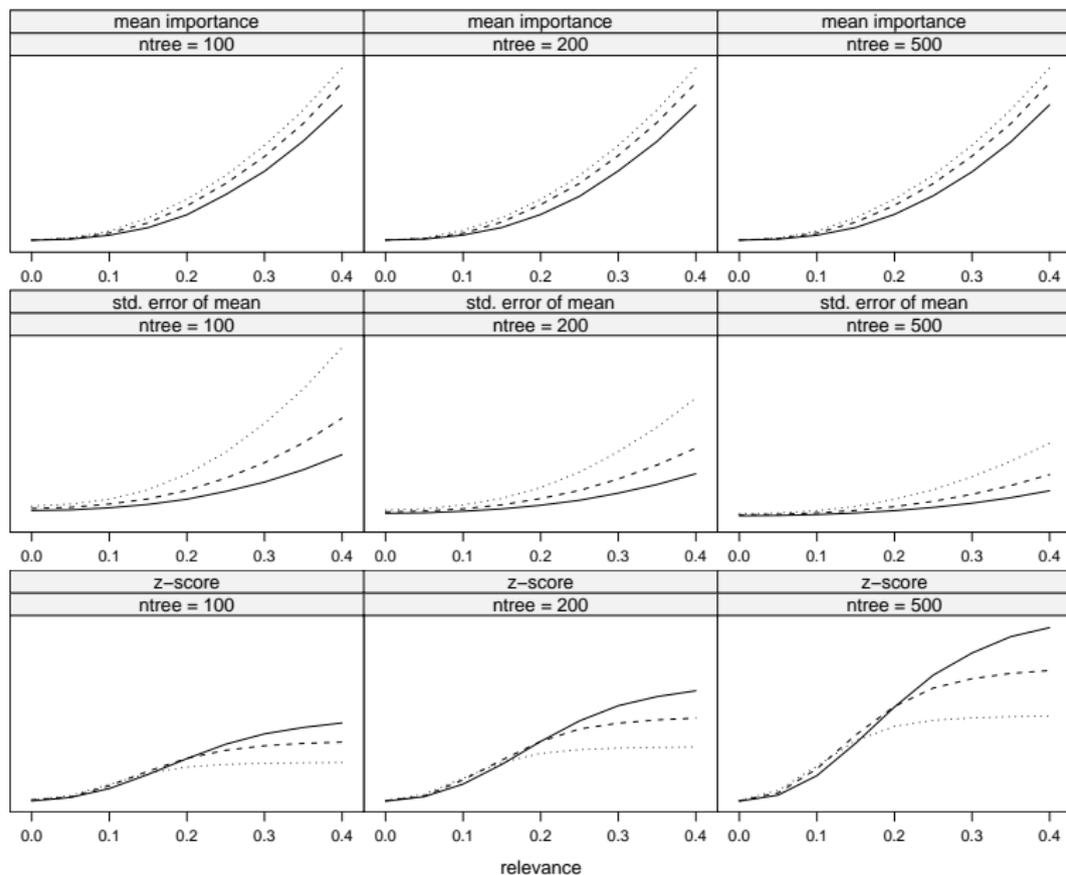
Specifying the null hypothesis

Summary and outlook

References

sample size

100 ———
200 - - - -
500 ·····



The permutation importance

The suggested test

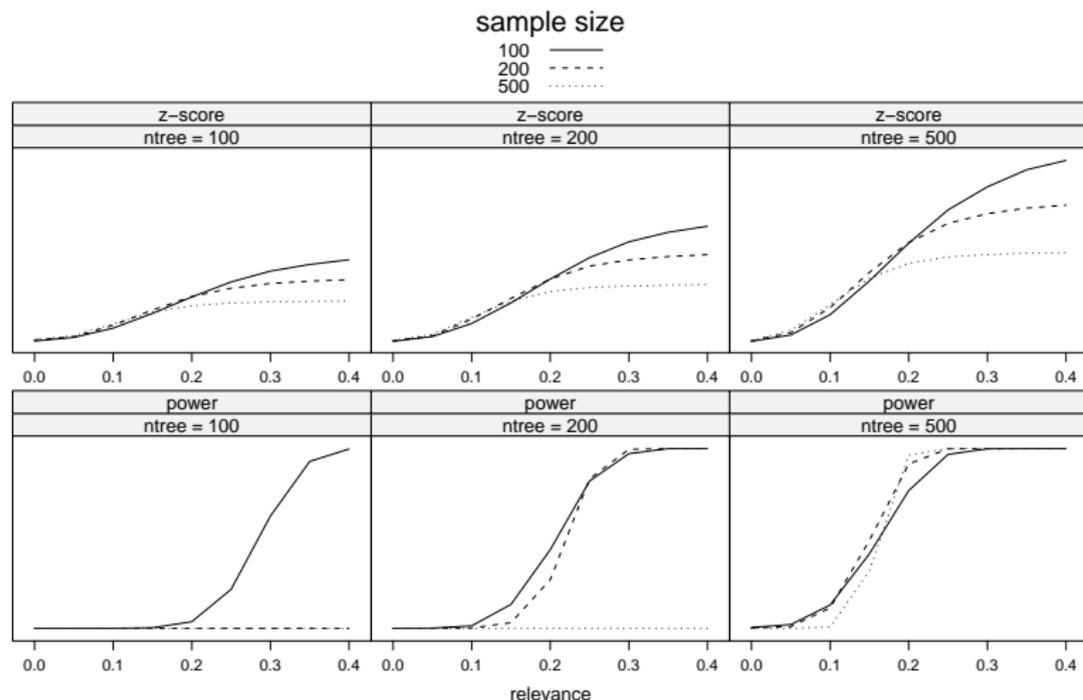
Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

The average z-score and the power



The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

Findings

z-score and power

- ▶ increase in the number of trees
- ▶ decrease in the sample size

The permutation
importance

The suggested test

**Investigating the
statistical properties**

Specifying the null
hypothesis

Summary and
outlook

References

What null hypothesis were we testing in the first place?

<i>obs</i>	<i>Y</i>	<i>X_j</i>	<i>Z</i>
1	<i>y</i> ₁	<i>x</i> _{$\pi_j(1),j$}	<i>z</i> ₁
⋮	⋮	⋮	⋮
<i>i</i>	<i>y</i> ₁	<i>x</i> _{$\pi_j(i),j$}	<i>z</i> _{<i>i</i>}
⋮	⋮	⋮	⋮
<i>n</i>	<i>y</i> ₁	<i>x</i> _{$\pi_j(n),j$}	<i>z</i> _{<i>n</i>}

$$H_0 : X_j \perp Y, Z \text{ or } X_j \perp Y \wedge X_j \perp Z$$

$$P(Y, X_j, Z) \stackrel{H_0}{=} P(Y, Z) \cdot P(X_j)$$

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

What null hypothesis were we testing in the first place?

the current null hypothesis reflects independence of X_j from both Y and the remaining predictor variables Z

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

What null hypothesis were we testing in the first place?

the current null hypothesis reflects independence of X_j from both Y and the remaining predictor variables Z

⇒ a high variable importance can result from violation of either one

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

Conditional permutation scheme

<i>obs</i>	<i>Y</i>	<i>X_j</i>	<i>Z</i>
1	<i>y</i> ₁	$X_{\pi_{j Z=a}(1),j}$	$z_1 = a$
3	<i>y</i> ₃	$X_{\pi_{j Z=a}(3),j}$	$z_3 = a$
27	<i>y</i> ₂₇	$X_{\pi_{j Z=a}(27),j}$	$z_{27} = a$
6	<i>y</i> ₆	$X_{\pi_{j Z=b}(6),j}$	$z_6 = b$
14	<i>y</i> ₁₄	$X_{\pi_{j Z=b}(14),j}$	$z_{14} = b$
33	<i>y</i> ₃₃	$X_{\pi_{j Z=b}(33),j}$	$z_{33} = b$
⋮	⋮	⋮	⋮

$$H_0 : X_j \perp Y | Z$$

$$P(Y, X_j | Z) \stackrel{H_0}{=} P(Y | Z) \cdot P(X_j | Z)$$

$$\text{or } P(Y | X_j, Z) \stackrel{H_0}{=} P(Y | Z)$$

The permutation importance

The suggested test

Investigating the statistical properties

Specifying the null hypothesis

Summary and outlook

References

to be continued...

The permutation
importance

The suggested test

Investigating the
statistical properties

**Specifying the null
hypothesis**

Summary and
outlook

References

Summary and outlook

the significance test suggested on the random forest website has strange properties:

The permutation importance

The suggested test

Summary and outlook

References

Summary and outlook

the significance test suggested on the random forest website has strange properties:

- ▶ the z-score and power increase in the number of trees and decrease in the sample size

The permutation importance

The suggested test

Summary and outlook

References

Summary and outlook

the significance test suggested on the random forest website has strange properties:

- ▶ the z-score and power increase in the number of trees and decrease in the sample size
- ▶ the null hypothesis may not reflect what you wanted

The permutation importance

The suggested test

Summary and outlook

References

Summary and outlook

the significance test suggested on the random forest website has strange properties:

- ▶ the z-score and power increase in the number of trees and decrease in the sample size
 - ▶ the null hypothesis may not reflect what you wanted
- ⇒ use conditional permutation scheme

The permutation importance

The suggested test

Summary and outlook

References

Summary and outlook

the significance test suggested on the random forest website has strange properties:

- ▶ the z-score and power increase in the number of trees and decrease in the sample size
- ▶ the null hypothesis may not reflect what you wanted
 - ⇒ use conditional permutation scheme
 - ⇒ use distribution over > 1 permutations

The permutation importance

The suggested test

Summary and outlook

References

Summary and outlook

the significance test suggested on the random forest website has strange properties:

- ▶ the z-score and power increase in the number of trees and decrease in the sample size
- ▶ the null hypothesis may not reflect what you wanted
 - ⇒ use conditional permutation scheme
 - ⇒ use distribution over > 1 permutations

for now: stick to the unscaled importance

The permutation importance

The suggested test

Summary and outlook

References

The permutation
importance

The suggested test

Summary and
outlook

References

Breiman, L. and A. Cutler (2008). Random forests -
classification manual (website accessed in 1/2008)
<http://www.math.usu.edu/~adele/forests/>.

Strobl, C., A.-L. Boulesteix, A. Zeileis, and T. Hothorn
(2007). Bias in random forest variable importance
measures: Illustrations, sources and a solution. *BMC
Bioinformatics* 8:25.