Advanced Statistical Topics in Health Research

Day 4 (October 23, 2025)

Trees and forests

Helene Charlotte Wiese Rytgaard (hely@sund.ku.dk)
Mark Bech Knudsen (mark.bech.knudsen@sund.ku.dk)

Section of Biostatistics, Department of Public Health

Outline of today's topics

Part 1: Modeling cultures, model (Mark Bech Knudsen) selection and decision trees

Part 2: From trees to forests, and (Helene Rytgaard) recap on classical machine learning techniques

Part 3: Tuning random forests (Mark Bech Knudsen)

Part 4: Variable importance and (Helene Rytgaard) interpretable machine learning tools

Software Overview

		Outcor	me		
Package	Conti-	Binary	Survival	Comp.	Method
	nuous			risks	
rpart ¹	Χ	Χ	Χ		Tree
${\sf randomForest}^2$	Χ	Χ			Forest
party ³	Χ	Χ	Χ		Tree/Forest
randomForestSRC ⁴	X	Χ	Χ	Χ	Forest
ranger ⁵	X	X	Χ		Forest

¹rpart Therneau, Atkinson and Ripley

²randomForest Liaw and Wiener (based on Breiman and Cutler)

³ctree, cforest Hothorn

⁴rfsrc Ishwaran

⁵ranger Wright and Ziegler

Two modeling cultures

The two cultures

In a very influential commentary⁶, Leo Breiman outlines two dominant modeling cultures.

For a response variable Y and predictors X, we imagine the world as the following diagram:



⁶Leo Breiman, "Statistical Modeling: The Two Cultures", Statist. Sci. 16(3), 199-231, 2001.

The two cultures

In a very influential commentary⁶, Leo Breiman outlines two dominant modeling cultures.

For a response variable Y and predictors X, we imagine the world as the following diagram:



The goal of a data analysis is generally one or both of

▶ **Prediction:** For a given *X*, what is *Y* likely going to be?

⁶Leo Breiman, "Statistical Modeling: The Two Cultures", Statist. Sci. 16(3), 199-231, 2001.

The two cultures

In a very influential commentary⁶, Leo Breiman outlines two dominant modeling cultures.

For a response variable Y and predictors X, we imagine the world as the following diagram:



The goal of a data analysis is generally one or both of

- ▶ **Prediction**: For a given *X*, what is *Y* likely going to be?
- ▶ **Information**: By which "mechanism" is nature associating *X* to *Y*?
 - Which components of X are most important, and how strong is the association?

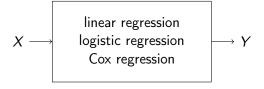
⁶Leo Breiman, "Statistical Modeling: The Two Cultures", Statist. Sci. 16(3), 199-231, 2001.

The data modeling culture

The data modeling culture starts by assuming a model

$$Y = f(X, parameters, noise)$$

where the function f is chosen by the analyst.

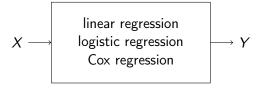


The data modeling culture

The data modeling culture starts by assuming a model

$$Y = f(X, parameters, noise)$$

where the function f is chosen by the analyst.

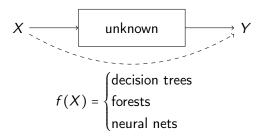


The goal is to estimate the parameters, which encode information about how X relates to Y. Afterwards, the model can also be used for prediction.

The data modeling culture has historically been very dominant.

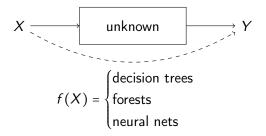
The algorithmic modeling culture

The algorithmic modeling culture regards nature as unknown. Uses data to find f(X) which predicts Y with high accuracy.



The algorithmic modeling culture

The algorithmic modeling culture regards nature as unknown. Uses data to find f(X) which predicts Y with high accuracy.



The algorithmic modeling culture has historically played a very minor role. Nowadays very important (machine learning, AI).

Epo case study

Anemia is a deficiency of red blood cells and/or hemoglobin and an additional risk factor for cancer patients.

Randomized placebo controlled trial⁷: does treatment with epoetin beta - epo - (300 U/kg) enhance hemoglobin concentration level and improve survival chances?

Henke et al. 2006 identified the c20 expression (erythropoietin receptor status) as a new biomarker for the prognosis of locoregional progression-free survival.

⁷Henke et al. Do erythropoietin receptors on cancer cells explain unexpected clinical findings? J Clin Oncol, 24(29):4708-4713, 2006.

Treatment

The study includes 149 head and neck cancer patients⁸ with a tumor located in the oropharynx (36%), the oral cavity (27%), the larynx (14%) or in the hypopharynx (23%).

	Resection		
	Complete	Incomplete	No
Placebo	35	14	25
Еро	36	14	25

⁸with non-missing blood values

Outcome

Blood hemoglobin levels were measured weekly during radiotherapy (7 weeks).

Treatment with epoetin beta was defined successful when the hemoglobin level increased sufficiently. For patient i set

$$Y_i = \begin{cases} 1 & \text{treatment successful} \\ 0 & \text{treatment failed} \end{cases}$$

Predictors

Age min: 41 y, median: 59 y, max: 80 y

Gender male: 85%, female: 15%
Base hemoglobin mean: 12.03 g/dl, std: 1.45
Treatment epo: 50%, placebo 50%

Resection complete: 48%, incomplete: 19%,

no resection: 34%

Epo receptor status neg: 32%, pos: 68%

Goal of analysis

Want to know: Is treatment with epo more likely to increase hemoglobin levels sufficiently compared to placebo?

Also: predict probability of hemoglobin increase given predictor variables.

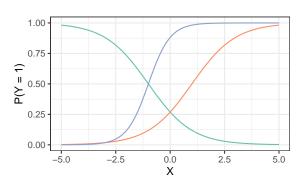
Logistic regression (data modeling culture)

Logistic regression

For binary outcome it is common to use logistic regression:

$$P(Y = 1) = f(X, parameters) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}}$$

Parameters β_k have interpretation as (log) odds-ratios.



Logistic regression results

We fit a logistic regression with all covariates (no interactions).

Covariate	${\sf OddsRatio}$	CI.95	pValue
(Intercept)	0.00		0.0040
age	0.97	[0.91; 1.03]	0.2807
sexmale	0.21	[0.038; 1.10]	0.0657
HbBase	3.26	[1.99; 5.91]	< 0.0001
TreatPlacebo	0.01	[0.0020; 0.042]	< 0.0001
ResectionIncompl	0.42	[0.083; 1.96]	0.2801
ResectionNo	0.24	[0.058; 0.89]	0.0395
Receptorpositive	5.81	[1.72; 23.39]	0.0076

Logistic regression results

We fit a logistic regression with all covariates (no interactions).

Covariate	${\sf OddsRatio}$	CI.95	pValue
(Intercept)	0.00		0.0040
age	0.97	[0.91; 1.03]	0.2807
sexmale	0.21	[0.038; 1.10]	0.0657
HbBase	3.26	[1.99; 5.91]	< 0.0001
TreatPlacebo	0.01	[0.0020; 0.042]	< 0.0001
ResectionIncompl	0.42	[0.083; 1.96]	0.2801
ResectionNo	0.24	[0.058; 0.89]	0.0395
Receptorpositive	5.81	[1.72; 23.39]	0.0076

Epo treated has much higher odds of hemoglobin level increase. Does that mean everyone should be treated?

The model provides information for a single patient

For example: the predicted probability that a 48 year old male with no tumor resection, negative receptor status and baseline hemoglobin level $10.8 \ g/dl$ reaches the target hemoglobin level $(Y_i = 1)$ is

Epo treatment: 16.9%

Placebo group: 0.2%

The model provides information for a single patient

For example: the predicted probability that a 48 year old male with no tumor resection, negative receptor status and baseline hemoglobin level $10.8 \ g/dl$ reaches the target hemoglobin level $(Y_i = 1)$ is

Epo treatment: 16.9%

Placebo group: 0.2%

If a similar patient has baseline hemoglobin level 12.8 g/dl then the model predicts:

Epo treatment: 68.4%

Placebo group: 2.3%

Model selection

Very many different "logistic regression models" can be constructed by selecting subsets of variables, transformations, and interactions of variables.

For example, should we include age, since it is not statistically significant?

Model selection

Very many different "logistic regression models" can be constructed by selecting subsets of variables, transformations, and interactions of variables.

For example, should we include age, since it is not statistically significant?

Backward elimination iteratively removes covariate with highest p-value (if above 0.05).

```
library(rms)
mod_all <- lrm(Y~age+sex+HbBase+Treat+Resection+Receptor,data=Epo)
fastbw(mod_all)</pre>
```

```
Deleted Chi-Sq d.f. P Residual d.f. P AIC age 1.16 1 0.2807 1.16 1 0.2807 -0.84 Resection 3.75 2 0.1532 4.92 3 0.1781 -1.08
```

Reduced model

Then refit model with reduced covariates:

```
mod_bw <- lrm(Y~sex+HbBase+Treat+Receptor,data=Epo)
```

Covariate	OddsRatio (Full)	CI.95	pValue
(Intercept)	0.00 (0.00)		<0.0001
sexmale	0.16 (0.21)	[0.032; 0.74]	0.0213
HbBase	3.22 (3.26)	[2.02; 5.62]	< 0.0001
TreatPlacebo	0.02 (0.01)	[0.0039; 0.054]	< 0.0001
Receptorpositive	4.49 (5.81)	[1.45; 15.90]	0.0129

Reduced model

Then refit model with reduced covariates:

```
mod_bw <- lrm(Y~sex+HbBase+Treat+Receptor,data=Epo)
```

Covariate	OddsRatio (Full)	CI.95	pValue
(Intercept)	0.00 (0.00)		<0.0001
sexmale	0.16 (0.21)	[0.032; 0.74]	0.0213
HbBase	3.22 (3.26)	[2.02; 5.62]	< 0.0001
TreatPlacebo	0.02 (0.01)	[0.0039; 0.054]	< 0.0001
Receptorpositive	4.49 (5.81)	[1.45; 15.90]	0.0129

Predictions for the 48 year old male from previously:

	Full model	Reduced model
Еро	16.9%	24.1%
Placebo	0.2%	0.5%

Is ad-hoc model selection a good idea?



Journal of Clinical Epidemiology

Journal of Clinical Epidemiology 57 (2004) 1138-1146

Automated variable selection methods for logistic regression produced unstable models for predicting acute myocardial infarction mortality

Peter C. Austina,b,c,*, Jack V. Tua,b,c,d,e

Abstract

Objectives: Automated variable selection methods are frequently used to determine the independent predictors of an outcome. T objective of this study was to determine the reproducibility of logistic regression models developed using automated variable selectimethods.

Study Design and Setting: An initial set of 29 candidate variables were considered for predicting mortality after acute myocard infarction (AMI). We drew 1,000 bootstrap samples from a dataset consisting of 4,911 patients admitted to hostital with an AMI. Usi each bootstrap sample, logistic regression models predicting 30-day mortality were obtained using backward elimination, forward selectic and stepwise selection. The agreement between the different model selection methods and the agreement across the 1,000 bootstrap samp were compared.

Results: Using 1,000 bootstrap samples, backward elimination identified 940 unique models for predicting mortality. Similar resu were obtained for forward and stepwise selection. Three variables were identified as independent predictors of mortality among all bootstr samples. Over half the candidate prognostic variables were identified as independent predictors in less than half of the bootstrap sample.

Conclusion: Automated variable selection methods result in models that are unstable and not reproducible. The variables selected independent predictors are sensitive to random fluctuations in the data. © 2004 Elsevier Inc. All rights reserved.

Keywords: Regression models; Multivariate analysis; Variable selection; Logistic regression; Acute myocardial infarction; Epidemiology

Sensitivity to small changes

We randomly select 130 patients out 149 and do as before (fit full model and run backward elimination).

Sensitivity to small changes

We randomly select 130 patients out 149 and do as before (fit full model and run backward elimination).

```
Deleted Chi-Sq d.f. P Residual d.f. P AIC age 0.96 1 0.3264 0.96 1 0.3264 -1.04 sex 2.07 1 0.1501 3.03 2 0.2193 -0.97
```

Factors in Final Model

[1] HbBase Treat Resection Receptor

Different variables are removed from the model!

Predict chance of treatment success for new patient

```
newpatient
```

```
age sex HbBase Treat Resection Receptor 1 48 male 10.8 Epo No negative
```

```
library(riskRegression)
mod_sub_bw <- lrm(Y~HbBase+Treat+Resection+Receptor, data=Epo_sub)
predictRisk(mod_sub_bw, newpatient)</pre>
```

[1] 0.2013964

Exercise

Load the Epo data into R:

```
Epo <- read.csv(
   "https://biostatistics.dk/teaching/advtopics/data/Epo.csv",
   stringsAsFactors=TRUE)
newpatient <- read.csv(
   "https://biostatistics.dk/teaching/advtopics/data/newpatient.csv",
   stringsAsFactors=TRUE)</pre>
```

- Choose your favorite seed to generate a subsample (n=130) of the Epo data
- Run backward elimination with function rms::fastbw
- Predict the outcome for the new patient
- Report the selected variables and the predicted risk

When to do regression?

"Standard" multiple (logistic) regression works if

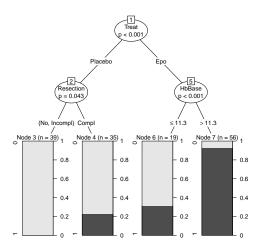
- the number of predictors is not too large, and substantially smaller than the sample size
- the decision maker has a priori knowledge about which variables to put into the model

Ad-hoc model selection algorithms, like automated backward elimination, do not lead to reproducible prediction models or selected covariate sets!

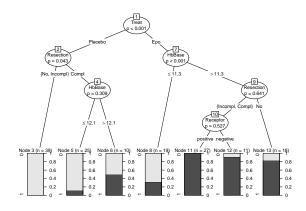
Decision trees (algorithmic modeling culture)

Classification and regression trees

```
library(party)
plot(ctree(Y~age+sex+HbBase+Treat+Resection+Receptor,data=Epo))
```



A deeper more greedy tree



Classification trees

A tree model is a form of recursive partitioning.

It lets the data decide which variables are important and where to place cut-offs in continuous variables.

In general terms, a tree-building algorithm attempts to determine a set of splits that permit accurate prediction or classification of cases.

In other words: a tree can be thought of as a sequence of many medical tests.

Roughly, the algorithm works as follows:

- Find the predictor so that the best possible split on that predictor optimizes some statistical criterion over all possible splits on the other predictors.
- 2. For ordinal and continuous predictors, the split is of the form X < c versus $X \ge c$.
- 3. Repeat step 1 within each previously formed subset.
- Proceed until fewer than k observations remain to be split, or until nothing is gained from further splitting, i.e. the tree is fully grown.
- 5. The tree is pruned according to some criterion.

Characteristics of classification trees

- Trees are specifically designed for accurate classification/prediction
- Results have a graphical representation and are easy to interpret
- No model assumptions
- Recursive partitioning can identify complex interactions
- One can introduce different costs of misclassification in the tree

But:

- Trees are not robust against even small perturbations of the data (like backward elimination)
- It is quite easy to over-fit the data
- Trees are weak learners

A Conversation of Richard Olshen with Leo Breiman

Olshen: What about arcing, bagging and boosting?

Breiman: Okay. Yeah. This is fascinating stuff, Richard. In the last five years, there have been some really big breakthroughs in prediction. And I think combining predictors is one of the two big breakthroughs. And the idea of this was, okay, that suppose you take CART, which is a pretty good classifier, but not a great classifier. I mean, for instance, neural nets do a much better job.

Olshen: Well, suitably trained? Breiman: Suitably trained. Olshen: Against an untrained

CART?

Breiman: Right. Exactly. And I think I was thinking about this. I had written an article on subset selection in linear regression. I had realized then that subset selection in linear regression is really a very unstable procedure. If you tamper with the data just a little bit, the first best five variable regression may change to another set of five variables. And so I thought, "Okay. We can stabilize this by just perturbing the data a little and get the best five variable predictor. Perturb it again. Get the best five variable predictor and then average all these five variable predictors." And sure enough, that worked out beautifully. This was published in an article in the Annals (Breiman, 1996b).

Random forests

Outline of today's remaining topics

Part 1: Modeling cultures, model (Mark Bech Knudsen selection and decision trees

Part 2: From trees to forests, and recap on classical machine learning techniques

(Helene Rytgaard)

Part 3: Tuning random forests

(Mark Bech Knudsen)

Part 4: Variable importance and interpretable machine learning tools

(Helene Rytgaard)

From trees to forests

Decision trees are nice because:

- They produce results that are easy to interpret
- They require no model assumptions

But:

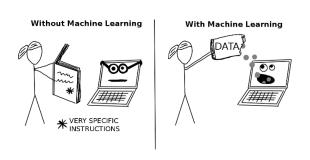
- They easily overfit the data
- Trees are weak learners

A random forest is a machine learning method that combines a large collection of decision trees to construct a strong learner

Machine learning versus classical statistics

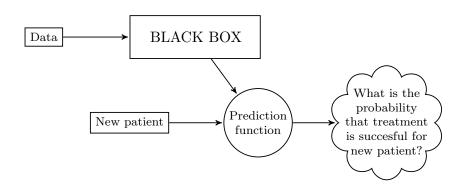
When does it make sense to apply "machine learning"?

- Little knowledge of the system we wish to analyze
- No prespecified hypotheses
- Focus on prediction rather than understanding



Machine learning as a powerful prediction tool

What is prediction?



Why do we need prediction in medical research? (Ex. 1)

Combined test at 12-week pregnancy scan

- the age of the mother, a blood sample and a measurement of fetus' neck are combined to provide a prediction of the risk of the baby having Down's syndrome, Edwards' syndrome or Patau's syndrome
- those with higher-risk results can have a subsequent diagnostic test that can tell for sure if the baby has Down's syndrome, Edwards' syndrome or Patau's syndrome but can in rare cases cause miscarriage

Why do we need prediction in medical research? (Ex. 2)

Early detection of diabetic retinopathy

- Diabetic retinopathy is a leading cause of blindness
- Diabetic retinopathy may go unnoticed until it is too late for effective treatment
- A prediction model based on fundus photography data can help detect patients with diabetic retinopathy in time for effective therapeutic intervention





Application of Random Forests Methods to Diabetic Retinopathy Classification Analyses

Ramon Casanova¹*, Santiago Saldana¹, Emily Y. Chew², Ronald P. Danis³, Craig M. Greven⁴, Walter T. Ambrosius¹

1 Department of Biostatistical Sciences, Wake Forest School of Medicine, Winston-Salem, North Carolina, United States of America, 2 National Eye Institute, National Institutes of Health (NIHI,) Bethesda, Maryland, United States of America, 3 Fundus Photograph Reading Center, University of Wisconsin, Madison, Wisconsin, United States of America, 4 Wake Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wake Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina Programme Vision States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina, United States of America, 4 Wase Forest School of Medicine Winston-Salem. North Carolina Sc

Why do we need prediction in medical research? (Ex. 3)

Prediction of long-term survival after esophagectomy

- Esophagectomy is a highly invasive surgical treatment
- A prediction model can combine multiple risk factors to provide personalized survival predictions
- This can further enable identification of high-risk patients for enhanced surveillance and/or treatment intensification

The AUGIS Survival Predictor

Prediction of Long-term and Conditional Survival after Esophagectomy Using Random Survival Forests

Rahman, Saqib A. MRCS*, "I; Walker, Robert C. MRCS*; Maynard, Nick FRCS1; Trudgill, Nigel MBBS1; Crosby, Tom FRCP5; Cromwell, David A. PhD1; Underwood, Timothy J. PhD1 on behalf of the NOGCA project team AUGIS

Author Information ⊗

Why do we need prediction in medical research? (Ex. 4)

Cancer class classification

- Accurate cancer classification can be used to target specific therapies to distinct tumor types
- A prediction model can be used to provide a data-based classification algorithm based on gene expression monitoring⁹

Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

T. R. Golub, 1.2*† D. K. Slonim, 1† P. Tamayo, 1 C. Huard, 1 M. Gaasenbeek, 1 J. P. Mesirov, 1 H. Coller, 1 M. L. Loh, 2 J. R. Downing, 3 M. A. Caligiuri, 4 C. D. Bloomfield, 4 E. S. Lander, 1.5*

⁹this is where we will end today, using a random forest (n = 38, p = 3051)

Random forests as a classical machine learning method

- 1. Let's try to understand what goes on inside the forest
 - ▶ Recap on basic machine learning techniques

- 2. Applying random forests
 - Hyperparameter selection/tuning
- 42 for this data point?

Why did you predict

- 3. Interpretability of random forests¹⁰
 - Variable importance
 - Partial Dependence Plots (PDPs)

¹⁰(and machine learning in general)

From trees to forests

From trees to forests

What is a forest 11 ...



A random forest combines the information from a collection of weak learners = randomized decision trees

- 1. Each tree is built on a bootstrap sample of the data
- 2. Only a small number of randomly selected predictor variables are used to find the best split of each node

The forest predictions are averages over the individual trees

¹¹Leo Breiman (2001). "Random Forests". Machine Learning 45 (1), 5-32,

Classical machine learning techniques utilized inside the forest

- 1. Bootstrap sampling
- 2. Nearest neighbor smoothing
- 3. Ensemble learning



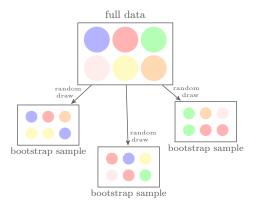
Classical machine learning techniques utilized inside the forest

- 1. Bootstrap sampling
- 2. Nearest neighbor smoothing
- 3. Ensemble learning



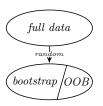
Trees are built on bootstrapped subsamples of the data

The purpose of bootstrapping is to create new pseudo samples, each of which will be used to fit a tree



Trees are built on bootstrapped subsamples of the data

Each time we draw a random bootstrap sample:



inbag: subjects in the bootstrap sample

oob: subjects not in the bootstrap sample

There are n = 149 subjects in the Epo data

```
n <- nrow(Epo)
```

Let's get a bootstrap sample (of same size) of these subjects

Everything depends on the seed:

```
set.seed(5)
```

We draw a bootstrap sample of size n:

```
bootstrap.sample <- sample(1:n, n, replace=TRUE)
```

Who is included in the bootstrap sample (look at first six)?

```
head(table(bootstrap.sample))
```

bootstrap.sample 2 3 4 5 6 8 1 1 3 1 1 3

Is subject i = 15 in this bootstrap sample?

```
15 %in% bootstrap.sample
```

[1] TRUE

Is subject i = 15 inbag or oob (out-of-bag)?

Each time a patient is left oob, we can compare the prediction for this patient with the outcome that was observed for sample patient

- model validation (which we get back to)
- variable importance measures (which we get back to)

Classical machine learning techniques utilized inside the forest

- 1. Bootstrap sampling
- 2. Nearest neighbor smoothing
- 3. Ensemble learning



Classical machine learning techniques utilized inside the forest

- 1. Bootstrap sampling
- 2. Nearest neighbor smoothing
- 3. Ensemble learning



Say we only have access to two predictors of the Epo dataset:

```
age HbBase
1 70 10.7
2 68 12.7
3 70 13.4
4 55 12.0
5 69 11.2
6 59 13.5
```

We want to estimate the probability:

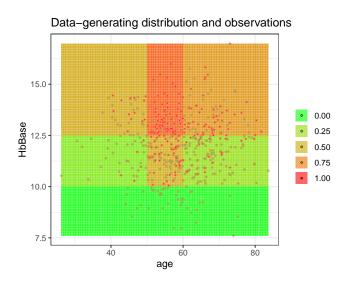
$$P(Y = 1 \mid age, HbBase)$$



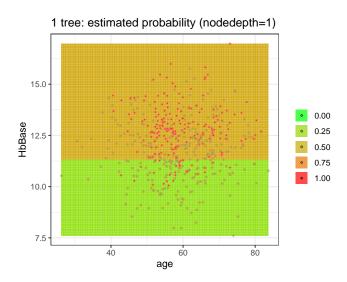
In the following, I have cheated

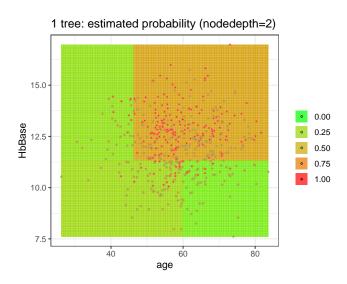
I have simulated data to imitate the Epo data

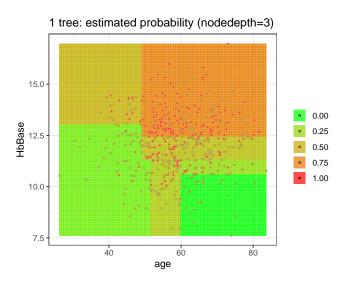
So I know the true P(Y = 1 | age, HbBase)

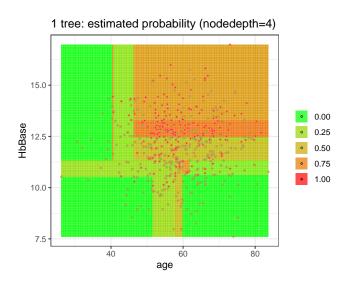


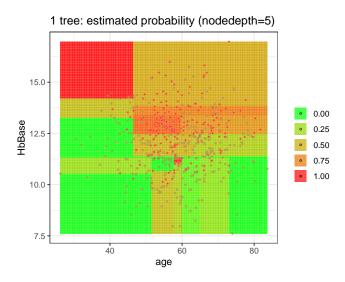
Grow a tree to fit P(Y = 1 | age, HbBase)

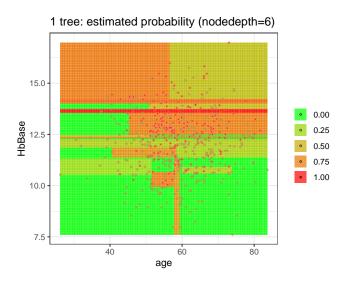


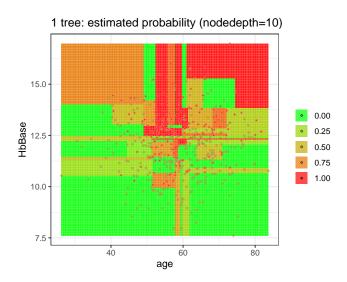




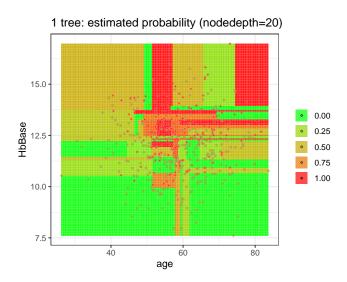








A tree is a nearest neighbor method



Classical machine learning techniques utilized inside the forest

- 1. Bootstrap sampling
- 2. Nearest neighbor smoothing
- 3. Ensemble learning

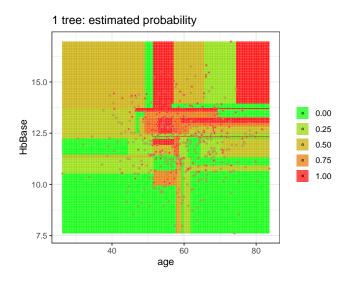


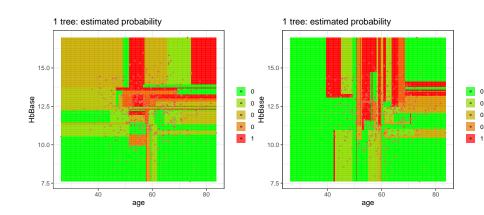
Classical machine learning techniques utilized inside the forest

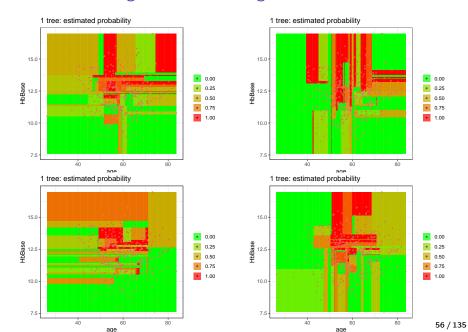
- 1. Bootstrap sampling
- 2. Nearest neighbor smoothing
- 3. Ensemble learning



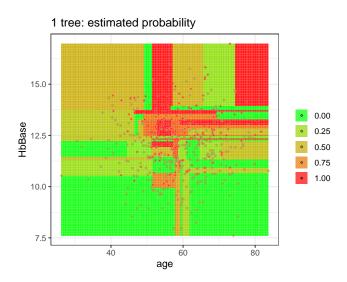
A forest takes the nearest neighbors from each tree (new tree, new seed, new bootstrap sample) to define "weighted nearest neighbors"

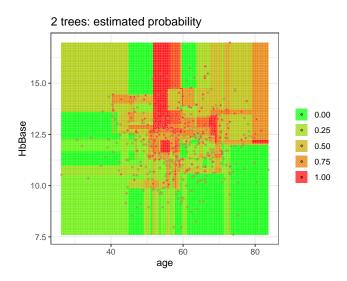


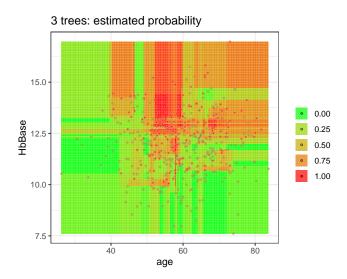


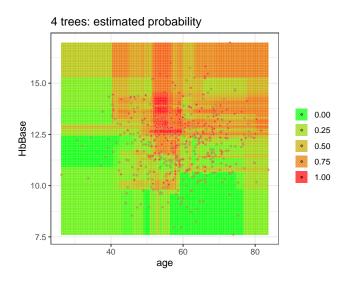


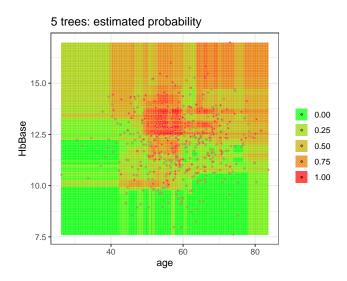
Now, combine trees to fit P(Y = 1 | age, HbBase)

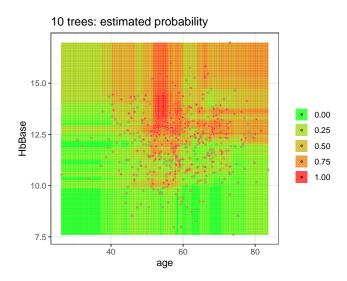


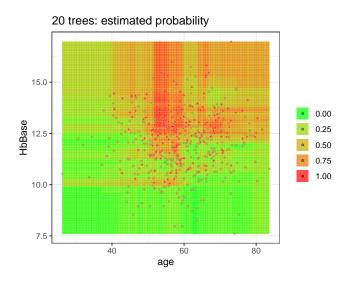


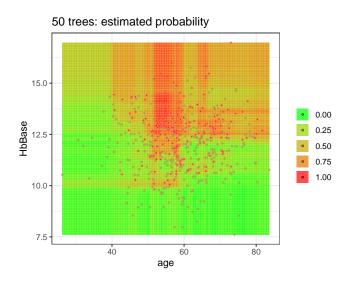


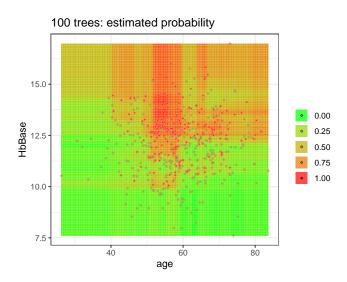


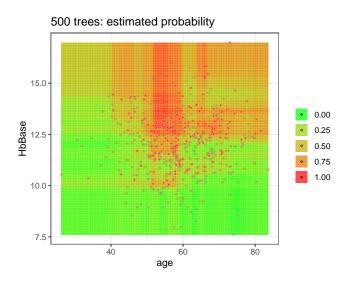












Random trees and forests in R

Load the package:

```
library("randomForestSRC")
```

We will start by using the software to grow single trees

```
tree1 <-
    rfsrc(Y~age+sex+HbBase+Treat+Resection,
    Epo, # data
    ntree=1, # only 1 tree!
    seed=1) # the result depends on seed</pre>
```

Random trees and forests in R.

Prediction and the oob prediction, e.g., for individual i = 89:

```
tree1$predicted[89]
```

[1] 0.7777778

```
tree1$predicted.oob[89]
```

[1] NA

... individual i = 89 was inbag

Exercise: From trees to forests

In this exercise, we will use the rfsrc() function from the randomForestSRC package to grow single trees

▶ The point is to assess stability of tree and forest predictions

The exercise is described in random-forest-exercises.pdf

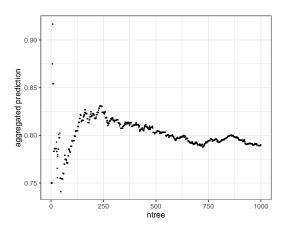
Exercise 1: From trees to forests

Exercise: From trees to forests (result plot)

```
M < -1000
pred \leftarrow rep(0, M)
for (ii in 1:M) {
    tree1 <- rfsrc(Y~age+sex+HbBase+Treat+Resection,</pre>
            Epo, ntree=1, seed=ii)
    pred[ii] <- tree1$predicted.oob[25]</pre>
pred.mean <- sapply(1:M, function(ii) {</pre>
    mean(na.omit(pred[1:ii]))
})
```

Exercise: From trees to forests (result plot)

plot(pred.mean)



Exercise: From trees to forests (remark)

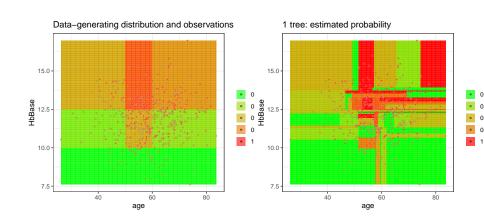
Here we produced the forest prediction ourselves across an increasing number of trees

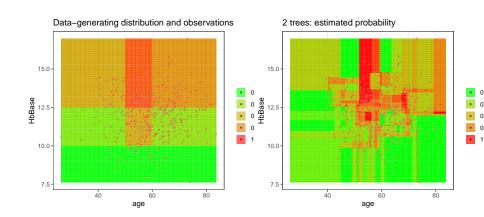
In real life we use the implementation in R to do this automatically. Here we use 1000 trees by specifying the argument ntree=1000:

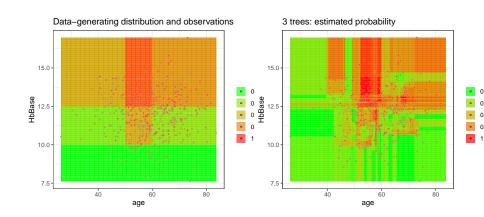
This gives us directly the forest prediction:

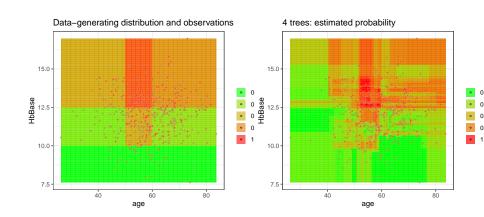
```
rf1$predicted.oob[25]
```

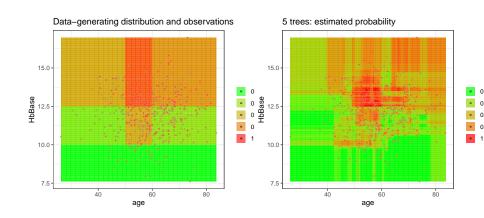
[1] 0.7528273

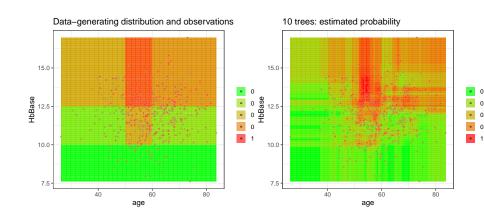


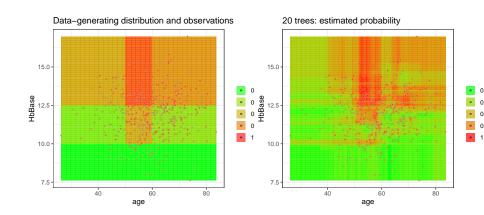


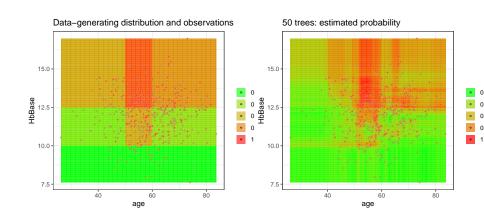


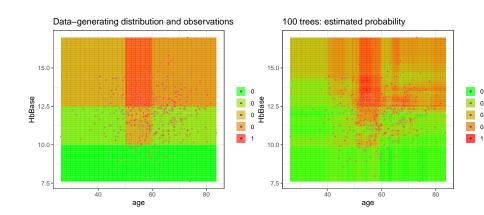


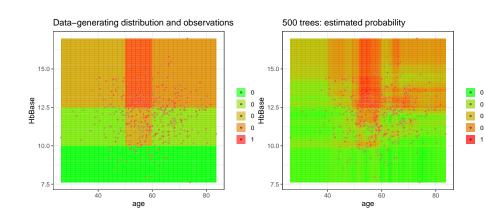












Prediction accuracy

Measuring and comparing performances of machine learning models

Outline of today's remaining topics

Part 1: Modeling cultures, model (Mark Bech Knudsen selection and decision trees

Part 2: From trees to forests, and recap on classical machine learning techniques

Part 3: Tuning random forests (Mark Bech Knudsen)

Part 4: Variable importance and (Helene Rytgaard) interpretable machine learning tools

Predictive accuracy

Combined test at 12-week pregnancy scan

 accurate prediction is important to avoid recommending unneeded invasive subsequent diagnostic test

Early detection of diabetic retinopathy

accurate prediction is important to discover as many patients as possible in time for effective treatment

Prediction of long-term survival after esophagectomy

 accurate prediction is important to correctly identify and attend to as many high-risk patients as possible

Predictive accuracy

Patient no.	Treatment successful	Predicted probability
1	0	P_1
2	0	P_2
3	1	P_3
4	1	P_4
5	0	P_5
6	1	P_6
7	1	P_7
		•
•	•	•

Prediction error is measured in terms of some distance¹² between:

- 1) the observed outcome: Y_i
- 2) and the predicted probability: $\hat{P}_i = \hat{P}(Y_i = 1 \mid age_i, HbBase_i, ...)$

One example of a loss function is the squared error loss:

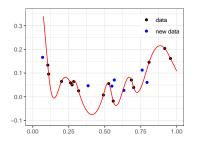
$$\mathcal{L}(Y_i,\hat{P}_i) = (Y_i - \hat{P}_i)^2$$

¹²Measured in terms of a loss function

Machine learning 101

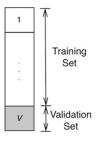
To estimate the prediction error correctly, we cannot train the model and assess the model on the same data

Overfitting happens when a model learns the detail and noise in the data too well so that it negatively impacts the performance of the model on new data



Evaluating a model on the same data results in overfitting.

To do it correctly, we can use sample splitting:



- 1. I create and fit my model on the training data: \hat{P}^{train}
- 2. I check the quality of my model on the validation data
 - Average of $\mathcal{L}(Y_i, \hat{P}_i^{\text{train}})$ in validation sample

Let's compare the predictions from 1 tree to those from a forest of 100 trees

Fix seed:

```
set.seed(5)
```

Take 10 % of original data to be our validation set:

```
val.set <- sample(1:n, n/10, replace=FALSE)</pre>
```

The rest comprise our training data:

```
train.set <- (1:n)[!(1:n) %in% val.set]
```

Fit 1 tree on the training data:

```
tree1.train <-
    rfsrc(Y~age+sex+HbBase+Treat+Resection,
    Epo[train.set,],
    ntree=1, seed=1)</pre>
```

Fit a forest of 100 trees on the training data:

```
forest.train <-
    rfsrc(Y~age+sex+HbBase+Treat+Resection,
    Epo[train.set,],
    ntree=100, seed=1)</pre>
```

Predict from the tree model on the validation set:

Predict from the forest model on the validation set:

We define the loss function:

```
loss.fun <- function(Y, Phat) mean((Y-Phat)^2)</pre>
```

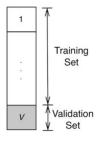
Now we can compare performance:

```
print(rbind(
    "1 tree " = loss.fun(Epo[val.set, ]$Y, tree1.val),
    "forest " = loss.fun(Epo[val.set, ]$Y, forest.val))
    )
```

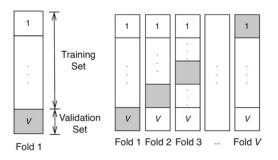
```
[,1]
1 tree 0.1443149
forest 0.0726101
```

Which one seems to perform best?

In practice, the splitting of data is not done once



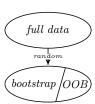
In practice, the splitting of data is not done once

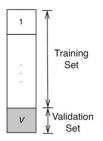


 \dots but several times. This is called V-fold cross-validation

The same technique is used inside the forest!

► The oob prediction for patient *i* only uses the trees built on bootstrap samples where patient *i* was left oob





The oob prediction error is estimated by:

$$\widehat{\mathsf{error}}_{\mathsf{oob}} = \frac{1}{n} \sum_{i=1}^{n} \mathscr{L}(Y_i, \hat{P}_i^{\mathsf{oob}})$$

```
Sample size: 149
                     Number of trees: 100
           Forest terminal node size: 5
       Average no. of terminal nodes: 12.85
No. of variables tried at each split: 2
              Total no. of variables: 5
       Resampling used to grow trees: swor
    Resample size used to grow trees: 94
                            Analysis: RF-R
                              Family: regr
                      Splitting rule: mse *random*
       Number of random split points: 10
                % variance explained: 58.37
```

Error rate: 0.1

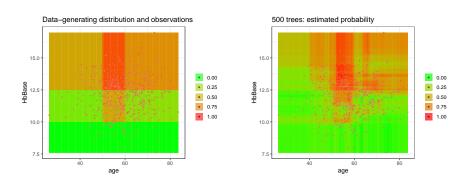
79 / 135

Picking the random forest model

Hyperparameter tuning

Picking the random forest model

The random forest algorithm automatically detects nonlinear effects, complex interactions, . . .



Picking the random forest model

But the algorithm involves some choices: hyperparameters!

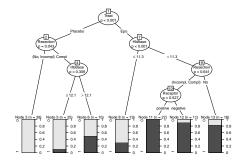
- These can be tuned and lead to different results.
- ▶ These can be tuned to optimize predictive performance

Hyperparameters of the random forest

ntree the number of trees

mtry only mtry randomly selected predictor variables are
 used to find the best split ("split-variable
 randomization")

nodesize is connected to the depth of each tree; it specifies the minimum number of observations that must be remain to perform a split



Applying a random forest

```
library("randomForestSRC")
```

Applying a random forest

We fit a random forest model on the Epo data with 1000 trees, mtry=3 and nodesize=3:

Applying a random forest

Look at predictions (here first 5):

```
rf1$predicted.oob[1:5]
```

[1] 0.029569892 0.042219020 0.971616712 0.984539768 0.00493

The oob prediction for patient 25:

```
rf1$predicted.oob[25]
```

[1] 0.7713155

Compare to what was observed for this patient:

```
Epo[25, "Y"]
```

[1] 1

Predictions on new data

newpatient

age sex HbBase Treat Resection Receptor 1 48 male 10.8 Epo No negative

Predictions on new data

Make predictions for this patient:

```
rf.pred.new <- predict(rf1, newdata=newpatient)</pre>
```

```
rf.pred.new$predicted
```

[1] 0.5763333

Exercise: Get predictions for newpatient

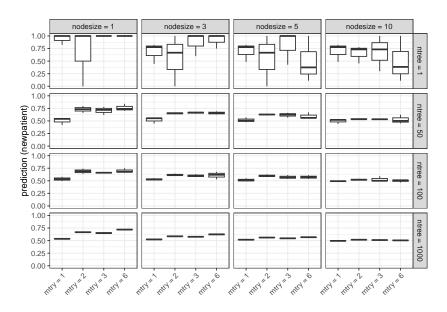
In this exercise, we will use the rfsrc() function from the randomForestSRC package to get forest prediction for newpatient for different values of hyperparameters

► The point is to see how sensitive the forest predictions are to the choice of hyperparameters

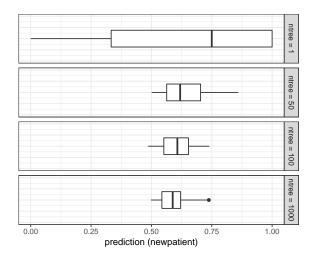
The exercise is described in random-forest-exercises.pdf

Exercise 2: Get predictions for newpatient

Exercise: Get predictions for newpatient



Exercise: Get predictions for newpatient



Let's tune the random forest = pick hyperparameters that optimize predictive accuracy

Let's tune the random forest = pick hyperparameters that optimize predictive accuracy

Look at the estimated error rate across the different choices of hyperparameters in the exercise:

```
[,1]
forest (ntree=50) 0.1229789
forest (ntree=100) 0.1225639
forest (ntree=1000) 0.1182131
forest (nodesize=3) 0.1097429
forest (nodesize=5) 0.1055143
forest (mtry=1) 0.1247416
forest (mtry=6) 0.1218394
```

Which one is the best one?

Should we tune the number of trees?

Consider the oob predicted errors for a number of different forests:

```
1 tree: 0.2199
5 trees: 0.1513
10 trees: 0.1366
50 trees: 0.1031
100 trees: 0.1074
150 trees: 0.1055
200 trees: 0.1031
500 trees: 0.1025
1000 trees: 0.1038
```

Should we tune the number of trees?

Consider the oob predicted errors for a number of different forests:

```
1 tree: 0.2199
5 trees: 0.1513
10 trees: 0.1366
50 trees: 0.1031
100 trees: 0.1074
150 trees: 0.1055
200 trees: 0.1031
500 trees: 0.1025
1000 trees: 0.1038
```

Error is minimized at ntree = 500. Better to use less trees?

Should we tune the number of trees?

One can show¹³ that increasing ntree always decreases the theoretical predictive error when using common loss functions (e.g. squared error).

Advice: do not tune ntree, just choose it "large enough" such that the accuracy has stabilized, but not so large that it becomes a computational issue.

¹³Philipp Probst and Anne-Laure Boulesteix. To Tune or Not to Tune the Number of Trees in Random Forest. Journal of Machine Learning Research. 18(181):1-18, 2018

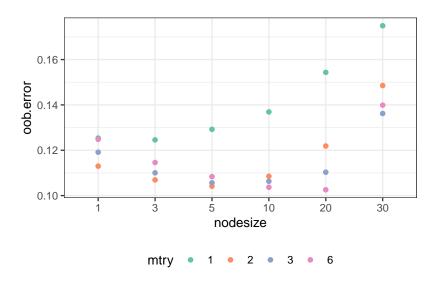
Tuning a model is **tedious work** . . . there are lot of possible combinations of the parameters

Having chosen a value for ntree, propose (relevant) combinations of values for mtry and nodesize

```
hyper.grid <- expand.grid(
  mtry = floor((ncol(Epo) - 1) / c(4, 3, 2, 1)),
  nodesize = c(1, 3, 5, 10, 20, 30),
  ntree = 1000,
  oob.error = NA
)</pre>
```

	mtra	nodogiza	ntroo	ooh orror
	шсту			oob.error
1	1	1	1000	NA
2	2	1	1000	NA
3	3	1	1000	NA
4	6	1	1000	NA
5	1	3	1000	NA
6	2	3	1000	NA
7	3	3	1000	NA
8	6	3	1000	NA
9	1	5	1000	NA
10	2	5	1000	NA
11	3	5	1000	NA
12	6	5	1000	NA
13	1	10	1000	NA
14	2	10	1000	NA
15	3	10	1000	NA
16	6	10	1000	NA
17	1	20	1000	NA

Compute the oob error for all combinations:



Tuning hyperparameters

Which combination gave the lowest estimated error rate?

```
hyper.grid[which.min(hyper.grid$oob.error),]
```

```
mtry nodesize ntree oob.error
20 6 20 1000 0.1025708
```

Tuning hyperparameters

Let's fit the corresponding, now tuned, forest:

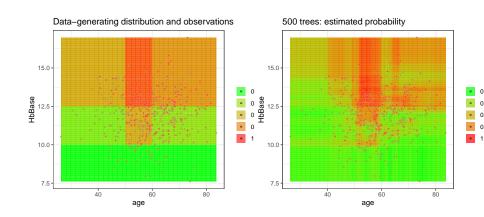
```
j <- which.min(hyper.grid$oob.error)</pre>
```

```
j
```

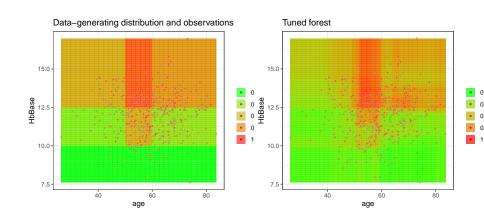
[1] 20

```
tuned.rf <-
    rfsrc(Y~age+sex+HbBase+Treat+Resection,
        Epo,
        mtry=hyper.grid[j, "mtry"],
        nodesize=hyper.grid[j, "nodesize"],
        ntree=hyper.grid[j, "ntree"], seed=1)</pre>
```

Tuning hyperparameters for the simulated data



Tuning hyperparameters for the simulated data



Interpretable machine learning

Variable importance

Outline of today's remaining topics

Part 1: Modeling cultures, model (Mark Bech Knudsen selection and decision trees

Part 2: From trees to forests, and (Helene Rytgaard) recap on classical machine learning techniques

Part 3: Tuning random forests (Mark Bech Knudsen)

Part 4: Variable importance and (Helene Rytgaard) interpretable machine learning tools

Logistic regression

Response: treatment successful yes/no

Covariate	${\sf OddsRatio}$	CI.95	pValue
(Intercept)	0.00		0.0040
age	0.97	[0.91; 1.03]	0.2807
sexmale	0.21	[0.038; 1.10]	0.0657
HbBase	3.26	[1.99; 5.91]	< 0.0001
TreatPlacebo	0.01	[0.0020; 0.042]	< 0.0001
ResectionIncompl	0.42	[0.083; 1.96]	0.2801
ResectionNo	0.24	[0.058; 0.89]	0.0395
Receptorpositive	5.81	[1.72; 23.39]	0.0076

Machine learning

The Algorithmic Modeling Culture

The analysis in this culture considers the inside of the box complex and unknown. Their approach is to find a function $f(\mathbf{x})$ —an algorithm that operates on \mathbf{x} to predict the responses \mathbf{y} . Their black box looks like this:



Model validation. Measured by predictive accuracy. Estimated culture population. 2% of statisticians, many in other fields.

Interpretable machine learning

- Decision trees produce results that are easy to interpret
- Random forest results, on the other hand, are not per se so easy to interpret



- What predictor variables were important for the prediction?
- What effect did the predictor variables have on the prediction?

How important was a given variable for building the forest model?

We consider two different approaches

- 1. "VIMP"
- 2. Minimal depth

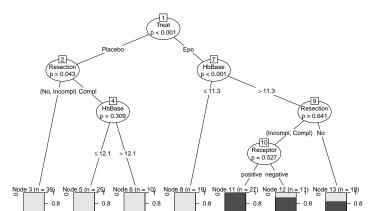
VIMP (Variable IMPortance) is measured by the difference prediction error between:

- running the forest with a "noised-up" version of X
- running the forest with X as was observed

If prediction performance decreases more for variable X_1 than for variable X_2 , then importance(X_1) > importance(X_2)

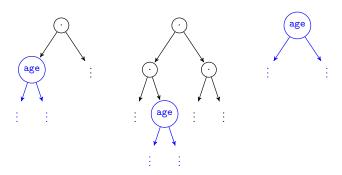
Recall:

- Trees are built by recursive partitioning
- They let the data decide which variables are important for splitting node



110 / 135

The minimal depth is the average distance from the root node to the first split on a specific variable



The smaller the minimal depth, the more important is the variable

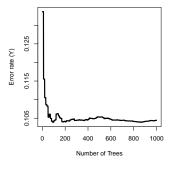
VIMP for the Epo data:

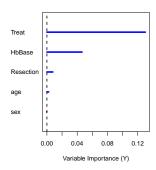
```
tuned.rf$importance
```

```
age sex HbBase Treat Resection 0.07420322 0.00589100 0.23239730 0.36622336 0.02681411
```

What variables are most important?

```
plot(tuned.rf)
```





Minimal depth for the Epo data:

age	sex	HbBase	Treat	Resection
1.550	3.688	1.300	0.992	2.072

The forest provides a threshold (cut-off) value:

[1] 2.343143

What variables are important?



What you should not do...



What you should not do...

```
mtry=1: mtry1
age 1.688
sex 2.017
HbBase 1.697
Treat 1.867
Resection 1.903
Threshold=
```

[1] 2.051348



What you should not do...

mtry=1:		mtry=2:	
-	mtry1	-	mtry2
age	1.688	age	1.550
sex	2.017	sex	3.688
HbBase	1.697	HbBase	1.300
Treat	1.867	Treat	0.992
Resection	1.903	Resection	2.072
Threshold= Threshold		Threshold=	
[1] 2.0513	348	[1] 2.343143	

Variable importance measures from random forests

Why is this?

mtry controls split-variable randomization:

- for each node only a small number of randomly selected predictors are used to find the best split of that node (= mtry)
- this is done as part of the randomization of trees
- (it ensures some of the theoretical properties of the forests)

In fact, if we are interested in variable importance (rather than predictive accuracy) we should choose a high value for this.

Variable importance on simulated data

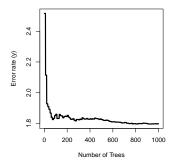
- Two uncorrelated variables x1 and x2 with the same effect
- One variable c1 correlated with x1 but with no effect
- ▶ Two correlated variables z1 and z2 with the same effect
- ► Ten noise variables w1,..., w10

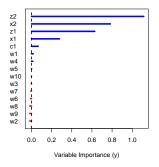
```
x1 <- runif(n)
x2 <- runif(n)
z1 <- rnorm(n, mean=0, sd=0.3)
z2 <- rnorm(n, mean=z1+0.1, sd=0.3)
c1 <- rnorm(n, mean=x1+0.1, sd=0.3)
w <- matrix(runif(n*10), ncol=10)
y <- rnorm(n, mean=0.1+2.5*x1+2.5*x2+2.5*z1+2.5*z2)</pre>
```

Variable importance on simulated data

Variable importance on simulated data

- Two uncorrelated variables x1 and x2 with the same effect
- One variable c1 correlated with x1 but with no effect
- ▶ Two correlated variables z1 and z2 with the same effect
- ► Ten noise variables w1,..., w10





Exercise: Identifying risk factors with variable importance

In this exercise we will look at the analysis of Hsich et al. (2011):

Identifying Important Risk Factors for Survival in Patient With Systolic Heart Failure Using Random Survival Forests

Eileen Hsich, MD; Eiran Z. Gorodeski, MD, MPH; Eugene H. Blackstone, MD; Hemant Ishwaran, PhD; Michael S, Lauer, MD

Background—Heart failure survival models typically are constructed using Cox proportional hazards regression. Regression modeling suffers from a number of limitations, including bias introduced by commonly userd variable selection methods. We illustrate the value of an intuitive, robust approach to variable selection, random survival forests (RSF), in a large clinical cohort. RSF are a potentially powerful extensions of classification and regression trees, with lower variance and bias.

Methods and Results—We studied 2231 adult patients with systolic heart failure who underwent cardiopulmonary stress testing. During a mean follow-up of 5 years, 742 patients died. Thirty-nine demographic, cardiac and noncardiac comorbidity, and stress testing variables were analyzed as potential predictors of all-cause mortality. An RSF of 2000 trees was constructed, with each tree constructed on a bootstrap sample from the original cohort. The most predictive variables were defined as those near the tree trunks (averaged over the forest). The RSF identified peak oxygen consumption, scrum urea nitrogen, and treadmill exercise time as the 3 most important predictors of survival. The RSF predicted survival similarly to a conventional Cox proportional hazards model (out-of-bag C-index of 0.705 for RSF versus 0.698 for Cox proportional hazards model).

Conclusions—An RSF model in a cohort of patients with heart failure performed as well as a traditional Cox proportional hazard model and may serve as a more intuitive approach for clinicians to identify important risk factors for all-cause mortality. (Circ Cardiovasc Outl Outcomes, 2011:4:39-45.)

Key Words: heart failure ■ prognosis ■ statistics ■ survival analyses

The exercise is described in random-forest-exercises.pdf

Exercise 3: Identifying risk factors

Effects of predictor variables on the final prediction

Plots can be useful to assess the effect of predictor variables on the final prediction.

Effects of predictor variables on the final prediction

Plots can be useful to assess the effect of predictor variables on the final prediction. There are different ways to do so:

Partial Dependence Plots (PDPs)

- Average forest predictions as a function of predictor variables
- Obtained by marginalizing the forest prediction over the other features/covariates
- Can show if the relationship is linear, monotonic or more complex

Individual Conditional Expectation (ICE) plots

 Looking at the individual predictions as a function of predictor variables

Say, we want to know how

$$\hat{P}(Y = 1 \mid age, Gender, HbBase, Treatment, Resection)$$

varies when HbBase varies

We can estimate this by:

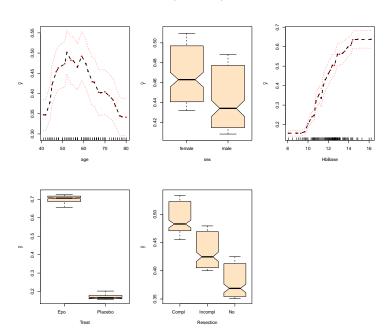
$$\hat{P}^{\mathsf{HbBase}}(b) = \frac{1}{n} \sum_{i=1}^{n} \hat{P}(Y_i = 1 \mid \mathsf{age}_i, \mathsf{Gender}_i, \mathsf{HbBase} = b,$$

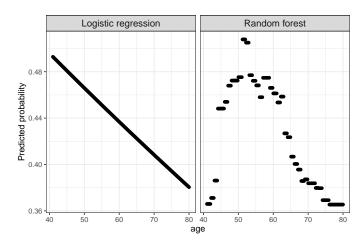
$$\mathsf{Treatment}_i, \mathsf{Resection}_i)$$

 We marginalize the forest prediction over the other features/covariates

In R, we can plot these estimates for all variables by simply writing:

```
plot.variable(tuned.rf, partial=TRUE, plots.per.page=3)
```





It is clear that the random forest captures a highly nonlinear effect of age on the predicted probability!

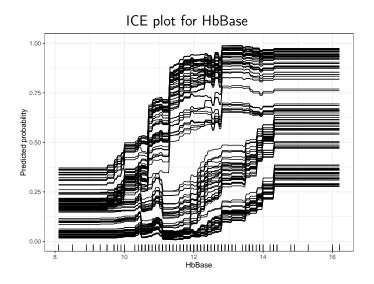
The ICE plot shows the variation of

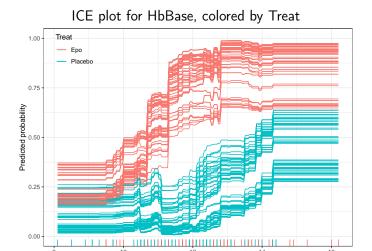
$$\hat{P}_i^{\mathsf{HbBase}}(b) = \hat{P}(Y = 1 \mid \mathsf{age}_i, \mathsf{Gender}_i, \mathsf{HbBase} = b,$$

$$\mathsf{Treatment}_i, \mathsf{Resection}_i)$$

for each individual i one by one.

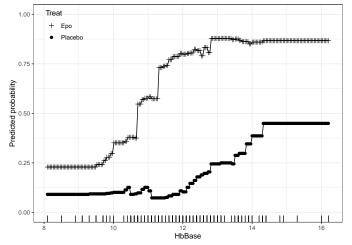
- This can very useful if there are interactions
- ▶ Do the curves follow the same course (e.g., changepoints, linearity, etc) for all individuals?
 - OBS: Intercept differences are not a sign of interactions

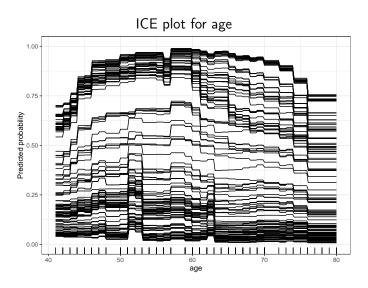


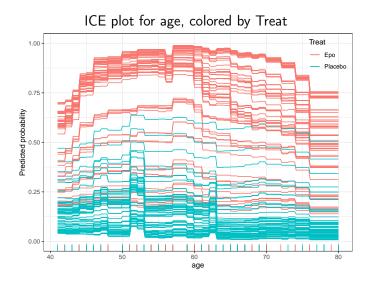


HbBase

PDP plot for HbBase, computed in groups defined by Treat







Interpretable machine learning

Variable importance measures can tell us what variables seem important for prediction

Beware of correlated predictors

PDPs and ICE plots can show us how predicted probabilities vary as a function of predictor values

- ▶ PDPs show the average variation
 - Beware of hidden interactions
- ICE show the individual variations

Logistic regression versus random forests







Logistic regression versus random forests

When utilizing a logistic regression approach:

We must specify the model¹⁴

$$\begin{split} \hat{P}(Y = 1 \mid \mathsf{age}, \mathsf{Gender}, \mathsf{HbBase}, \mathsf{Treatment}, \mathsf{Resection}, \mathsf{Epo}) \\ &= \mathsf{expit}(\beta_0 + \beta_1 \mathsf{age} + \beta_2 \mathsf{age} : \mathsf{female} + \cdots) \end{split} \tag{*}$$

Based on the model we may:

- Predict $\hat{P}(Y=1)$ for a new patient
- Interpret odds ratios, conditional on holding the other features fixed (p-values, confidence intervals, etc)

But all inference relies on (*) being correct and prespecified

¹⁴Interactions, quadratic terms (e.g., age²), ...

Logistic regression versus random forests

When utilizing a random forest approach:

- The forest automatically detects nonlinear effects and complex interactions
- "Model selection" is comprised by hyperparameter tuning

Based on the model we may:

- ▶ Predict $\hat{P}(Y = 1)$ for a new patient often with high accuracy
- Obtaining interpretable measures from the random forest¹⁵ is applied after model training, e.g.:
 - ▶ Variable importance, PDPs, ICEs, ...

But, inference (confidence intervals, p-values) is not so obvious And, beware that everything depends on the random seed

¹⁵And other machine learning methods

Exercise: Predicting tumor class (Golub et al., 1999)

Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

T. R. Golub, ^{1,2*}† D. K. Slonim, ¹† P. Tamayo, ¹ C. Huard, ¹ M. Gassenbeek, ¹ J. P. Mesirov, ¹ H. Coller, ¹ M. L. Loh, ² J. R. Downing, ³ M. A. Caligiuri, ⁴ C. D. Bloomfield, ⁴ E. S. Lander^{1,3}*

Although cancer classification has improved over the past 30 years, there has been or general approach for identifying new cancer classes (class adiscovery) been or general paperad for identifying new cancer classes (class prediction). Here, a generic paper classification based on gene expression monitorion), and the control of the co

- Accurate cancer classification can be used to target specific therapies to distinct tumor types
- We could use a random forest model to provide a data-based classification algorithm based on gene expression monitoring

Exercise: Predicting tumor class (Golub et al., 1999)

In this practical we will work with a dataset containing information on 38 tumor mRNA samples from 38 individuals and the gene expression values from 3051 genes

We will go through the steps on the lectures slides to explore these data

▶ The goal of the analysis is to predict the tumor class

The exercise is described in random-forest-exercises.pdf

Exercise 4: Predicting tumor class

That was it ...

Comments and suggestions for this material are very much welcome at hely@sund.ku.dk \circledcirc