

LUDWIG-MAXIMILIANS-UNIVERSITY
MUNICH

INSTITUTE FOR STATISTICS

Master Thesis

A comparison study of prediction
approaches for multiple training
data sets and test data with
block-wise missing values

Author:

Frederik
LUDWIGS

Supervisor:

Dr. Roman
HORNUNG



April 14, 2020

Abstract

Contents

List of Figures	3
List of Tables	4
1 Introduction	5
2 Methods	8
2.1 Block-wise missingness in multi-omics data	8
2.2 Random Forest for classification	11
2.2.1 Decision Tree	11
2.2.2 Random Forest Model	16
2.2.3 OOB error	17
2.2.4 Variable importance	19
2.3 Complete-Case Approach	21
2.4 Single-Block Approach	23
2.5 Imputation Approach	25
2.6 Block-wise Approach	25
2.7 Fold-wise Adaption	30
3 Benchmark Experiments	35
3.1 Data	35
3.1.1 Own data	35
3.1.2 Data from Hagenberg	35
3.1.3 Real data	35
3.2 Accessing the Performance	35
3.2.1 CV - Test-Situations	35
3.2.2 Metrics	35
4 Results	36
4.1 Own data	36
4.1.1 Scenario 1	36
4.1.2 Scenario 2	36
4.1.3 Scenario 3	36
4.1.4 Scenario 4	36
4.2 Data from Hagenberg	36
4.3 Real data	36
5 Discussion and Conclusion	37

6 Bibliography	38
7 Attachment	41

List of Figures

1	Example for block-wise missingness, when concatenating data from diverse sources.	10
2	Example for the recursive binary splitting of decision tree on a two-dimensional feature space.	14
3	Corresponding decision tree for the segmented feature space on the rightmost plot in figure 2.	15
4	The data used to grow the M different decision trees of a random forest. Below each decision tree the in-bag observations are labelled in pink, while out-of-bag observations are labelled in grey [[31], p. 13].	18
5	Calculation of the variable importance of x_1 for a random forest model consisting of M decision trees [[31], p. 16].	20
6	Two examples for the 'Complete Case' processing of the training data according to the available feature-blocks in the test set.	22
7	The 'Complete Case' processing results in a empty training set, such that no model can be trained. In these settings, predictions can not be generated for the test set with the 'Complete Case' approach.	23
8	'Single-Block' processing of the training data so a random forest model can be regularly trained with each of these processed data sets.	25
9	Example for the training of random forest models with the 'block-wise' approach.	27
10	Example for the prediction on test data with the 'block-wise' approach. The fitting of the random forest models was described with figure 9.	28
11	Example for the training of random forest models with the 'fold-wise' approach.	31
12	Example for the pruning of a single decision tree. The decision tree was originally introduced in figure 3	32
13	Example for the prediction on test data with the 'fold-wise' approach. The training of the models was described with figure 11.	33
14	The different impurity functions (3), (4) and (5) plotted for a given fraction of a binary target variable within any node N [[27], p. 13]	41

List of Tables

- | | | |
|---|--|---|
| 1 | Example for a data set with block-wise missingness - consisting of three blocks, three folds and the binary target variable 'Y'. | 9 |
|---|--|---|

1 Introduction

On October 1, 1990 the international scientific research project named *Human Genome Project* was launched, with the aim to sequence the first complete human genome ever [1]. After investments of totally \$2.7 billion and 13 years of research the sequencing was officially finished in 2003 [2]. Since then, on the one hand, there have been biomedical advances that have led to the identification of disease genes, “leading to improved diagnosis and novel approaches in therapy” [[3], p. 14]. On the other hand there has been an “extraordinary progress [...] in genome sequencing technologies [itself]” [[4], p. 333], leading to a sharp drop in sequencing prices. Nowadays whole genome sequencing is available and affordable for almost everyone - e.g. ‘Veritas Genomics’, offers whole genome sequencing for \sim \$700 [5].

Besides the ‘genome’ that carries the whole genetic material of an organism, there are also other types of ‘-omes’, such as ‘epigenomes’, ‘transcriptomes’, ‘proteomes’ and ‘microbiomes’. The time and costs to collect data from these different types of ‘-omes’ have been reduced drastically ever since the completion of the Human Genome Project [[6], [7], [8], [9], [10], [11]]. The methods for “fast, automated analyses of large numbers of substances including DNA, RNA, proteins, and other types of molecules” [12] are summarized under the term ‘High Throughput Technologies’. These technologies made data from molecular processes available for many patients on a large scale.

The collected data from any type of ‘-omes’ is commonly referred to as ‘omics data’. In the clinical context it is of utmost interest to incorporate such omics data into different statistical approaches. A common example in this context is the survival time prediction for cancer patients, where in addition to regular clinical data (e.g. ‘weight’, ‘height’, ‘blood pressure’), gene expression data has been incorporated into the survival models. This additional omics data has “often been found to be useful for predicting survival response[s]” [[7], p. 1]. In “the beginning, only data from single omics was used to build such prediction models, together or without [...] clinical data” [[13], p. 1]. The usage of multiple distinct types of ‘-omes’ in a single prediction approach was the next logical step and coined the term ‘multi-omics data’. The theoretical aspects of integrating multiple omics types into a single prediction model and how to deal with the block-wise structures have been topic of several papers already - e.g. [13], [14], [15], [16], [17].

This thesis deals with a special type of missing data “that is common in practice, particular in the context of multi-omics data” [18] - the so called block-wise missingness. Data with block-wise missingness consists of different folds and feature-blocks. While a feature-block stands for a collection of

associated covariates, a fold represents a set of observations with the same observed feature-blocks. In data sets with block-wise missingness there is always at least one fold, where at least one feature-block is missing, so that not all observations have the same observed blocks.

Most statistical methods require fully observed data for training and predictions. In data with block-wise missingness this requirement is clearly not met, so that either the approaches need methodical adjustment or the data needs to be processed. This foundational problem with block-wise missingness raises the following challenges and questions: How can we fit a model on the block-wise missing data, without removing observations or whole feature-blocks? How does a model that uses complete cases only perform in comparison? Does imputation work properly in these settings? How does a model that uses single feature-blocks only perform in comparison? How can a model predict on observations with missing feature-blocks?

In addition to the problem of block-wise missingness, there is also the challenge of “inherent high dimensionality” [[19] p. 93], when working with multi-omics data. Data from a single omics type can easily exceed thousands of covariates and the corresponding data sets usually consist of less observations than features [13]. Besides the predictive performance of an approach it is furthermore important for the approach to be sparse. “Sparsity is [...] an important aspect of the model which contributes to its practical utility” [[15], p. 3], as it makes the model much more interpretable than models including several thousands of variables.

A method that handles high dimensional data, even if the number of observations is lower than the amount of features, is the random forest method [13]. Additionally the method handles different input types, does not need a lot of tuning and yields comparable predictive performances [20]. The only drawback is that it is not as interpretable as “models yielding [*in*] coefficient estimates of few relevant features” [[13], p. 35], as penalised regression approaches for example. Nevertheless variable importance measures can be extracted with the random forest method, as well as partial dependencies. Furthermore it has been used successfully in various articles dealing with multi-omics data - e.g. [13], [14]. Moreover there have been proposals by Hornung et al. [18] and Krautenbacher [19] that modify the random forest approach, hence it can directly handle block-wise missing data.

The different adaptations of penalised regression, as for example the IPF- & Priority-Lasso [[7], [15]] can be modified so they can directly deal with block-wise missing data. The theoretical aspects of these approaches are not part of this thesis, but of Hagenberg’s [21]. Nevertheless the performances of the different random forest approaches and the penalised regression adaptations are compared in this thesis as well.

Even though the problem of block-wise missingness is common in multi-omics data there are, to my knowledge, no comparison studies of such prediction approaches yet. Krautenbacher has already stated that “reliable analysis strategies for multi-omics data [...] [with block-wise missingness are] urgently needed” [[19] p. 94]. The thesis at hand aims to provide such a large scale comparison study of prediction approaches capable of dealing with block-wise missingness and shall help finding a reliable strategy.

To investigate a reliable analysis strategy for multi-omics data with block-wise missingness this paper compares the predictive performance of two naive random forest approaches, a random forest approach on imputed data, two random forest adaptations and the adaptations of penalised regression. In the second chapter, firstly the term ‘block-wise missingness’ is defined in more detail and how it can arise in multi-omics data. Then a theoretical explanation of the random forest method for classification is given. Following three data processing approaches are explained - these process the block-wise missing data such that a regular random forest can be trained with it. Moreover two methodological adaptations of the random forest method are illustrated. These adaptations allow the random forest approach to directly deal with block-wise missing data and do not rely on any data processing. In the third chapter the different data sources used to validate the performances of the various approaches are described and investigated. The metrics and tactic used for the evaluation of the models are given afterwards. In the penultimate chapter all approaches are analysed and compared. In the last chapter of this thesis all findings are discussed, conclusions are drawn and an outlook is given.

2 Methods

This section deals with the theory of the random forest model and the different adaptations of it to handle data with block-wise missingness.

In the beginning block-wise missingness is defined in more detail and it is shown how it can arise in multi-omics data. Afterwards the theory of the random forest method for classification is illustrated. Subsequent three approaches that process the data with block-wise missingness, such that a regular random forest can be fit on them, are described. The last two sections of this chapter present two different adaptations of the random forest method. These adaptations enable the random forest method to directly deal with block-wise missing data.

2.1 Block-wise missingness in multi-omics data

Collecting omics data has become significantly cheaper and faster ever since the completion of the Human Genome Project. As a result, this type of data is used more and more frequently in the biomedical research - e.g. risk prediction of childhood asthma [19]. Even though the integration of multiple types of '-omes' into a single prediction approach seems promising there are still challenges to face. One of these challenges is a special type of missingness that is common in the context of multi-omics data, the so called block-wise missingness [18].

The term block-wise missingness needs to be defined in more detail, before clarifying how it can arise in multi-omics data. Table 1 shows a minimalist example for a data set with block-wise missingness, whereby the data consists of eight observations and 105 covariates in total. While the covariates 'weight', 'height', 'income' and 'education' are pretty much self-explanatory, the features ' g_1 ', ..., ' g_{100} ' could be any type of omics data. Data sets with block-wise missingness always consist of different blocks and folds. On the one hand, a **block** describes a set of covariates containing all features collected on the basis of a characteristic - basically all covariates that are related in content. The example data in table 1 has three blocks in total. 'Block 1' consists of the variables 'weight' and 'height' representing the physical properties. 'Block 2' contains the variables 'income' and 'education' standing for economic properties. 'Block 3' includes the remaining variables ' g_1 ', ..., ' g_{100} ' and could represent **biological/ genetic properties**. On the other hand, a **fold** represents a set of observations with the same observed feature-blocks - basically all observations with the same observed features. The data set in table 1 consists of three folds in total. 'Fold 1' holds the



observations 1, 2 and 3, as these have the same observed feature-blocks ('Block 1' & 'Block 2'). 'Fold 2' holds the observations 4 and 5, while 'Fold 3' consists of the remaining observations 6, 7 and 8. As each fold has different observed feature-blocks, each fold is unique and every observation belongs to exactly one of them. The only variable all folds must have in common is the target variable - 'Y' in table 1.

<i>ID</i>	<i>weight</i>	<i>height</i>	<i>income</i>	<i>education</i>	g_1	\dots	g_{100}	<i>Y</i>	
1	65.4	187	2.536	<i>Upper</i>				1	} Fold1
2	83.9	192	1.342	<i>Lower</i>				0	
3	67.4	167	5.332	<i>Upper</i>				1	
4			743	<i>Lower</i>	-0.42	\dots	1.43	1	} Fold2
5			2.125	<i>Lower</i>	0.52	\dots	-1.37	0	
6	105.2	175			-1.53	\dots	2.01	0	} Fold3
7	71.5	173			0.93	\dots	0.53	0	
8	73.0	169			0.31	\dots	-0.07	1	
<i>Block1</i>			<i>Block2</i>		<i>Block3</i>				

Table 1: Example for a data set with block-wise missingness - consisting of three blocks, three folds and the binary target variable 'Y'.

Multi-omics data with block-wise missingness have a structure as displayed in table 1, but the single feature-blocks are usually much higher dimensional than in the given example. When working with multi-omics data this type of missingness is a common problem. There are two main reasons for this: The first one is related to the costs of collecting omics data. Even though these have been reduced drastically over the last 15 years, collecting omics data is still more complex and expensive than obtaining clinical data for example. As a consequence, due to financial or even technical constraints, omics data can not always be collected for all participants of a study. Therefore participants from the same study can end up with different observed feature-blocks, such that the data for the whole study contains block-wise missingness.

The second reason is related to the collection of training sets from different sources - e.g. various hospitals. Even though the different sources do research regarding the same target variable - e.g. person has asthma (yes/ no) - the surveyed feature-blocks can still differ. Therefore the concatenation of such data sets can easily result in a data set with block-wise missingness. This scenario is illustrated in figure 1. In the top of the figure the three different data sources (Hospital 1, \dots , Hospital 3) are displayed. Each source consists

of the target variable 'Y' and two feature-blocks - e.g. 'Hospital 2' consists of the target variable 'Y' and of the feature-blocks 'RNA' and 'Clinical'. The feature-blocks 'RNA', 'MIRNA' and 'CNV' represent high dimensional omics data, while the 'Clinical' feature-block stands for several clinical features. Even though the target variable 'Y' is the same in all sources, the collected feature-blocks in these still differ. The concatenation of these data sets - bottom of figure 1 - results in data with block-wise missingness. In the concatenated data an observed block is marked with a tick and a missing block with a cross. The fold 'Hospital 2' only has 'RNA' and 'Clinical' as observed feature-blocks, so that the observations from this fold miss all the features from the blocks 'CNV' and 'MIRNA'. The concatenated data totally consists of three unique folds and four different feature-blocks.

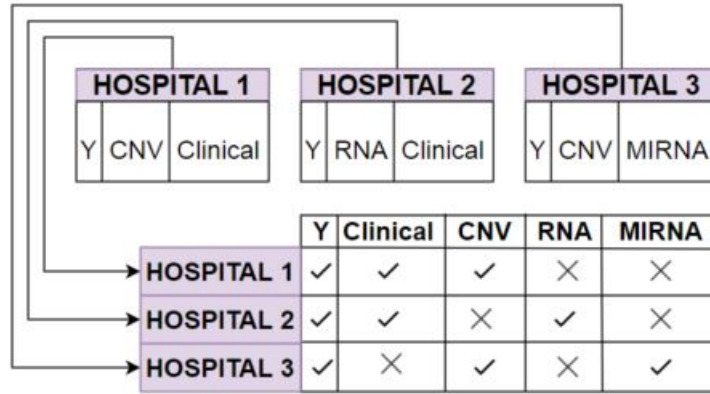


Figure 1: Example for block-wise missingness, when concatenating data from diverse sources.

Training a prediction model, as for example a random forest, is not directly possible on data with block-wise missingness. Either the methods have to be adopted or the data processed. As block-wise missingness can also affect the test data it raises the following question. How can a model do a prediction for a observation that misses feature-blocks the model has originally been trained with? This challenge has to be taken into account when proposing methods capable to deal with block-wise missingness.

The remaining sections in this chapter focus on approaches and adaptations to make model fitting on such data possible. Firstly the concept of the random forest for classification is explained and then the different approaches and adaptations to handle data with block-wise missingness.

2.2 Random Forest for classification

This chapter illustrates the random forest method that has already been applied in several articles dealing with multi-omics data [[13], [14], [15]]. It is a “powerful prediction method [...] able to capture complex dependency patterns between the outcome and the covariates” [[14] p. 2]. Furthermore it does not need a lot of tuning and naturally handles high-dimensional data, with more covariates than observations [13]. The random forest method can be applied to classification-, regression- and even survival-problems. Latter was added in 2008 by Ishwaran et al. [22]. As this thesis focuses on classification tasks, only the random forest for classification is explained. Nevertheless all of the approaches and adaptations described in sections 2.3 to 2.7 can also be applied for non-categorical target variables.

The random forest is a tree-based ensemble method that was introduced by Breiman in 2001 [23]. An ensemble is a concept from machine learning that “train[s] multiple models using the same learning algorithm” [24]. Therefore an ensemble consists of η identical so called base learners. The base learner of the random forest method is a decision tree. This is an excellent base learner for an ensemble, as a decision tree can capture complex interactions and have a relatively low bias, if grown sufficiently deep. Especially as single decision trees are known to be noisy, they benefit from the ensemble [20]. Since decision trees are the basis of random forest method it is crucial to understand how these work in order to properly understand the random forest method.

2.2.1 Decision Tree

A decision tree is a supervised learning method that was introduced by Breiman et al. in 1984 [25]. It has a hierarchical nature, is easy to interpret and non-model based [26]. It applies recursive binary splitting to “partition the feature space into a set of rectangles” [[20] p. 305], so that the resulting rectangles are as pure as possible in terms of the target variable. A prediction is generated by assigning an observation to one of the rectangles in the partitioned feature space. The prediction then equals the distribution of the target variable within the assigned rectangle - e.g. a observation that falls into a rectangle with three negative and seven positive responses has predicted probability of 70% for a positive response.

To partition the feature space into the purest rectangles possible the algorithm iterates over all possible split variable/ split value combinations. For each of these possible splits, the observations from the parent node N are divided - with respect to the split variable x_j at split point t - into the

child nodes N_1 and N_2 [[27], p. 10]:

$$N_1(x_j, t) = \{(x, y) \in N : x_j \geq t\} \quad (1)$$

$$N_2(x_j, t) = \{(x, y) \in N : x_j < t\} \quad (2)$$


N_1 therefore contains all observations from the parent node N with $x_j \geq t$, while N_2 contains all observations from the parent node N with $x_j < t$. The point (x_j, t) therefore creates a binary split and partitions the data from parent node N in the two subspaces N_1 and N_2 . The split variable x_j and split point t are chosen such that the resulting child nodes N_1 and N_2 have the greatest possible purity [27]. To measure the impurity of a node N regarding a categorical response with g classes the 'Gini-Index' (3), 'Missclassification-Error' (4) or 'Shannon-Entropy' (5) can be used [[27], p. 12]:

$$I(N) = \sum_{k=1}^g \hat{\pi}_{k,N} \cdot (1 - \hat{\pi}_{k,N}) \quad (3)$$

$$I(N) = 1 - \max_k \hat{\pi}_{k,N} \quad (4)$$

$$I(N) = - \sum_{k=1}^g \hat{\pi}_{k,N} \cdot \log(\hat{\pi}_{k,N}) \quad (5)$$

- $\hat{\pi}_{k,N}$: Relative frequency of category k in node N

For all of these impurity measures applies: The lower $I(N)$ the purer the node N and a node N is completely pure, when it only contains observations of the same response class $\rightarrow I(N) = 0$. The plotted impurity functions for a binary target variable can be seen in figure 14 in the attachment. 

The reduction of the impurity when splitting the parent node N into the child nodes N_1 and N_2 is calculated by [[27], p. 10]:

$$I(N) - \frac{|N_1|}{|N|} \cdot I(N_1) - \frac{|N_2|}{|N|} \cdot I(N_2) \quad (6)$$

- $|N|$: Number of observations in the parent node N
- $|N_1|$: Number of observations in child node N_1
- $|N_2|$: Number of observations in child node N_2

This equation basically calculates how strong the impurity from the parent node N is reduced for a given splitting point that divides the observations to the child nodes N_1 and N_2 . This impurity reduction is calculated for every possible split. The final split variable x_j and split point t are chosen, such that the impurity is maximally reduced.

For illustrative purposes the single partition steps of a classification tree are shown in figure 2. The figure consists of three plots in total, whereby each is a scatter plot of 'weight' and 'height' for the observations from 'Fold 1' and 'Fold 3' in table 1. Observations with a positive outcome are marked in blue, while negative outcomes are marked in red.

In the very beginning all observations are within the same feature space that has not been divided yet - the so called root node. This is displayed in the leftmost plot of figure 2. The root node contains three observations with a positive and three with a negative response - hence the class distribution in this node is 50|50. The node is not pure regarding its responses and all possible impurity measures [(3), (4), (5)] are rather high. The algorithm now iterates over all features, and for each feature over all possible split points and calculates the impurity of the resulting child nodes for each of these possible splits. The split variable and corresponding split value are chosen, such that the impurity reduction according to equation (6) is maximised. In the example of figure 2 the first split variable is chosen as 'weight' with the split value 69. Therefore the data from the root node is split into the two child nodes 'N2' and 'N3' - central plot in figure 2. 'N2' contains the observations with a weight ≥ 69 , while 'N3' consists of the observations with a weight < 69 . The distribution of the target variable in 'N2' is 25|75 and in 'N3' 100|0. Hence both resulting child nodes are purer than their parent node 'N1'. The node 'N3' only contains observations with a positive response, therefore it is completely pure and can not be split any further - all possible impurity measures [(3), (4), (5)] for this node are 0. The node 'N2' on the other hand is not completely pure yet, and can be split further. 'N2' is now the parent node and the algorithm tries all possible splits on this segmented feature space. The highest impurity reduction of 'N2' is archived with the split-variable 'height' on the value 171. 'N2' is therefore further split into 'N4' - all observations from 'N2' with a height ≥ 69 - and 'N5' - all observations from 'N2' with a height < 69 . As well 'N4', as 'N5' are completely pure and the impurity of these can not be reduced any further. The final partitioned feature space is displayed on the rightmost plot in figure 2. Based on this final partitioned feature space predictions can be done by assigning observations to one of the segments in the feature space. An observation with weight = 90 and height = 185 for example falls into the segment 'N4' and has a predicted class probability of 100% for response class 0.

In summary: The decision tree algorithm tries to split the feature space, such that the resulting child nodes maximally gain purity regarding the target variable. This is done with an exhaustive search, trying all possible split variables and corresponding split points, choosing the one that maximises the reduction of impurity.

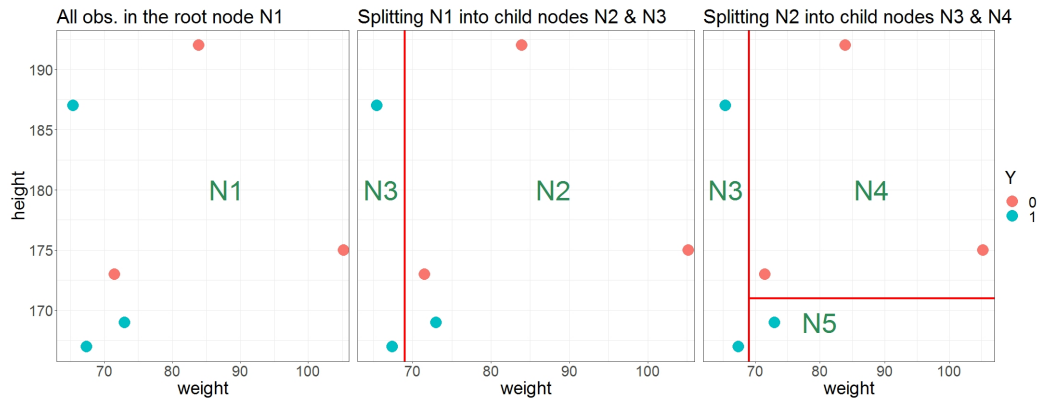


Figure 2: Example for the recursive binary splitting of decision tree on a two-dimensional feature space.

A very handy property of a decision tree consists of its natural graphical display, which makes it extremely easy to interpret - even for people without mathematical background. This visualisation is especially useful, when the training data for the decision tree holds more than two covariates and can not be displayed as scatter plot [26]. The segmentation of the feature space from figure 2 can be displayed much easier as graphical decision tree and is displayed in figure 3.

Each square in the figure represents a node of the decision tree. Each of these nodes displays the response class with the highest proportion (top), the distribution of the response classes (mid) and the fraction of observations they contain (bottom). The split variables and split values are displayed below each node - nodes without a split variable/ value are terminal nodes. The prediction for a test observation with figure 3 is very easy and intuitive. The test observation is simply passed down the decision tree until it reaches a terminal node. This is shown for a observation with weight = 90 and height = 185. The first node splits on the variable 'weight' with the value 69. As the test observations has a weight ≥ 69 it is send down to the left child node. The next node splits the on the variable height with the value 171. As the test observation is taller than 171cm it is send to the left child node. The observation is then in the node on the leftmost in the bottom of figure 3. This node is a terminal node and can not be divided further. The distribution of this node equals 100|0 and the prediction for the observations is therefore class 0 with 100% probability.

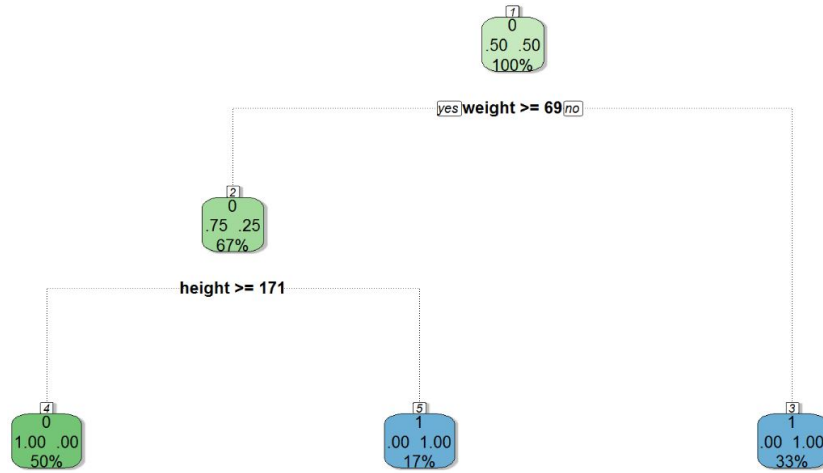


Figure 3: Corresponding decision tree for the segmented feature space on the rightmost plot in figure 2.



A decision tree repeats the binary partition of the feature space recursively until the resulting nodes are as pure as possible - even to the extent that the terminal nodes only consist of a single observation. The complexity of a decision tree grows with the number of used splits and resulting terminal nodes [27]. The more complex a decision tree, the higher the chances of overfitting. On the other hand, a tree with not enough complexity “might not capture the important structure[s]” [[26], p. 20]. So when should a tree stop with the binary partition of the feature space? There are multiple stopping criteria to control this, whereby the two most common used arguments are [26]:

- MinSplit:
“The minimum number of observations that must exist in a node in order for a split to be attempted” [[28], p. 22]
- Complexity:
“Split tree nodes only if the decrease in impurity due to splits exceeds some threshold” [[26], p. 20]

Both arguments have a huge impact on the complexity of a decision tree, as they control when the tree stops the partition of the feature space. The ‘MinSplit’ argument forces the tree to stop the partition, as soon as a potential parent node contains less than ‘MinSplit’ observations. The higher this argument is the earlier the tree has to stop growing and therefore the less complex is the resulting tree. ‘Complexity’ on the other hand only allows splits that lead to a decrease of impurity of a given threshold when splitting

the parent node N to N_1 and N_2 . The drawback of this argument is that it is rather short-sighted, as a “seemingly worthless split might lead to a very good split below” [[26], p. 20]. Hence in this thesis the ‘MinSplit’ argument to control the complexity of a decision tree is preferred over the ‘Complexity’ argument.

The advantages of the decision tree method are numerous. It is easy to interpret, has no problems with outliers in the features, captures interaction effects between features, handles categorical features and scales well with larger data [26]. Besides all these advantages, unfortunately there is also a huge disadvantage. A decision tree is highly unstable meaning that “small changes in the data could lead to completely different splits, thus, to a completely different tree” [[27], p. 26]. Even the removing of one single observation/ feature from the train data can lead to a completely different decision tree.



2.2.2 Random Forest Model

As already mentioned in the beginning of this chapter, the random forest is an ensemble method that uses the decision tree as a base learner. The random forest model therefore consists of multiple decision trees. To meaningfully train these multiple decision trees, the random forest method uses a modified version of bagging that was originally proposed by Breimann in 1996 [29]. As bagging is an important component of the random forest method it is now explained in more detail.

To meaningfully train M base learners on a single data set, each base learner needs to be fit on a modified data set, else all the resulting learners are completely identical. To generate a different data set for each of the M base learners, bagging - short for **B**ootstrap **A**ggregation - is applied to the original data. In the first step bootstrapping is applied. It is a “type of



resampling where large numbers of smaller samples of the same size are repeatedly drawn, with replacement, from a single original sample” [30]. The bootstrapping therefore generates M different bootstrap samples of the original data D and trains any base learner B on these M bootstrapped data sets. To obtain a prediction from these, each of M fitted base learners is asked for a prediction $B_m(x)$. These M different predictions are then aggregated via their average $\hat{f}(x) = \frac{1}{M} \sum_{i=1}^M B_m(x)$ [[31], p. 4]. Bagging works best for learners with a high variance - e.g. a decision tree - as it reduces the variance of the base learner and increases only the bias in return [31].

The random forest method uses a slight modification of the bagging algorithm to construct bootstrapped decorrelated decision trees [20]. To do so the random forest algorithm does not only fit each decision tree on a sepa-

rate bootstrapped data set, but decreases the correlation of these as well by randomly drawing 'mtry' features as possible split candidates at each split, instead of having all 'p' features as possible split candidates [31]. The standard value for 'mtry' with a categorical response class is $\lceil \sqrt{p} \rceil$, whereby p equals the number of variables [32]. Hence the decision trees only see a subset of the available features, when partitioning a given node. This modification of the original bagging algorithm therefore ensures that the trees are grown more differently and the resulting trees are less correlated as with the regular bagging algorithm. The modified bagging algorithm to fit a random forest is the following [[20], p. 588]:

Algorithm 1: Random forest

input : $D \leftarrow$ data of n observations & p features
 $M \leftarrow$ number of trees in the forest
 $n_{min} \leftarrow$ 'MinSplit' argument of a decision tree
 $mtry \leftarrow$ number of variables to draw at each split

for $m \leftarrow 1$ **to** M **do**

1. Draw a bootstrap sample \mathbf{Z}^* of size ' n ' from ' D ';
2. Grow a decision tree of the random forest, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached;
 - 2.1 Randomly draw 'mtry' of the ' p ' available variables;
 - 2.2 Pick the best splitting point among the 'mtry' variables;
 - 2.3 Split the node into two daughter nodes;

The procedure to receive a prediction from a the random forest model is the same as in the original bagging algorithm. The input x is passed to each of the decision trees in the random forest model and each of these trees creates a prediction based on the input x - details to that in the subsection before. The final prediction in case of a categorical response can either be the average of all M predicted class probabilities or the label that was predicted by the majority of the trees.

2.2.3 OOB error

A very handy property of the random forest method is the so called out-of-bag error (OOB error). The random forest model consists of multiple decision trees, whereby the data for each of these trees is obtained by drawing observations with replacement from the original data. For each tree, the

average probability for an observation not to be drawn is ~ 0.37 [[31], p. 12]:

$$P(\text{Obs. not drawn}) = \left(1 - \frac{1}{n}\right)^n \xrightarrow{n \rightarrow \infty} \frac{1}{e} \approx 0.37 \quad (7)$$

- n : Amount of observations in the data

Therefore each tree is grown with $\sim 63\%$ of the available observations. The remaining $\sim 37\%$ of available observations were not used to grow the decision tree and can hence be used to get an estimate of the predictive performance of the **whole** random forest model - the so called OOB error. Before explaining the OOB error, let's have a look at figure 4. The figure displays the M different decision trees of a random forest model that was originally supplied with data of n observations and p features. Under each tree the data used for growing is displayed - a pink background indicates that the observation is **in-bag** and could be used to grow the tree, while a grey background indicates that the observation is out-of-bag. It should be noticed that the observations are drawn with replacement, so that an observation can enter the in-bag samples more than once.

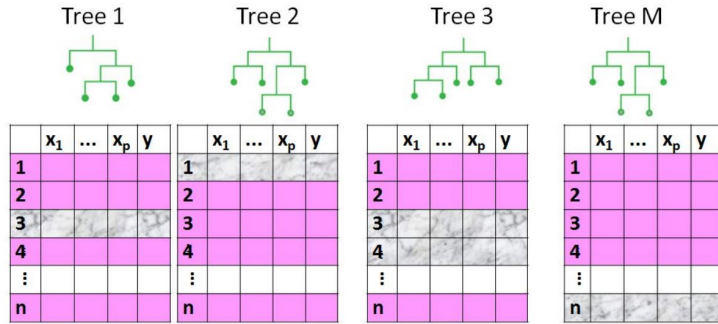


Figure 4: The data used to grow the M different decision trees of a random forest. Below each decision tree the in-bag observations are labelled in pink, while out-of-bag observations are labelled in grey [[31], p. 13].

To receive the OOB error of a random forest model, each of the M decision trees is asked for a prediction for the current observation i , but only if the observation i is an out-of-bag observation for the tree. This results in $\psi \leq M$ predictions for the observation i , whereby the final out-of-bag estimation for observation i equals the average of the ψ predictions. After receiving this out-of-bag prediction for all n observations, the final OOB error can be calculated. To compare the predicted classes and the true response

class of the n observations, any metric - e.g. Accuracy, F-1-Score - can be used. The OOB error “is almost identical to that obtained by N-fold cross validation” [[20], p. 593]. Therefore, unlike most other prediction models, the random forest can be fit and evaluated in one single step - an extremely handy property.

2.2.4 Variable importance

In most applications, not all feature variables are equally important and mostly only a few have a relevant influence. Therefore the property of variable importance in the random forest method has a high practical usage. Even though the single decision trees of a random forest model are highly interpretable, the random forest model itself “lose[s] this important feature, and must therefore be interpreted in a different way” [[20], p. 367].

Firstly lets have a look on how to calculate the variable importance of a variable x_l for the decision tree T . A decision tree T consists of J internal nodes, whereby the feature-space of each node t is split into two sub-regions. To each node t in the decision tree T , the used split variable l as well as the improvement of purity i_t can be extracted. The “relative importance of variable x_l is the sum of such squared improvements over all internal nodes for which it was chosen as the splitting variable” [[20], p. 368]:

$$I_l^2(T) = \sum_{t=1}^J \hat{i}_t^2 \mathbb{1}(v(t) = l) \quad (8)$$

- $v(t)$: Used split variable in node t



This equation and the corresponding logic is easily extendable to a whole random forest model. The variable importance for a variable x_j is then the average importance in every of the M single trees the random forest consists of [[20], p. 368]:

$$I_l^2 = \frac{1}{M} \sum_{m=1}^M I_l^2(T_m) \quad (9)$$

Another possibility to measure the variable importance in a random forest model is based on the permutations of the out-of-bag observations. All out-of-bag observations of a decision tree T_m are passed to their tree for a prediction and the accuracy of the decision tree is recorded - acc_m , without permutation. To obtain the importance of a variable x_l , all out-of-bag observations of the decision tree T_m are permuted in the variable x_l , such that all out-of-bag

observations receive a different value for the variable x_l . The permuted out-of-bag observations are then passed to the decision tree T_m and the accuracy of the decision tree is recorded again - $acc_{m, \text{ with permutation in } x_l}$. The difference between the regular OOB accuracy - $acc_{m, \text{ without permutation}}$ - and the OOB accuracy with permuted variable x_l - $acc_{m, \text{ with permutation in } x_l}$ - is used as measure for the importance of the l -th variable in the decision tree T_m . The average importance of the l -th variable over all M decision trees equals the variable importance for x_l for the whole random forest model [31].

This technique to access the importance for the different variables is displayed in figure 5 for the variable x_1 . For the decision trees 1 and M the data used to train these is displayed below. A grey background marks the out-of-bag observations. Based on these observations the out-of-bag accuracy can be calculated for each of the M trees. For each tree this results in $acc_{m, \text{ without permutation}}$. Then the values of the variable x_1 are permuted for the out-of-bag observations for each decision tree. Following the out-of-bag accuracy is calculated with the permuted variable x_1 resulting in $acc_{m, \text{ with permutation in } x_1}$. The difference $diff_m$ between $acc_{m, \text{ with permutation in } x_1}$ and $acc_{m, \text{ without permutation}}$ represents the importance of variable x_1 in the decision tree T_m . The final importance of variable x_1 then equals the average over all these differences $\frac{1}{M} \sum_{i=1}^M diff_i$ [[31], p. 16].

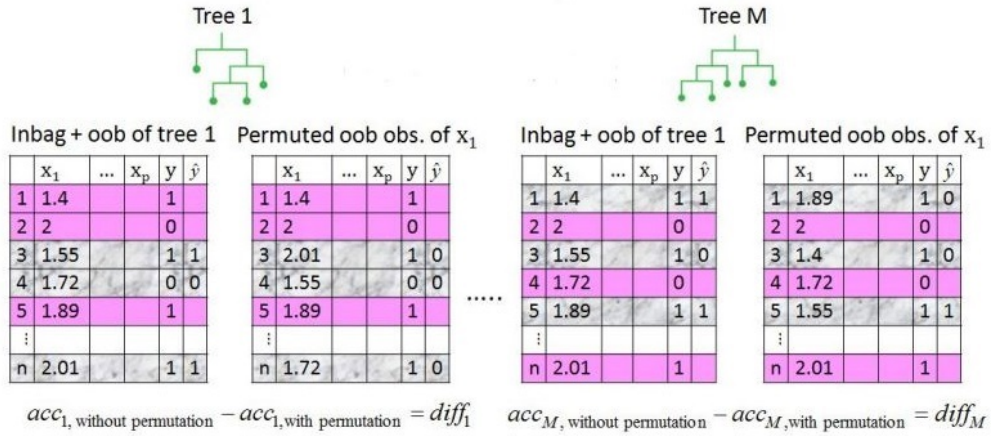


Figure 5: Calculation of the variable importance of x_1 for a random forest model consisting of M decision trees [[31], p. 16].

2.3 Complete-Case Approach

In this section the first baseline approach to handle data with block-wise missingness is explained - the so called 'Complete-Case' approach. This approach does not modify the random forest model itself, but processes the training data, such that it does not contain any missing values afterwards. This has the advantage that every prediction model - e.g. a random forest model - can be regularly trained on the processed data and the disadvantage of not using all available folds and feature-blocks. Therefore it is a rather simple approach and shall serve as a first baseline. The results from this first baseline approach are a hurdle to overcome for the more sophisticated approaches from the sections 2.5 - 2.7.

Let's have a look the approach itself. As block-wise missingness can affect the test data as well as the training data it is possible that the test observations are missing feature-blocks - even if these are available in the training set. The 'Complete-Case' approach removes all folds from the training data that miss at least one of the available feature-blocks from the test data. The feature-blocks of the training data that are not available in test set are removed as well. After the processing, the training data only consists of the feature-blocks available in the test set with observations that were completely observed in these blocks. Based on this processed training set a random forest can be regularly trained. The prediction on the test observations with such a fitted model can be then be done completely regular, as the model does not use any split variables that are not available for the test observations, as well as the test observations do not contain any features the model has not been trained with. To make the processing of the training data easier to understand two examples are shown in figure 6. In these examples the concatenated data with block-wise missingness from figure 1 is used as as a exemplary training data:

1. Example: This example is displayed in the top of figure 6. Even though the training data consists of four feature-blocks, the test observations have only two observed feature-blocks - 'Clinical' and 'CNV'. The 'Complete-Case' approach processes the training data, such that it removes all observations that miss at least one of the available feature-blocks of the test set. Therefore only observations from the fold 'Hospital 1' can be used, as the observations from the other folds either miss the feature-block 'Clinical' or 'CNV'. The fold and feature-blocks that can be used for the model fitting are marked with a green box. On this processed data a regular random forest model can be trained and used to create predictions for the test observations then. The processed training data contains two feature blocks and two folds less than the original training data, as these were removed by the processing

of the 'Complete-Case' approach.

2. Example: This example is displayed in the bottom of figure 6. The available training data originally consists of four feature blocks, while the test observations were only observed in the single feature-block 'CNV'. The 'Complete-Case' approach removes now all observations from the training data that do not have an observed 'CNV' feature-block. The observations from the fold 'Hospital 1' and 'Hospital 3' can be used as training data, whereby only feature-block 'CNV' of these observations are used. The other feature-blocks from the folds 'Hospital 1' and 'Hospital 3' need to be discarded, as they are missing in the test set. The folds and feature-block that can be used for the model fitting are marked with the green box in the figure. On this data a regular random forest model can be trained and used to create predictions for the test observations then. As in the example above, **most** of the original training data is discarded by the 'Complete-Case' approach.



		Y	Clinical	CNV		
<u>Test Set</u>		?	✓	✓		
		Y	Clinical	CNV	RNA	MIRNA
<u>Processed Data</u>	HOSPITAL 1	✓	✓	✓	×	×
	HOSPITAL 2	✓	✓	×	✓	×
	HOSPITAL 3	✓	×	✓	×	✓

		Y	CNV			
<u>Test Set</u>		?	✓			
		Y	Clinical	CNV	RNA	MIRNA
<u>Processed Data</u>	HOSPITAL 1	✓	✓	✓	×	×
	HOSPITAL 2	✓	✓	×	✓	×
	HOSPITAL 3	✓	×	✓	×	✓

Figure 6: Two examples for the 'Complete Case' processing of the training data according to the available feature-blocks in the test set.

Besides the generous discarding of training data, the method has another big disadvantage. As the 'Complete-Case' approach removes all observations from the training set that miss at least one of the available feature-blocks of the test set, it may happen that there are no training observations left after the processing. This situation is displayed in figure 7. The test set in the figure contains the feature-blocks 'RNA' and 'MIRNA' as observed feature-blocks. But as no fold in the training data was observed with both

of these feature-blocks, the 'Complete Case' approach is not applicable, as the processing of the training data results in an empty data frame.

<u>Test Set</u>		Y	RNA	CNV
		?	✓	✓

	Y	Clinical	CNV	RNA	MIRNA
HOSPITAL 1	✓	✓	✓	×	×
HOSPITAL 2	✓	✓	×	✓	×
HOSPITAL 3	✓	×	✓	×	✓

<u>Processed Data</u>	
-----------------------	--

Figure 7: The 'Complete Case' processing results in a empty training set, such that no model can be trained. In these settings, predictions can not be generated for the test set with the 'Complete Case' approach.

In summary, the 'Complete Case' approach removes all observations that miss at least one of the observed feature-blocks in the test set. Also all feature-blocks from the training data that are not available in the test set are removed. This data processing approach can discard a big part of the original training data and therefore does not handle the data very efficiently.

2.4 Single-Block Approach

The second baseline approach to handle data with block-wise missingness is introduced in this section. It is called 'Single-Block' approach and as well as the 'Complete-Case' approach, it does not modify the random forest model itself but processes the training data, such that it does not contain any missing data anymore. As well as the 'Complete-Case' approach, the 'Single-Block' approach discards much of the available data. As this approach is rather naive it is the second baseline approach and shall serve as another lower limit for the performances of the more sophisticated approaches from the coming sections 2.5 - 2.7.

As the name of the approach already suggests, it only uses a single feature-block to train a random forest model and predicts on the test set then. In order to create predictions on the test set the model must be trained with a feature-block that is also available in the test set. Else the fitted model could not predict on the test set, as it uses split variables that are not available in the test data. Hence the single feature-blocks from the training data that can be used to train a model depend on the observed feature blocks in the

test set. The concept of this approach is now explained with the example shown in figure 8. The training data in this example was already introduced in the section 2.1 and has been used as an example in the previous section as well:

Example: The test set for this example is displayed in top of figure 8 and contains two different feature blocks - 'Clinical' and 'CNV'. The training data consists of four different feature blocks and three different folds in total. The 'Single-Block' approach processes the training data in multiple ways to get rid of the block-wise missingness in the training set. For each available feature-block in the test set it is checked, whether the training data consists of the feature-block as well. For each feature-block that the test and training set have in common a random forest is fitted and used to predict the outcome of the test observations. In the current example it is firstly checked, whether the training data consists of a 'Clinical' or a 'CNV' feature-block. In this example the training data consists of both feature-blocks of the test set. For each of the feature-blocks the test and training set have in common a separate processed data set is created.

Processed Data 1: As the test and training set have the feature-block 'Clinical' in common, the first processed training data only consists of the response Y and the feature-block 'Clinical' for the observations that have been observed in this block. This subset of the data is displayed in the middle of figure 8 and marked with a green box. Based on this subset, a random forest model can be regularly fit and used to create predictions for the test set. For the predictions on the test set only the features from the 'Clinical' feature-block can be used.

Processed Data 2: As the test and training set also have the feature-block 'CNV' in common, another processed training set is created. The processed training data then only consists of the response Y and the feature-block 'CNV' for the observations that were observed in the 'CNV' block. This subset of the data is displayed in the bottom of figure 8 and marked with a green box. Based on this data, a random forest model can be regularly fit and used to create predictions for the test set. For the predictions on the test set only the features from the 'CNV' feature-block can be used.

Predictions: As the processing of the training data with the 'Single-Block' approach results in two processed training sets, this approach consists of two different fitted models then. One random forest model that was fitted on the 'Clinical' feature-block and one random forest model that was fitted on the 'CNV' feature-block. Both of the fitted models can create predictions for the test set based in the features they have been trained with. The 'Single-Block' approach can therefore result in multiple predictions for the observations in the test set.

<u>Test Set</u>		Y	Clinical	CNV		
		?	✓	✓		

	Y	Clinical	CNV	RNA	MIRNA
HOSPITAL 1	✓	✓	✓	×	×
HOSPITAL 2	✓	✓	×	✓	×
HOSPITAL 3	✓	×	✓	×	✓

<u>Processed Data 1</u>		Y	Clinical	CNV	RNA	MIRNA
		✓	✓	✓	×	×
		✓	✓	×	✓	×
		✓	×	✓	×	✓

	Y	Clinical	CNV	RNA	MIRNA
HOSPITAL 1	✓	✓	✓	×	×
HOSPITAL 2	✓	✓	×	✓	×
HOSPITAL 3	✓	×	✓	×	✓

<u>Processed Data 2</u>		Y	Clinical	CNV	RNA	MIRNA
		✓	✓	✓	×	×
		✓	✓	×	✓	×
		✓	×	✓	×	✓

Figure 8: 'Single-Block' processing of the training data so a random forest model can be regularly trained with each of these processed data sets.

In summary, the 'Single Block' approach creates an own processed training set for each of the feature-blocks the test and training set have in common. Each of the resulting processed training sets consist of only one single feature-block and do not contain any missing data, as the observations with missing values are removed. On each of these processed training sets a random forest model can be trained and used for predictions on the test set. As the 'Complete-Case' approach, the 'Single-Block' does not handle the data very efficiently. It also has the drawback that no predictions can be generated with this approach if the training and test set do not have any feature-blocks in common. In these cases the processing of the train data results in an empty data frame.

2.5 Imputation Approach

TBD

2.6 Block-wise Approach

This section introduces the 'block-wise' approach that was originally proposed by Krautenbacher in 2018 [19]. Other than the approaches from the previous sections 2.3 - 2.5 this approach does not modify the training data, but the random forest model itself. The 'block-wise' approach can

directly handle block-wise missingness in the training data and does not need to process the data at all. Therefore it uses the available training data efficiently and does not discard any observations or feature-blocks. Furthermore a 'block-wise' fitted random forest is flexibly applicable and can provide predictions for a test observation based on a single feature-block, but also for a test observation on the basis of multiple different feature-blocks.

As the name of the approach already suggests, the random forest model is fitted in a 'block-wise' manner to the training data. In the beginning all available feature-blocks of the training data are extracted. On each of these feature-blocks a random forest model is separately fitted. This enables "all observations per [feature-block] [...] to be utilised for learning" [[19], p. 102] and no observation or feature-block has to be left out. With the 'block-wise' approach as many separate random forest models are fitted as the training data has feature-blocks. To create a prediction for a test observation then, each block-wise fitted random forest model is asked for prediction. The models that were fitted on a feature-block that is not available for the test observation can not create a prediction, as these use split variables that are not available for the test observation. The remaining random forest models can create a prediction for the test observation by using the features from the test observation the model has originally been trained with. The predictions from the separate block-wise fitted models can then be aggregated - e.g. weighted average - to obtain a final prediction. The separate model fitting is explained in more detail with the example in figure 9. The training data in this example was already introduced in the section 2.1 and has been used as example in the previous sections as well:

Model Fitting: The training data is displayed in the top of figure 9 and consists of four feature-blocks and three folds. To fit a separate random forest model on each feature-block, the training data needs to be split, such that each feature-block can be used to train a random forest model. This is done by merging each feature-block and the response Y to a separate training set. In figure 9 these separate training sets are displayed as data frames with a green background below the original training data. From each of these separate training sets, all folds that contain missing values in the corresponding feature-block are removed - e.g. in the separate 'clinical' training set all observations from 'Hospital 3' had to be removed, as this hospital did not collect any clinical data. The folds that had to be removed from these separate training sets are marked with a red horizontal line, while the usable folds are marked with a tick. Based on each of these four different training sets a random forest model can be trained. This results in four random forest models in total - $RF_{Clinical}$, RF_{CNV} , RF_{RNA} and RF_{MIRNA} . Each of these models has been trained with a single feature-block only - e.g.

RF_{RNA} was only trained with the feature-block 'RNA'.

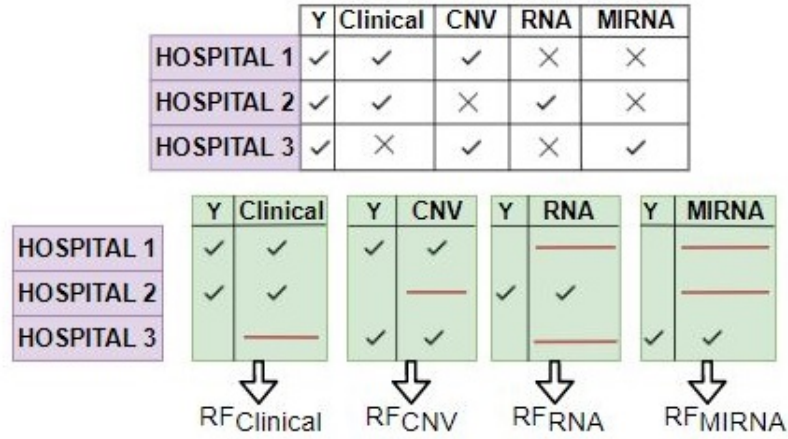


Figure 9: Example for the training of random forest models with the 'block-wise' approach.

The 'block-wise' approach trains the separate random forest models on the distinct feature-blocks in the training data and has as many separate random forest models as the training data consists of distinct feature-blocks. But how can we these models be used to create a prediction? As already mentioned, the block-wise predictions from the different random forest models need to be aggregated for a final prediction. This is explained in the following paragraph based on the example in figure 10.

Predictions: Assume that the four block-wise fitted random forest models from figure 9 can be used for the example in figure 10. The test set is displayed at the top of figure 10. For the observations in this test set the outcome 'Y' needs to be predicted. Other than the training data from figure 9, the test set only contains three feature-blocks and misses the 'CNV' feature-block from the training data. To create predictions for the observations in the test set the four block-wise fitted random forest models from figure 9 can be used. Each of the four block-wise fitted random forest models is asked to create predictions for the test observations. As the test data contains the feature-blocks 'Clinical', 'RNA' and 'MIRNA' the corresponding random forest models $RF_{Clinical}$, RF_{RNA} and RF_{MIRNA} can be used to predict on the test data. The random forest model RF_{CNV} can not create predictions on this test set, as the feature-block 'CNV' is not available. Each of the three block-wise fitted models $RF_{Clinical}$, RF_{RNA} and RF_{MIRNA} create a prediction for each observation in the test set, by only using the variables from the feature-block the models have originally been trained with. Therefore each

model creates a prediction for each observation in the test set, such that there are three predicted outcomes for each observation - $\text{Preds}_{\text{Clinical}}$, $\text{Preds}_{\text{RNA}}$ and $\text{Preds}_{\text{MIRNA}}$. These predictions represent the probabilities for each of the possible response classes. The final predictions for the target variable 'Y' equals a weighted average of these predictions.

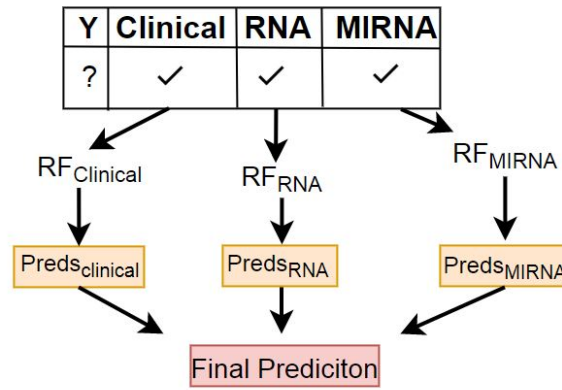


Figure 10: Example for the prediction on test data with the 'block-wise' approach. The fitting of the random forest models was described with figure 9.

To create a meaningful weighted average of the different block-wise predictions, different techniques can be applied. The simplest method is giving each block-wise fitted model the same weight and return the plain average over all block-wise predictions. But as the block-wise fitted models have been trained on different feature-blocks this might not always be optimal, as the models might differ strongly in their prediction quality. To make this clear let us assume that the feature-block 'RNA' is not related at all to the outcome 'Y', while the feature-block 'Clinical' is strongly related to it. In this case it can be assumed that the predictions based on the 'MIRNA' feature-block are worse than the predictions based on the 'Clinical' feature-block. Therefore it would be meaningful to put a higher weight on the predictions from the RF_{Clinical} model than on the predictions from the RF_{MIRNA} model. Usually the true strength of the relation between a feature-block and the target variable is unknown. Therefore the predictive quality of the different feature-blocks needs to be estimated. This can be done with the out-of-bag error of the block-wise fitted models. For each block-wise fitted random forest model the predicted classes for all out-of-bag observations are generated - see chapter 2.2.3 for details. Based on the predicted outcomes and the



true response values any metric can be calculated to judge the predictive quality of a block-wise fitted random forest model then. In this thesis either the accuracy or F-1-Score is used as metric to judge the predictive quality - details to the metrics in chapter 3.2.2. The better the out-of-bag accuracy/ F-1-Score of a block-wise fitted model, the higher the estimated predictive quality of the model. The higher the predictive quality of a model, the higher its weight and therefore the higher its contribution to the final prediction. The reason to use the F-1-Score besides the accuracy is that the F-1-Score is sensitive to class imbalances in the target variable, while the accuracy only represents the fraction of correctly classified observations.

Let us have a look at a minimalist example to make the idea of the weighted average clearer. Assume the block-wise fitted random forest models from figure 10 have the following out-of-bag accuracy and predicted probabilities for a **postive** response for the observation i :

$OOB_{Acc}(RF_{Clinical}) = 0.67$	$Preds_{Clinical}(Obs_i) = 0.19$
$OOB_{Acc}(RF_{RNA}) = 0.86$	$Preds_{RNA}(Obs_i) = 0.33$
$OOB_{Acc}(RF_{MIRNA}) = 0.21$	$Preds_{MIRNA}(Obs_i) = 0.99$

The predictions of the models represent the probability for a positive response, such that all probabilities < 0.5 result in a predicted negative class, while probabilities ≥ 0.5 result in a positive predicted response class. The true response class for the observation i is negative, as well as the predictions of $RF_{Clinical}$ and RF_{RNA} . Only the RF_{MIRNA} model predicts the response for observation i wrong. When calculating the plain average of all these predicted probabilities it results in 0.503. Therefore the final predicted probability is ≥ 0.5 and the predicted class is the positive class - which is wrong for observation i . If we use the out-of-bag accuracy of the models as weights for a weighted average **the the** final predicted probability is $0.355 < 0.5$ and therefore the predicted class is the negative class - which is correct for observation i . So instead of giving all block-wise predictions the same weight, the predictive power of the single feature-blocks can be estimated with any out-of-bag metric and used to weight the predictions. This results in a higher weight for the predictions from models with a better out-of-bag metric.

In summary: With the 'block-wise' approach a separate random forest model is fitted on each feature-block of the training data. For a prediction on a test observation, all block-wise fitted models are asked for a prediction. Only those models that have been trained with a feature-block that is available for the test observation can create a prediction. These predictions can then be averaged in a weighted/ unweighted way to create a final prediction for the test observation.

2.7 Fold-wise Adaption



This section introduces the 'fold-wise' approach that was originally proposed by Hornung et al. [14]. Other than the approaches from the sections 2.3 - 2.5 this approach does not modify the training data, but the random forest model itself. The 'fold-wise' approach can directly handle block-wise missingness in the training data and does not need to process the data at all. Therefore it does not discard any of the available observations or feature-blocks and uses the available training data efficiently. Furthermore a 'fold-wise' fitted random forest is flexibly applicable and can "provide predictions for test data that do not feature all covariates available from training" [14].



As the name of the approach already suggests, the random forest model is fitted in a 'fold-wise' manner to the training data. In the beginning all available folds of the training data are extracted. On each of these folds a random forest model is then separately fitted. This results in as many fold-wise fitted random forest models as the training data has folds. As the different folds of the training data usually consist of multiple feature-blocks, each fold-wise fitted random forest model incorporates the covariates from multiple feature-blocks. To receive a prediction on a test observation only "the subsets of covariates included in the test data that are also included in at least one of the training data sets" [14] are used. The prediction of a single fold-wise fitted random forest model is then obtained as follows: (1) Remove all trees from the fold-wise fitted random forest that use a split variable as first split that is not available for the test observation. These trees can not even split the test data once, as the first split variable is not available for the test observation. Therefore these decision trees are of no value for the given test observation. (2) For each remaining decision tree "follow each branch of the tree and cut the branch as soon as a covariate is used for splitting that is not available" [14] for the test observation. This process of cutting branches is called 'pruning'. A node that had to be pruned is a new terminal node of the decision tree then. After these two steps have been applied to the fold-wise fitted model the predictions can be obtained as for a standard random forest model. The predictions from the separate fold-wise fitted models can then be aggregated - e.g. weighted average - to obtain a final prediction. The fold-wise model fitting is explained in more detail with the example in figure 11. The training data in this example was already introduced in the section 2.1 and has been used as example in the previous sections as well:



Model Fitting: The training data is displayed at the top of figure 11 and consists of four feature-blocks and three folds. To fit a separate random forest model on each fold, the training data needs to be split, such that each fold can be used to train a random forest model. This is done by merging

the feature-blocks of a fold and the corresponding response Y to a separate training set. The feature-blocks that were not observed for a certain fold are removed from the fold-wise training data - e.g. the feature-blocks 'RNA' and 'MIRNA' were not observed for the fold 'Hospital 1' and had to be removed from the training data of the fold. In figure 11 these separate training sets are displayed as data frames with a green background below the original training data. Based on each of these three different training sets a random forest model can be trained. This results in three random forest models in total - $RF_{Hospital1}$, $RF_{Hospital2}$ and $RF_{Hospital3}$. Each of these models has only been trained with the observed feature-blocks of the different folds - e.g. $RF_{Hospital1}$ was trained with the feature-blocks 'Clinical' and 'CNV'.

	Y	Clinical	CNV	RNA	MIRNA
HOSPITAL 1	✓	✓	✓	×	×
HOSPITAL 2	✓	✓	×	✓	×
HOSPITAL 3	✓	×	✓	×	✓

	Y	Clinical	CNV
HOSPITAL 1	✓	✓	✓

→ $RF_{Hospital1}$

	Y	Clinical	RNA
HOSPITAL 2	✓	✓	✓

→ $RF_{Hospital2}$

	Y	CNV	MIRNA
HOSPITAL 3	✓	✓	✓

→ $RF_{Hospital3}$

Figure 11: Example for the training of random forest models with the 'fold-wise' approach.

The 'fold-wise' approach trains the separate random forest models on the distinct folds of the training data and consists of as many separate random forest models as the training data consists of folds. But how can we use these models to create a prediction? As already mentioned, the fold-wise predictions from the different random forest models need to be aggregated for a final prediction. To receive a prediction from a fold-wise fitted random forest model the single decision trees of such a model might be pruned. Before explaining the aggregation of the fold-wise predictions it is important to understand pruning. It is explained in the following paragraph with the help of figure 12:

Pruning: Pruning is a technique that can be applied to decision trees, if these contain split variables that are not available for a test observation for



which a prediction is asked for. This process is explained with the help of figure 12. In the top of the figure, the original decision tree from figure 3 can be seen. It was grown on the basis of the two feature variables 'weight' and 'height'. To obtain a prediction the observation is simply passed down the tree until it reaches a terminal node. The predicted probabilities equal the distribution of the target variable in the terminal node. But how can this decision tree be used to predict on a observation with a unknown 'height'? To receive such a prediction the original decision tree has to be pruned. For this all nodes that split with the variable 'height' needs to be cut. This is displayed in the bottom of figure 12. The scissors indicate the pruning at the node that uses 'height' as split variable. This node is a terminal node then. The pruned tree has one terminal node less than the original decision tree and can create predictions for observations without a 'height' variable, as the pruned decision tree does not use this variable as split variable anymore.

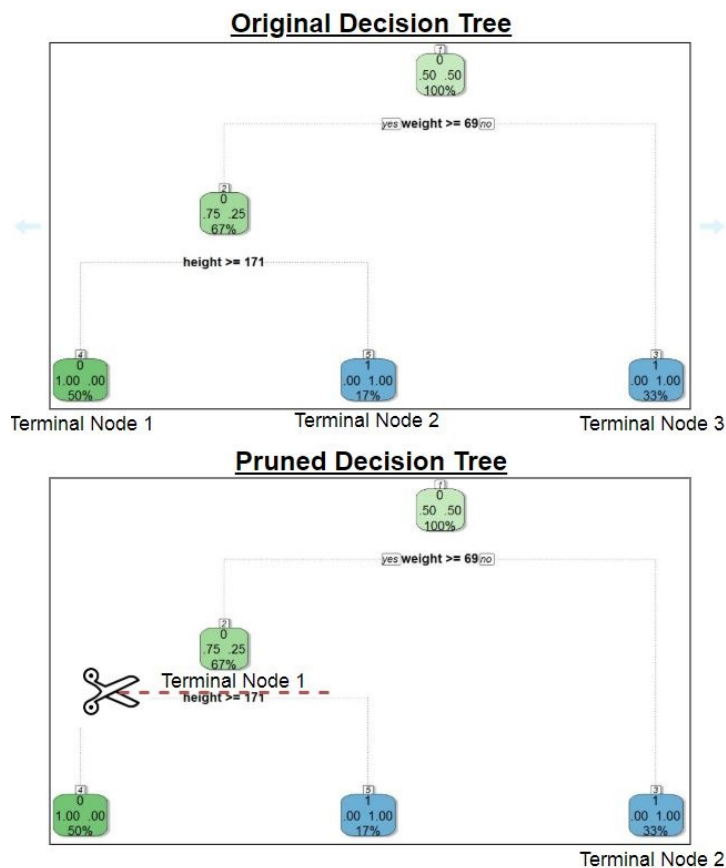


Figure 12: Example for the pruning of a single decision tree. The decision tree was originally introduced in figure 3

The process of receiving a final prediction for a observation based on the predicted classes from multiple fold-wise fitted random forest models, is explained in the next paragraph based on the example in figure 13.

Predictions: Assume that the three fold-wise fitted random forest models from figure 11 can be used for the example in figure 13. The test set is displayed at the leftmost of figure 13, whereby the outcome 'Y' needs to be predicted. Other than the training data from figure 11, the test set only contains three feature-blocks and misses the 'CNV' block. To create predictions for the observations in the test set the three fold-wise fitted random forest models from figure 11 can be used. Each of these models is asked for a prediction. As $RF_{Hospital2}$ was only trained on feature-blocks that are also available in the test set - 'Clinical' and 'RNA' - no tree needs to be pruned, as all of the used split variables in this random forest model are available for the test set. The prediction on the test set with $RF_{Hospital2}$ is therefore completely regular. The fold-wise fitted random forest models $RF_{Hospital1}$ and $RF_{Hospital3}$ were both trained with the feature-block 'CNV'. This feature-block is not available for the test observations from figure 13. Therefore the single decision trees in $RF_{Hospital1}$ and $RF_{Hospital3}$ need to be pruned, as these trees could contain nodes with split variables that are not available for the test observations. Firstly all decision trees of these models that use a 'CNV' covariate as first split variable have to be removed, as these trees can not even partition the test data once. Secondly all remaining trees are pruned as explained in the paragraph before. After applying these two steps to the models $RF_{Hospital1}$ and $RF_{Hospital3}$, the predictions for the test observations can be obtained "as in the case of a standard" [14] random forest. In the end, each fold-wise fitted model creates a prediction for each observation in the test set, such that there are three predicted outcomes for each observation - $Preds_{Hospital1}$, $Preds_{Hospital2}$ and $Preds_{Hospital3}$. These predictions represent the probabilities for each of the possible response classes. The final predictions for the target variable 'Y' equal a weighted average of these predictions then.

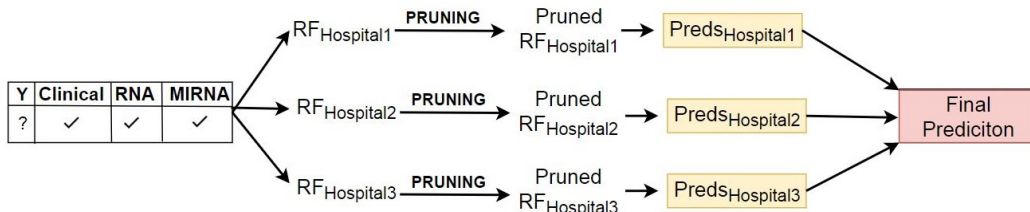


Figure 13: Example for the prediction on test data with the 'fold-wise' approach. The training of the models was described with figure 11.

To create a meaningful weighted average of the different fold-wise predictions, different techniques can be applied. The simplest method is giving each fold-wise fitted model the same weight and return the plain average over all the fold-wise predictions. This might not always be optimal, as the fold-wise fitted random forest models were trained with different combinations of feature-blocks and might differ in their predictive quality. As for the 'block-wise' approach from section 2.6 the predictive quality of the different prediction models can be estimated with a out-of-bag metric - e.g. accuracy or F-1-Score. The out-of-bag metrics of the fold-wise fitted random forest models can then be used to weight the predictions of the different models. The higher the out-of-bag metric for a model, the higher the contribution of the model to the final prediction. There is only one difference in the out-of-bag calculation between the two approaches. As the fold-wise fitted random forest models might be pruned - depending on the test data. The out-of-bag predictions are generated with the pruned trees of the random forest then. This is meaningful, as pruning might reduce the predictive power of a model and we want to obtain a realistic estimation of the predictive power of a random forest model on the test set.

In summary: The fold-wise approach fits a separate random forest model on each fold of the training data. Depending on the available feature-blocks in the test set, the fold-wise fitted random forest models might be pruned. After the pruning process each fold-wise fitted model can generate predictions for the test set. These predictions can then be averaged in a weighted/unweighted way to create a final prediction for the test observation.

3 Benchmark Experiments

3.1 Data

3.1.1 Own data

Subsetting the omics blocks

Inducing blockwise missingness

3.1.2 Data from Hagenberg

3.1.3 Real data

3.2 Accessing the Performance

3.2.1 CV - Test-Situations

3.2.2 Metrics

4 Results

4.1 Own data

4.1.1 Scenario 1

4.1.2 Scenario 2

4.1.3 Scenario 3

4.1.4 Scenario 4

4.2 Data from Hagenberg

4.3 Real data

5 Discussion and Conclusion

6 Bibliography

- [1] Francis S Collins. “Medical and societal consequences of the Human Genome Project”. In: *New England Journal of Medicine* 341.1 (1999), pp. 28–37.
- [2] *National Human Genome Research Institute*. <https://www.genome.gov/about-genomics/fact-sheets/Sequencing-Human-Genome-cost>. Accessed: 2020-01-07.
- [3] Belinda JF Rossiter and C Thomas Caskey. “Impact of the Human Genome Project on medical practice”. In: *Annals of surgical oncology* 2.1 (1995), pp. 14–25.
- [4] Sara Goodwin, John D McPherson, and W Richard McCombie. “Coming of age: ten years of next-generation sequencing technologies”. In: *Nature Reviews Genetics* 17.6 (2016), p. 333.
- [5] *Veritas - The Genome Company*. <https://www.veritasgenetics.com/myGenome>. Accessed: 2020-01-19.
- [6] Ke Bi et al. “Transcriptome-based exon capture enables highly cost-effective comparative genomic data collection at moderate evolutionary scales”. In: *BMC genomics* 13.1 (2012), p. 403.
- [7] Anne-Laure Boulesteix et al. “IPF-LASSO: Integrative-penalized regression with penalty factors for prediction based on multi-omics data”. In: *Computational and mathematical methods in medicine* 2017 (2017).
- [8] Valeria D’Argenio. “The high-throughput analyses era: are we ready for the data struggle?” In: *High-throughput* 7.1 (2018), p. 8.
- [9] Gregory B Gloor et al. “Microbiome profiling by illumina sequencing of combinatorial sequence-tagged PCR products”. In: *PloS one* 5.10 (2010).
- [10] Shrutii Sarda and Sridhar Hannenhalli. “Next-generation sequencing and epigenomics research: a hammer in search of nails”. In: *Genomics & informatics* 12.1 (2014), p. 2.
- [11] Forest M White. “The potential cost of high-throughput proteomics”. In: *Sci. Signal.* 4.160 (2011), pp. 8.
- [12] *National Institutes of Health*. <https://commonfund.nih.gov/arra/highthroughput>. Accessed: 2020-01-30.
- [13] Moritz Herrmann. “Large-scale benchmark study of prediction methods using multi-omics data”. PhD thesis. 2019.

- [14] Roman Hornung and Marvin N Wright. “Block Forests: random forests for blocks of clinical and omics covariate data”. In: *BMC bioinformatics* 20.1 (2019), p. 358.
- [15] Simon Klau et al. “Priority-Lasso: a simple hierarchical approach to the prediction of clinical outcome using multi-omics data”. In: *BMC bioinformatics* 19.1 (2018), p. 322.
- [16] Stefanie Hieke et al. “Integrating multiple molecular sources into a clinical risk prediction signature by extracting complementary information”. In: *BMC bioinformatics* 17.1 (2016), p. 327.
- [17] Qing Zhao et al. “Combining multidimensional genomic measurements for predicting cancer prognosis: observations from TCGA”. In: *Briefings in bioinformatics* 16.2 (2015), pp. 291–303.
- [18] Roman Hornung et al. “Random forests for multiple training data sets with varying covariate sets”. manuscript - unpublished yet. - in prep.
- [19] Norbert Krautenbacher. “Learning on complex, biased, and big data: disease risk prediction in epidemiological studies and genomic medicine on the example of childhood asthma”. PhD thesis. Technische Universität München, 2018.
- [20] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- [21] Jonas Hagenberg. “Penalized regression approaches for prognostic modelling using multi-omics data with block-wise missing values”. manuscript - unpublished yet. - in prep.
- [22] Hemant Ishwaran et al. “Random survival forests”. In: *The annals of applied statistics* 2.3 (2008), pp. 841–860.
- [23] Leo Breiman. “Random forests”. In: *Machine learning* 45.1 (2001), pp. 5–32.
- [24] *What is the difference between Bagging and Boosting?* <https://quantdare.com/what-is-the-difference-between-bagging-and-boosting/>. Accessed: 2020-03-06.
- [25] Leo Breiman et al. *Classification and regression trees*. CRC press, 1984.
- [26] Wenbin Lu. *Lecture 21: Classification and Regression Trees*. Department of Statistics North Carolina State University, 2019.
- [27] Bernd Bischl and Christoph Molnar. *Introduction to Machine Learning - Chapter 13: Trees*. Department of Statistics - LMU Munich, WinterTerm 2017/18.

- [28] Terry Therneau and Beth Atkinson. *rpart: Recursive Partitioning and Regression Trees*. R package version 4.1-15. 2019. URL: <https://CRAN.R-project.org/package=rpart>.
- [29] Leo Breiman. “Bagging predictors”. In: *Machine learning* 24.2 (1996), pp. 123–140.
- [30] *Bootstrap Sample: Definition, Example*. <https://www.statisticshowto.com/bootstrap-sample/>. Accessed: 2020-03-04.
- [31] Bernd Bischl and Christoph Molnar. *Introduction to Machine Learning - Chapter 15: Bagging and Random Forests*. Department of Statistics - LMU Munich, WinterTerm 2017/18.
- [32] H. Ishwaran and U.B. Kogalur. *Fast Unified Random Forests for Survival, Regression, and Classification (RF-SRC)*. R package version 2.9.3. manual, 2020. URL: <https://cran.r-project.org/package=randomForestSRC>.

7 Attachment

Figures

Plots of the impurity functions (3), (4) and (5) for a binary target variable

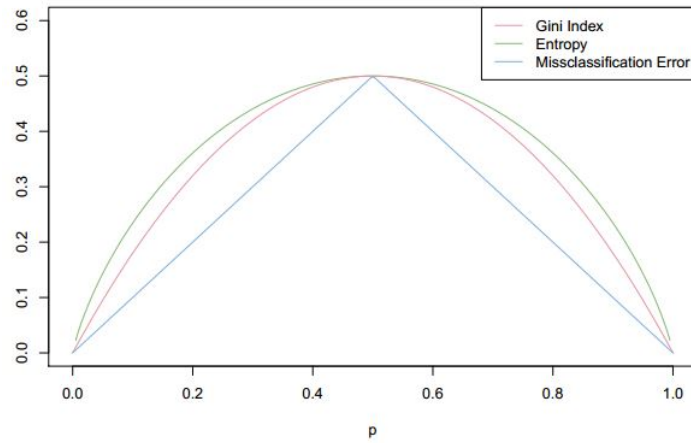


Figure 14: The different impurity functions (3), (4) and (5) plotted for a given fraction of a binary target variable within any node N [[27], p. 13]

Tables