

Chapter 1

Automatic feature selection for high-dimensional biomedical data

1.1 Introduction

For a biomedical engineering framework dimensionality reduction is important. Some biomedical datasets contain a high number of features, whereas only a selection of those features are interesting. A direct link of features to the output is desired, therefore especially dimensionality reduction based on feature selection is important. Multiple projects attempted to reduce the number of features [19, 20] which resulted in multiple algorithm proposals [84, 89, 85, 90] and tests on their performance [91, 92]. The research goal is *to evaluate the performance of feature selection methods and make a choice on which methods should be added to the framework*. In this document several of those algorithms are presented and test their quality.

Dimensionality TO BE ADDED

The dimensionality of a dataset in literature is the number of features of that dataset [77, 84, 85]. A high dimensionality therefore indicates a high number of features. The package *metalearn* calls dimensionality the number of features divided by the number of instances. If this definition of dimensionality is used it is specifically stated. This definition is used to show possible imbalance between the number of features and instances and normally the dimensionality should be very low, at least < 1 .

Dataset issues TO BE ADDED The choice of a data analysis approach heavily depends on the nature of the data. The amount of data in the biomedical world is growing at an enormous rate, faster than biomedical engineers can analyse. Several additional challenges came up with this uncontrollable growth. These challenges are in many cases focused on data volume, dimensionality, complexity, heterogeneity and quality [35, 36].

Scientists tend to collect abundant data, which makes data sets bigger than needed. Both in number of instances and features, data sets are harder to understand or analyse when more instances and features are available [35]. This volume problem usually is tackled by taking sub-populations of the complete set. Sub-sets can either be focused around a part of the population (e.g. gender, age, race) or taken at random to still represent all of it. Due to the efficiency of analysis techniques and the rise in computational speed of servers [37], volume on its own becomes less of an issue. Volume does however become a challenge when it is combined with heterogeneity and quality [38, 39].

Not all data sets have a high number of instances that causes a big data volume. Sometimes there are relatively few instances, while the number of features is disproportionately high [40]. This is also called high dimensionality (high number of features and low number of samples). Usually many of those features are not relevant enough for research, however are still used for testing. Trying to remove features that are not important significantly helps finding relations between the other features and create more knowledge about the research topic. Lowering the number of

Table 1.1: An overview of the symbols used to explain the data

Value	Description
m	The number of features
F	A list of (m) features
n	The number of samples
S	A list of (n) samples
X	A matrix of $(n \times m)$ dataset values
y	A vector of (n) output labels
W	A list of feature weights

features also makes the data volume decrease, therefore making analysis easier. Essentially an optimal features set should be selected to obtain the best results [41].

1.1.1 Data Symbols

TO BE DISCUSSED - ONLY TABLE?

Often in this research project, equations and pseudo-algorithms are discussed. In these equations and pseudo-algorithms, several symbols were used to explain the data. A quick overview of these symbols is given (Table 1.1) as well as a more detailed description:

- F - A list that contains all the names of the features. It is a list of size m with m being the number of features. A subset x of F is written as F_x .
- S - A list of all samples in a dataset. It is a list of size n with n being the number of samples. A subset x of S is named S_x and a single sample is usually called s .
- X - A matrix that contains all sample values for every feature. It is an n by m matrix with n being the number of samples and m being the number of features. If the values of a specific feature f or a subset of features F_x are used this is written as X_f and X_{F_x} respectively, hence the column f or columns F_x are selected. Similarly for X_s and X_{S_x} the row s or rows S_x respectively are collected from X .
- y - A vector that contains all class labels for every sample. It is a vector of size n with n being the number of samples.
- W - Weights given to a feature after training a feature selection technique. The weight of feature f or set of features F_x is called W_f and W_{F_x} respectively.

1.2 Software

TO BE DISCUSSED - ONLY MENTIONS OF THEM?

For this research one programming language is used, *Python* version 3.6. This language was chosen, due to the numerous possibilities in preprocessing and analysis in the open source language [70]. *Python* makes use of packages containing additional classes and methods to be used in experiments. To easily import these packages the *Anaconda* version 3 environment is used [71]. Packages available by *Anaconda* and often used in the analyses are:

- *NumPy* [72]
NumPy provides data structures efficient in use for data analysis. On top of that it also provides multiple methods to alter and preprocess these data structures, improving data analysis efficiency significantly.

- *SciPy* [73]
SciPy provides analysis methods, mainly focused on statistics. These methods are compatible with the *NumPy* data structures and provide a way to simply perform statistical analyses.
- *Scikit-learn* [74]
Scikit-learn, also known and shortened as *sklearn*, is focused around machine learning. It provides numerous simple machine learning algorithms as well as several other algorithms and classes that are often used in analysis, such as preprocessing methods and environment testing creation.
- *Pandas* [75]
Pandas is also used in parts of this research. *Pandas* creates a very intuitive way of distributing data, presenting datasets that are easy to understand for users and in analysis. It is an overlay of the *NumPy* data structure and often used in *Python*.

Two other packages were used in this research that are not part of *Anaconda*. These packages have to be downloaded separately:

- *TPOT* [76]
TPOT is an abbreviation for "Tree-based Pipeline Optimization Tool" and uses automated machine learning to find the best machine learning pipeline. It makes use of several other packages, such as the *NumPy*, *SciPy* and *Scikit-learn*, all of these being available in *Anaconda*.

1.3 Methods

Before testing several feature selection methods, first the background of the study is explained. Firstly, the datasets and their characteristics are explained. Secondly, feature selection methods are discussed. Thirdly the concept automated machine learning is discussed as a possible technique to add to the framework, combined with a tool that implements automated machine learning.

1.3.1 Datasets

Four datasets were used as a case study for the feature selection algorithms. Two sets are microarray datasets that are used for research on psoriasis [25, 26, 27, 29] and cancer [55]. Two other sets are mass spectrometry data sets, used for research on cancer [59] and micro organisms [60] (Subsection ??). These four datasets all have a high number of features varying from 1300 to 54675 features with a number of samples varying from 200 to 580 samples (Table 1.2). All of these datasets are based on classification as tests are done for different test subject groups, therefore the focus on feature selection algorithms will also be based on classification. Many features are expected to be irrelevant in this classification and therefore could be removed with feature selection.

Micro-array Datasets

TO BE ADDED, ONLY GENERAL MENTIONS?

Micro-array data is a fairly new type of biomedical data. Datasets of this kind have information on expression levels of a high number of genes [2]. These expression levels can be used for gene level research, such as genome mapping, transcription factor activity and pathogen identifications. The data is gathered by fixing DNA molecules on a glass slide on specific spots. The nucleotide base pairs of those fixed DNA molecules can form hydrogen bonds with RNA corresponding to genes, testing the expression of those genes by how many bonds are made [53] (Figure 1.1).

Micro-array data is known to present challenges in data quality and needs normalisation for the data to be useful [53]. The thousands of features that are based in the data also indicate that

Table 1.2: A schematic overview of the four datasets.

Dataset focus	Data type	Features	Samples	Classes	Remarks
Psoriasis	Micro-array	54675	580	3	- Derived from five different datasets [25, 26, 27, 29] –
Cancer	Micro-array	54675	383	9	- Used in a data mining challenge [55] –
Cancer	Mass Spectrometry	10000	200	2	- Created for the NIPS conference [59] - Several probe features are present –
Micro-organisms	Mass Spectrometry	1300	571	20	- Originates from a micro organisms study [60]

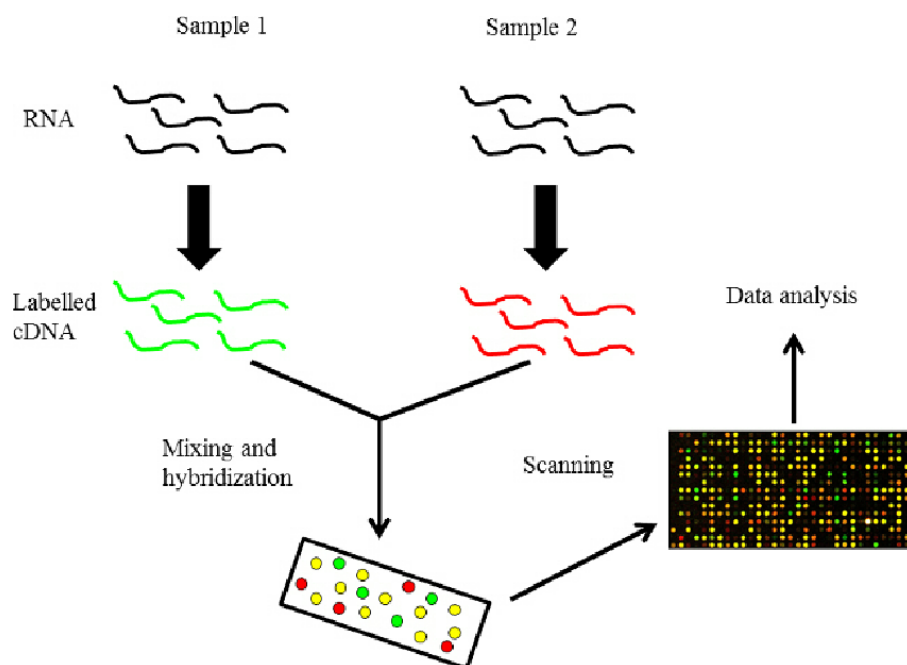


Figure 1.1: A picture showing the creation of microarray datasets [54].

feature selection is very important, so the irrelevant genes can be removed. Aside from that, size may also be an issue. Due to the sheer size of the data not all analyses will perform optimally in both time and quality. In this research two microarray datasets are used for research:

- *Psoriasis microarray dataset*

This dataset is comprised of five different data sets [25, 26, 27, 29]. These five different datasets consist of 54675 features, all corresponding to gene expression. Samples were collected from three different test subject groups: affected skin from test subjects suffering from psoriasis (214 samples), unaffected skin from test subjects suffering from psoriasis (209 samples) and skin from healthy test subjects (85 samples). Combining these three sample types gives 508 samples. Since the data comes from five different experiments, the data is normalized for every experiment.

- *Arcene: Cancer microarray dataset*

This dataset, called Arcene dataset, is used in a challenge focussing on classification problems with a low number of samples, but a high number of features [55]. It consists of the same number of features as the Psoriasis data set, 54675 features corresponding to gene expression. It also has 383 samples corresponding to nine different test subject groups. The challenge did not provide labels for the test subject groups. Also these groups differ in size, one group corresponding to 150 samples and the others varying from 16 to 47 samples.

Mass Spectrometry Datasets

TO BE ADDED, ONLY GENERAL MENTIONS?

Mass spectrometry data contains information on proteins and peptides [3, 4]. Analysis of this information is used in studies about proteins, also known as proteomics [56]. Mass spectrometry is a technique that can be used to find how much of a certain protein is present. This is done by measuring the mass of gas-phase ions originating from molecules of that protein, showing that protein is present. On top of that it can also measure the quantity of the mass, therefore showing how much of a protein is present [57].

Challenges in the analysis of mass spectrometry datasets are mainly found in the way this data is produced. The expression is given for several proteins and this expression can differ significantly between proteins. Therefore some type of normalisation is useful. Aside from that, in this technique usually a high number of proteins is tested at the same time. Therefore feature selection can be helpful to remove irrelevant genes. Two mass spectrometry datasets were used in this research:

- *RSCTC: Cancer mass spectrometry dataset*

This dataset, called RSCTC dataset, was created as a classification problem to distinguish cancer patterns from normal patterns [59]. It is created for the 'Neural Information Processing Systems' conference by merging three mass spectrometry datasets. It consists of 10000 features corresponding to either spectra of the mass spectrometry or probe variables without any predictive power. Samples from two groups are taken, from patients with ovarian or prostate cancer and from control patients. No labels are given to the groups, however it is known that one of the groups has 88 samples and the other 112 samples, combined in a total of 200 samples.

- *Micro organisms mass spectrometry dataset*

This dataset is created to back up a proposed method for routinely performing direct mass spectrometry based bacterial species identification [60]. It consists of 1300 features corresponding to different spectra of the mass spectrometry data and 20 test subject groups corresponding to Gram positive and negative bacterial species. Gram classification is a result of a Gram stain test [61]. The groups differ in size varying from 11 to 60 samples, making a total of 571 samples.

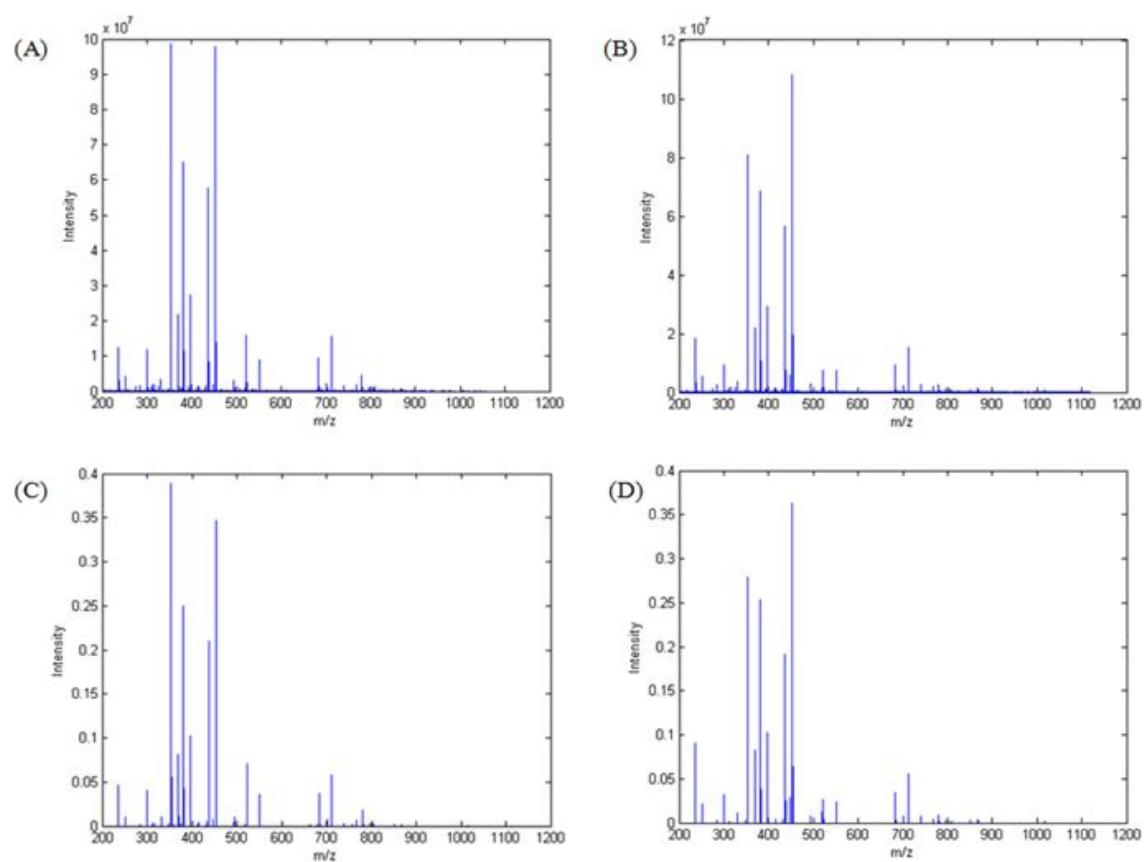


Figure 1.2: Four example figures of the data created by mass spectrometry. The measured molecule masses correspond to the peaks in the figures [58].

1.3.2 Feature Selection Methods

Feature selection is a way to perform dimensionality reduction. In feature selection a subset of features is chosen to represent the complete sample space [93]. Several techniques are available to choose a representation subset and the effectiveness of these techniques has been tested multiple times [94, 95]. These techniques can be grouped accordingly in three different categories: *filter methods*, *wrapper methods* and *embedded methods* [96]. Each of these methods is explained in what follows. For better understanding of the different feature selection methods, pseudo-algorithms are created for each of them.

Filter Methods

Filter methods are based on giving relative ranks to the features in a feature space. All features are given a value based on their performance and are ranked by those values [97, 96]. Several methods are used to rank the features. These methods are usually based on splitting data matrix X twice. First they are split by feature ($X_{f_1}, X_{f_2} \dots X_{f_n}$) for computation. Secondly the columns are split by output classes $c_1, c_2 \dots c_n$ ($X_f^{c_1}, X_f^{c_2} \dots X_f^{c_p}$) to distinguish the predictive power for features between classes. A ranking method would quantify the significance according to these values. This results in ranking method $R(X_f, y)$ in which sample label y is used to split the data per class. For illustration a ranking method collection is given:

- *T-test and ANOVA*

Statistics can be used to compare groups with each other. In statistics these groups are seen as separate distributions, which may or may not be seen as independent distributions. If there are two groups a t-test can be done to find the chance for these groups to originate from the same distribution [98]. For sample classifications the groups would consist of samples with the same label ($X_f^{c_1}$ and $X_f^{c_2}$). For example, the t-test can be done between group with class c_1 and group with class c_2 . If the values of a certain feature from these groups are highly unlikely to follow the same distribution ($R(X_f^{c_1}, X_f^{c_2})$), the p-value will be very low and therefore the rank will be higher.

If there are more than two groups that must be checked whether they are from the same distribution ($R(X_f^{c_1}, X_f^{c_2} \dots X_f^{c_p})$), a t-test cannot be used. In this case the option is available for a multiple group testing, also called analysis of variance (ANOVA). ANOVA computes whether not only two, but multiple groups can originate from the same distribution. ANOVA needs equal group sizes, otherwise the results of the tests are less reliable. The less powerful Kruskal Wallis test can be used if this reliability is needed. Another disadvantage is that t-test and ANOVA should only be used if the groups follow a normal distribution, which might not be the case [98].

- *Mutual information*

Mutual information is another way of matching features with the results. The relevance of using one variable to predict the other variable is used in the equation to compute mutual information.

Mutual information TO BE COMBINED

Mutual information can be used on two different variables. The relevance of using one variable to predict the other variable is used in the equation to compute mutual information (Equation 1.1). In this equation the probability density function p is used to find the mutual information (MI) between variables x and y [86]. Mutual information can be used for both classification and regression. It was initially meant for use in communication channels, however its statistical decision making capabilities made it a good filter method [87].

$$MI(x, y) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (1.1)$$

These ranking methods R to rank the features are only the first step in the filter methods. The next step is choose which features to filter out based on R . A quick approach would be to choose a number or fraction n and select the top n ranked features from the complete space (Algorithm 1). Another approach would be to use thresholds derived from literature [99]. For the t-test, a p-value of 0.05 is often chosen as the threshold value [100, 101].

Algorithm 1 A basic top n filter algorithm [97]

```

1: procedure FILTERSELECTION( $X, y, F, R, n$ )
2:    $F_{selected} \leftarrow \emptyset$  ▷ Start with empty feature set
3:    $Z \leftarrow \emptyset$  ▷ Start with empty set of ranking values
4:   for  $f$  in  $F$  do ▷ For all features in  $F$ 
5:      $Z_f \leftarrow R(X_f, y)$  ▷ Compute the ranking value between feature and output
6:      $Z \leftarrow Z \cup \{Z_f\}$  ▷ Add the ranking value to the set of ranking values
7:   end for
8:   Sort  $F$  by  $Z$  ▷ Sort the features by their ranking value
9:    $F_{selected} \leftarrow F_{[1, n]}$  ▷ Select the top  $n$  features from the sorted feature set
10:  return  $F_{selected}$ 
11: end procedure

```

Using filter methods has its advantages and disadvantages. These filter methods are very computationally efficient. Every feature is given a value for its rank and then a subset is selected based on those ranks. These filter methods however do not take into account dependencies between features. These dependencies could make the final feature subset worse, as maybe some features are related. Other methods are better at handling those dependencies [97, 96].

Wrapper Methods

Filter methods take into account the direct relation between features and the output classes, whereas wrapper methods focus more on the subsets of features and their ability to classify the data. Wrapper methods try to find the best combination of features to classify the given data and take into account the change when adding or removing features from a candidate subset. The choice of subset selection is "wrapped" around the quality computation, therefore this method is called a wrapper method [102]. Since an exhaustive search of trying out all possible subsets would span a computation time of 2^n with n being the number of features [103]. This combinatorial explosion should be avoided and therefore less computationally intensive approximation concepts have been constructed. Three of those concepts are explained, focusing on basic sequential search, extensions of sequential search and stochastic search.

Before explaining the possible wrapper methods, first the evaluation function J should be discussed. To find out which subset of features can classify the data the best way, a function should be used to evaluate the performance of the subset [102]. Evaluation functions can be based on conditional independence [102, 104], showing the difference in performance for a subset with and without a certain feature. Other evaluation functions are based on machine learning techniques [105, 96], rating the ability to classify the outcome. Several interesting approaches for evaluation functions are shown in maximum relevance, minimum redundancy algorithms [106, 107, 108] (MRMR algorithms). In these MRMR algorithms, several different evaluation functions are used based on the ability of features to classify the outcome (relevance) and the presence of correlation between features (redundancy). In this project, the evaluation function is fixed on the machine learning algorithm Naive Bayes, due to its understandability.

- *Sequential search*

The first concept to be explained is sequential search. This wrapper method tries to improve a candidate subset by evaluating the change of the sequential addition or removal of a specific feature. Therefore, two types of sequential search are possible, known as forward selection

and backward selection. Forward selection starts with the empty subset. Every feature is iteratively evaluated and added to the feature subset if the evaluation function shows an increase in prediction. A maximum feature subset selection size l can be defined if needed, as well. The layout of forward selection is shown as a pseudo algorithm (Algorithm 2) for understanding [102].

Backward selection does the opposite of forward selection. It starts with the complete feature set in which every feature is iteratively evaluated by the evaluation function for its contribution. If the evaluation function shows that its contribution is very small, it is removed. This way only features with a high impact will remain in the subset. The maximum number features that can be removed r can be defined if needed, as well. A layout of backward selection is shown as a pseudo algorithm (Algorithm 3) for understanding, as well [102].

In both forward and backward selection the sequence order is important. Different feature subsets will be made for different order of features. The ordering can be changed in a specific way if needed. Examples would be using the ranking concepts used by filter methods. Another possibility would be to randomize the order and run the algorithm multiple times to find the best subset [102].

Algorithm 2 A forward selection sequential search algorithm [102]

```

1: procedure SEQUENTIALFORWARDSELECTION( $X, y, F, J, \alpha, l$ )
2:    $F_{selected} \leftarrow \emptyset$  ▷ Start with empty feature set selection
3:   for  $f$  in  $F$  do ▷ For all features in  $F$ 
4:     if  $J(X_{F_{selected}}, y, X_f) > \alpha$  then ▷ Check if evaluation is higher than  $\alpha$  with feature  $f$ 
5:        $F_{selected} \leftarrow F_{selected} \cup \{f\}$  ▷ Add  $f$  to the feature set selection
6:     end if
7:     if  $\text{length}(F_{selected}) = l$  then ▷ Stop if size of feature set selection has been reached
8:       break
9:     end if
10:  end for
11:  return  $F_{selected}$ 
12: end procedure

```

Algorithm 3 A backward selection sequential search algorithm [102]

```

1: procedure SEQUENTIALBACKWARDSELECTION( $X, y, F, J, \alpha, r$ )
2:    $F_{selected} \leftarrow F$  ▷ Start with complete feature set selection
3:   for  $f$  in  $F_{selected}$  do ▷ For all features in  $F_{selected}$ 
4:     if  $J(X_{F_{selected} \setminus \{f\}}, y, f) < \alpha$  then ▷ Check if evaluation is higher than  $\alpha$  without  $f$ 
5:        $F_{selected} \leftarrow F_{selected} \setminus \{f\}$  ▷ Remove  $f$  from the feature set selection
6:     end if
7:     if  $\text{length}(F \setminus F_{selected}) = r$  then ▷ Stop if feature removal limit has been reached
8:       break
9:     end if
10:  end for
11:  return  $F_{selected}$ 
12: end procedure

```

- *Sequential search extension*

The two basic sequential search algorithms forward and backward selection can be altered for better use. An intuitive idea would be to combine the two algorithms, creating an algorithm that first selects features with the evaluation function followed by removing them according

to the evaluation function or the other way around. Also there is no restriction on only doing one iteration of both forward and backward selection, giving rise to "plus l-take away r" selection (PTA). $\text{PTA}(l, r)$ adds l features with forward selection and removes r features with backward selection per iteration, eventually converging to an optimum. These found optima in both sequential search algorithms and possible extensions can converge to local optima. Beam search is an example that also collects suboptimal branches for possible better optima, instead of only keeping the best possible outcome at all times [102].

A last example of a sequential search extension is called floating search. Instead of adding and removing a set number of features as in PTA, floating search continues to add and remove features until the best subset is found, doing a forward and backward search iteratively. Since feature relevance and redundancy changes for every new subset, all features are continuously evaluated to find out if they must be added to or removed from the subset. This way an optimal subset can be found [102].

The major advantage of using wrapper methods is that they also take into account possible dependencies between features. The computation time of wrapper methods usually is higher than of filter methods, but still relatively short as it should always converge to an optimum. These optima can be local however, so the result may not be the optimal, due to its greedy character. Also, these methods are dependent on the evaluation function and are known to be prone for over fitting [102, 96].

Embedded Methods

In both filter and wrapper methods, machine learning plays little to no role in selecting features. Filter methods do not use learning at all and a wrapper method can only use machine learning for evaluating feature subsets. Several machine learning algorithms contain attributes that can be used directly to select or eliminate features from the ideal feature set selection, for example the coefficients in linear SVM. Embedded methods combine feature selection with a machine learning algorithm M , in contrast to wrapper methods in which these are separated [111]. Usually this combination involves using the weights given to features by machine learning algorithms [112].

One approach to solve feature selection is based on contribution relaxation minimization. Since the theoretical approach of embedded methods can be very computationally intensive an approximation of the minimization problem is given, showing an example embedded method [111].

- *Contribution relaxation minimization*

Contribution relaxation is based on giving every feature f a contribution factor $\sigma_f \in [0, 1]$. This contribution factor shows the contribution of a feature to the preferred outcome. The final goal is to select a subset of features, though, and not use all features with a contribution factor. To achieve that a minimization function is used in which all contribution factors become $\sigma_f \in \{0, 1\}$. With contribution factors in $\{0, 1\}$, a subset with all features having $\sigma_f = 1$ can be selected as the ideal subset [111].

The minimization can be implemented with an embedded method. This method would however be computationally intensive and therefore approximations are used. The most used approximation is similar to the filter methods, ranking each feature and choosing features based on rank (Algorithm 4). This embedded method is also called a forward selection method. A machine learning algorithm M assigns a weight to each feature, which can be used as a rank (Algorithm ??). The difference between a ranking method R of the filter methods and the weights of a machine learning algorithm M is that M also takes into account the dependency between features [111].

Algorithm 4 An embedded forward selection algorithm [111]

```

1: procedure EMBEDDEDFORWARDSELECTION( $X, y, F, M, \alpha$ )
2:    $F_{selected} \leftarrow \emptyset$  ▷ Start with empty feature set
3:    $W \leftarrow M(X, y)$  ▷ Extract the weights for every feature
4:   for  $f$  in  $F$  do ▷ For all features in  $F$ 
5:     if  $W_f > \alpha$  then ▷ Check if weight is higher than threshold  $\alpha$ 
6:        $F_{selected} \leftarrow F_{selected} \cup \{f\}$  ▷ Add  $f$  to selected features
7:     end if
8:   end for
9:   return  $F_{selected}$ 
10: end procedure

```

The weights given by machine learning are different for every algorithm. A first and most obvious example of an algorithm giving weights are algorithms based on linear support vector machines (linear SVM, from now on called SVM). The coefficients given by SVM give the features a contribution factor in the outcome [113, 114, 115, 116]. A second example would be using decision trees and random forests. Decision trees use features to reduce the entropy between a set by splitting it in two subsets with a threshold and a feature. These splitting features can be used as a subset to create an effective feature selection outcome [117, 118, 97]. To overcome an over fitting problem, commonly occurring in decision trees, random forests can be used instead and the best splitting features can be used [119].

Tree-based Pipeline Optimization Tool

TO DO: Add automated machine learning and meta-learning paragraphs and simplification

Tree-based Pipeline Optimization Tool (TPOT) is a tool that implements autoML. It uses machine learning pipelines and genetic programming to find the best solution for every data set. TPOT makes use of this genetic programming with using the machine learning pipelines in a tree. *TPOT* consists of preprocessing and machine learning algorithms, that form the backbone of the pipelines. Their hyper-parameters are the variables. *TPOT* makes mostly use of the machine learning and preprocessing algorithms of *scikit-learn* from *Python* in which it is also written, but also several additional algorithms are present (e.g. a hot encoding algorithm).

Considering the capabilities from *TPOT* to cope with challenges in biomedical data, several methods are available. It has several different normalisation/standardisation scalers (StandardScaler, RobustScaler, MinMaxScaler) to tackle feature heterogeneity between different data sets and errors. It also has some feature selection operators to tackle errors (VarianceThreshold, SelectKBest, SelectPercentile). Wrapper methods however are not included, which may perform better in some situations. At last, if missing values are present, it imputes the median for those values before starting the evolutionary algorithm. However, it does not have any possibilities to handle non-numeric data. If non-numeric data is present, the user first needs to preprocess himself.

A big variety of algorithms are present for TPOT, as well as numerous hyper parameter values for those algorithms. This is both an advantage as an disadvantage as probably most algorithms and hyper parameter combinations are not useful as an outcome. Since TPOT randomly chooses mutations, possibly better algorithms could be chosen with less random choices. For example meta-features could be used to improve the algorithm selection.

1.4 Methods

For empirical evaluation the collection of feature selection methods experiments are done. The quality of the feature selection is tested by defining a definition of quality first. After that an exploration of feature selection is done with filter methods, followed by comparisons of multiple

feature selection methods. The four example datasets (subsection ??) were used as example datasets for the experiments.

1.4.1 Feature Selection Quality

To find out the quality of the feature selection, multiple machine learning algorithms are used [129]. Classification of the datasets can be done with several machine learning algorithms and validation and tests scores show how well the data can be classified after feature selection. The quality will be described with the accuracy of machine learning algorithms: the number of right classifications divided by the total number of classifications. Five different machine learning classifiers are used from *scikit-learn*: logistic regression, decision trees, nearest neighbour, support vector machines and Naive Bayes. Better feature selection algorithms have relatively higher accuracy, as they are better at preserving the right features. Since over fitting sometimes happens when fitting data in a machine learning algorithm, both a validation and a test score is computed for the accuracy. For every experiment a training set and a test set is created, making the test set 20% of the complete dataset. A validation score is computed by using the "leave one out" technique on the training set and a test score is computed by testing the classification score of the test set. Since the samples are not evenly distributed over the classes the precision, recall and F1 scores are also computed to find potential bias in the result.

In some cases the quality of the machine learning algorithm lacks a factor of feature subset size. The standard score in those cases is not sufficient enough to evaluate the result. To incorporate the quality of feature selection a new modified quality term *FS_score* is created (Equation 1.4.1). In *FS_score* a modified version of the original score, in this thesis being accuracy, is given by first multiplying it by a factor dependent on the number of features. This factor consists of a constant β , a value in range $[0, 1]$ which can be chosen for the influence of the value that represent the number of features. In practice $\beta > 0.95$ so the reduction does not have too much influence on the score.

$$FS_score = score \times \beta^{\#features}$$

The factor $\#features$ is an absolute count, instead of a relative count (for example a percentage of features). The reason for using an absolute count can be found in the goal of the data analysis. In data analysis the number of features that can be relevant are very limited. Useful results need relations between the input and the output, and these relations need to be as simple as possible. If the number of features is relative to the total number of features, the input size can change significantly. Take for example the Micro-organisms and the RSCTC datasets provided for this stud. 2% of the total features would be 26 and 1100 features for these datasets respectively. Whereas relations between 26 features and the output is possible to understand, relations between 1100 features and the output is much harder. Taking an absolute number of features allows for a more active choice on the input size.

To put the impact of the *FS_score* in perspective, an example figure is made how the outcome is computed (Figure 1.3). This figure shows the impact on the performance of a method when using a certain number of features. After using a correction factor, an optimum is created for which the number of features is important, the old optimal score did not include the number of features. The example (with $\beta = 0.99$) seems to give a good indication of the desired outcome. For every 10 features, the *FS_score* is reduced by about 10%, which meant for the example filter method that for 50 features an optimum was reached (Figure 1.3). Because of this empirically found desired trade-off, $\beta = 0.99$ will be chosen in this project when using *FS_score*.

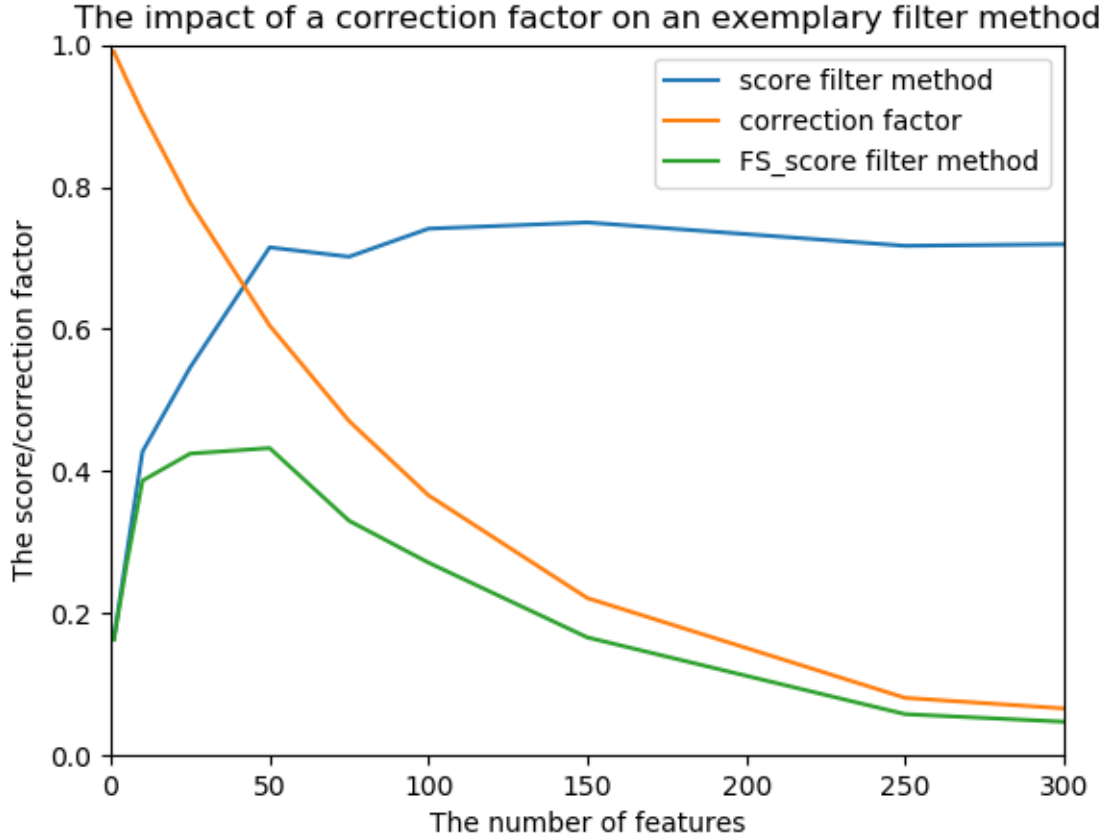


Figure 1.3: An example of the impact of the correction factor on the score, in this case accuracy. The shown correction factor uses $\beta = 0.99$. On the x-axis the number of features is shown and on the y-axis the value for the original accuracy, the correction factor and the *FS_score*.

1.4.2 Feature Selection Exploration

The first set of experiments is used to explore a combination of the basic filter method combined with the quality measurements. The basic filter method algorithm selecting the top n features (Algorithm 1) is used. The changing variables are the dataset, the ranking method, the feature preservation values and the accuracy computation methods (Table 1.3). The range of chosen feature preservation values comes from both its ability to show impact of separate features (more impact from fewer features) and the relevance of keeping that number of features (irrelevant feature selection when having more than 1000 features). All of this together results in a total of $4(\text{datasets}) \times 2(\text{ranking methods}) \times 11(\text{top } n \text{ features}) \times 5(\text{accuracy computation methods}) = 440$ experiments. These experiments are visualized in eight plots, one plot for every combination of data set and ranking method. These plots then show the change in quality for different number of preserved features.

1.4.3 Feature Selection Algorithms Evaluation

The second set of experiments compares the feature selection methods in both feature selection and quality. For this experiment the four datasets are used with the logistic regression quality measurement, since the logistic regression gave the most consistent result of the five machine learning algorithms, showing the smallest variation across the datasets. The average performance for

Table 1.3: The four meta-parameters with their possible values in the first experiment.

Variable	Description	Values
Dataset	The datasets used (subsection 1.3.1)	Psoriasis RSCTC Arcene Micro-Organisms –
Ranking method	The method used for ranking the features (Subsection 1.3.2)	T-test (<i>SciPy</i>) Mutual Information (<i>scikit-learn</i>) –
Feature preservation values	The fixed size of the feature subset after feature selection	1, 5, 10, 25, 50, 75, 100, 150, 250, 500, 1000 –
F1-score computation	The machine learning algorithms used to compute the F1-score (subsection 1.4.1)	Naive Bayes Logistic Regression Support Vector Machine Decision Tree Nearest Neighbours

all datasets is also computed for clarification purposes. An overview of feature selection methods is made (Table 1.4) and all separate meta-parameters are explained here. Spectra are made with the results of the experiment that show the performance of all different combinations, showing the accuracy, precision, recall and F1-score. On top of that the computation time is computed and shown per algorithm, as well, in a separate bar chart. At last a combination of the earlier proposed *FS_score* and the computation time is shown as well for $\alpha = 0.99$ to show the relation between computation time on one hand and the *FS_score* on the other hand. Several discussed algorithms are chosen to be excluded from the evaluation. This choice is based on their poor scalability with regards to computation time and therefore unfit for datasets with this many features. These excluded algorithms are the backwards elimination sequential method, the simulated annealing stochastic search method and the embedded backwards elimination method (Subsection 1.3.2).

1.4.4 TPOT Feature selection integration

TPOT is an effective tool to find the best machine learning pipeline for a certain dataset. Two restrictions hinder optimization regarding feature selection:

1. *Lack of warm start*

TPOT has a vast array of machine learning and preprocessing algorithms to find the best possible pipeline. Due to the number of possibilities being very high, a lot of time may be wasted due to searching in wrong directions. For feature selection, a pre-defined selection of pipelines (also known as a warm start) would improve efficiency.

2. *Feature selection possibilities*

Several filter and embedded methods are present (Subsection 1.3.2). All of these select percentages of feature selection, though. Still a high number of features can be present in the result after using percentages. On top of that, no wrapper methods are available, either.

For both of these optimization restrictions, two additions are made to the *TPOT* code structure for increased functionality with regards to feature selection, in the same order as the restrictions:

1. *Focused feature selection addition*

The possibility is created to always start the original population with a feature selection algorithm in the feature selection pipeline. Due to this start the expected search for a good

Table 1.4: The methods that are evaluated in the second experiment setup.

Type	Method	Parameters
Filter methods	Basic Filter Methods (Algorithm 1)	- Rank: T-test, Mutual Information - Thresholds: 50, 100, 150 features
Wrapper methods	Forward selection (Algorithm 2)	- Order: Random, Mutual Information - Evaluation: Naive Bayes - Alpha: 0.01, 0.001 —
	PTA (Algorithm ??)	- Order: Random, Mutual Information - Evaluation: Naive Bayes - $[l, r] = [20, 10], [5, 2]$ - Alpha: 0.01, 0.001 —
	Floating search (Algorithm ??)	- Order: Random, Mutual Information - Evaluation: Naive Bayes - Alpha: 0.01, 0.001
	Forward selection (Algorithm 4)	- Machine Learning: SVM, RandomForests - Threshold: 50, 100, 150 features
Embedded methods	Forward selection (Algorithm 4)	- Machine Learning: SVM, RandomForests - Threshold: 50, 100, 150 features

feature selection method is bypassed immediately, which should result in more optimized final pipeline.

2. Alternative feature selection algorithm set

The possibility is created to start with an alternative feature selection algorithms set. This selection consists of filter, wrapper and embedded methods and the hyper-parameters are predefined to create an upper bound of 200 features.

The experiment consists of multiple runs of *TPOT* and all of these steps are also shown in an explanatory table (Table 1.5). All four datasets are tested (Subsection 1.3.1) and the accuracy is changed to the feature sensitive *FS_score* with $\beta = 0.99$ (Equation 1.4.1), as previously discussed (Subsection 1.4.1). For testing *TPOT* an optimization time of 120 minutes (two hours) was chosen as a reasonable time constraint to run each experiment 5 times, an algorithm was not allowed to run for longer than ten minutes, a population size of 12 was chosen to not be too selective at the start and a training set size of 0.90 which is a general training set size when not many samples are present. Furthermore tests were done for pipeline selection both with and without focused feature selection and for both with and without alternative feature selection algorithms set, as these additions must be tested for their quality. This gives a total of $4(\text{datasets}) \times 2(\text{pipeline selection}) \times 2(\text{feature selection set}) \times 5(\text{experiment reruns}) = 16$ different experiments.

1.5 Results

The different experiments are all explained in their own subsections. First the results of the minimum feature preservation experiment were shown, followed by the feature selection algorithms evaluation.

1.5.1 Feature Selection Exploration Results

A difference in score quality for the datasets was visible (Figure 1.4). The difference between using Mutual Information and T-test/ANOVA was not, as for only the Arcene dataset there was a significant difference between the two. Therefore it seems that the ranking method type has less influence on the measurement quality. One interesting aspect was that methods using Mutual Information had a longer computation time than methods using the T-test/ANOVA.

Table 1.5: The experiment details for testing the non-trivial changes in *TPOT*. This experiment is rerun 5 times.

Experiment factors	Detailed values	Remarks
Datasets	Micro-organisms Arcene RSCTC Psoriasis	High number of features (Subsection 1.3.1)
Performance measurement	One type: <i>FS_score</i> - $\beta = 0.99$	Addition of correction factor (Equation 1.4.1)
TPOT input parameters	One set of input values: - max. optimization time = 120 min - max. alg. evaluation time = 10 min - pop. size = 12 - train size = 0.9	Explanation of input values: - The time one run should take (2 hours) - The evaluation time of one pipeline - The number of pipelines in one generation - The number of samples used for training
Pipeline selection	Regular selection Feature selection focused	Possible obligatory addition of a feature selection algorithm
Change in feature selection algorithm set	Regular feature selection algorithm set New feature selection algorithm set	A change between several basic feature selection algorithms to feature selection algorithms designed for at most 200 features preservation

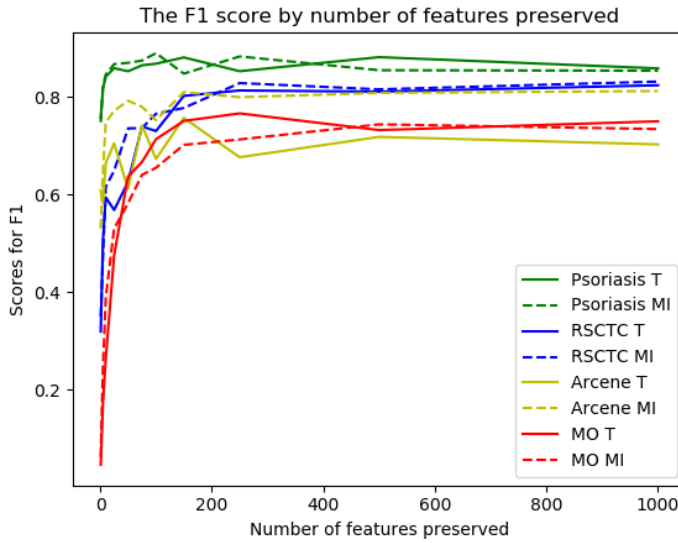


Figure 1.4: The average validation F1-scores shown per dataset and rank.

1.5.2 Feature Selection Algorithms Evaluation Results

The figure showed that all wrapper algorithms preserved on average less than 61 features for these settings, whereas the performance seems to average around 75%. The filter and embedded methods performed worse, with an overall lower performance score than the wrapper methods, even when more features were present. Also, no immediate observations can be made by looking only at the filter and embedded algorithms.

When only looking at the wrapper algorithms, some other observations could be done, as well. Ordering the features before using a wrapper method structurally gave a better result than using

a random ordering. Also a threshold of $\alpha = 0.001$ usually resulted in more features and in a higher scores in comparison with a threshold of $\alpha = 0.01$. Comparing the algorithms, the floating search with ordering did best in performance, whereas the other algorithms are performing closer together.

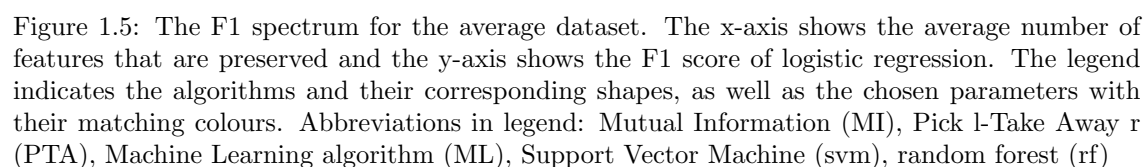


Table 1.6: The performance of the final pipeline for the different types of *TPOT* and the mentioned dataset. The five reruns are averaged into this one result.

Algorithm	selection availability	regular selection		always feature selection	
		regular algorithms	feature selection algorithms	regular algorithms	feature selection algorithms
Datasets	Micro-organisms	0.45	0.71	0.64	0.59
	Arcene	0.46	0.69	0.62	0.69
	RSCTC	0.00	0.22	0.17	0.44
	Psoriasis	0.00	0.24	0.16	0.54

1.5.3 TPOT Feature Selection Integration Results

To visualize the results of the experiment, a table with the final results is made (Table 1.6). The optimization process is also recorded and shown in a plot for the four datasets (Figure 1.6). The five experiment reruns are averaged to one complete result.

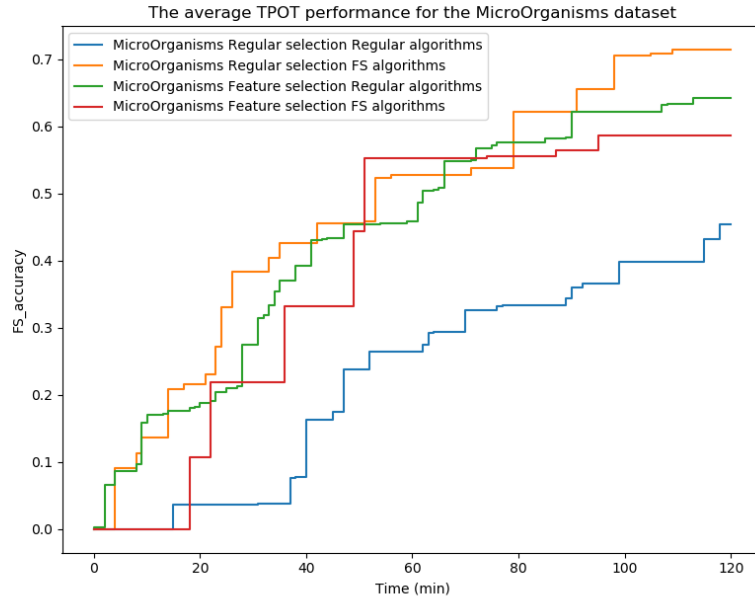


Figure 1.6: The optimization process for the different *TPOT* algorithms for the micro-organisms dataset. - TO DO: Make one for all 4 datasets

The results show that after two hours of running *TPOT*, regular *TPOT* always performs worst. According to the processes the high feature selection algorithm set shows much more big stepwise improvements, accompanied with fewer smaller stepwise improvements. This shows the earlier observed trade-off between computation time and feature selection quality for these specific feature selection algorithms. It takes longer to compute those algorithms, but it also gives a much better result.

The choice for the best *TPOT* algorithm is not conclusive. For high feature selection, all three new possibilities show improved results. The *TPOT* variant that is feature selection focused and has the new feature selection algorithm set performs the best for datasets with a very high number of features (RSCTC and Psoriasis), but the trade-off in computation time seems to hurt the performance for smaller datasets. The algorithm that is only feature selection focused has a

steady performance for datasets with a lower number of features, but this performance does not seem to be better than the algorithm with only the new feature selection algorithm set.

1.6 Discussion

There are several works focused on the possibility of feature selection, both with a data analytical [91] as well as a biomedical [19, 20, 92] point of view. These research projects are usually focused on datasets with a low number of features, for which feature selection is less relevant. This research for datasets that have a significantly higher number of at least 1000 features is new, even though newly available datasets become bigger over time due to new and improved techniques of measuring and storing data. With this research a beginning is made on how to approach feature selection on these bigger datasets.

Whereas multiple feature selection methods are discussed, several were not tested due to computation time constraints. The wrapper methods backwards elimination and simulated annealing and the embedded backwards elimination all were too computationally intensive to become relevant for the research. In datasets with fewer features, these methods may be showing better results and could be possible candidates for feature selection.

After evaluation of the results in the first experiment set-up, it could be concluded that there was not a big difference between using T-test/ANOVA or Mutual Information as a ranking method. Both gave similar results, with the exception of one dataset.

A second conclusion was drawn from looking at the accuracy with the number of features preserved. After a threshold of 200 features, additional features did not raise the validation score as much as the first 200 features did. This indicated a second rule of thumb, that at least 200 features should be preserved after using a filter method. This rule of thumb was according to the hypothesis H1 notion that at most 1000 features should be preserved to still show an efficient result, however was much lower than the expected upper bound of 1000 features. This absolute number was not intuitive and much lower than expected. More research could be done on these preserved 200 features and why no other features were needed to predict the output. Cluster analyses for example may show insights in this phenomenon.

After evaluation of both filter, wrapper and embedded methods, wrapper methods were significantly better at selecting a smaller fraction of features while preserving a similar test score. Whereas hypothesis H2 stated wrapper methods to be best in efficiency, embedded methods were expected to outperform both filter methods and wrapper methods when looking at both quality and computation time. This was not the case, however, as embedded methods had a near identical quality to filter methods and performed much worse than wrapper methods. The outcome indicated that wrapper methods are significantly different in results and filter and embedded methods being nearly identical.

Since wrapper methods take dependencies between features into account, they kept these dependencies at a minimum. If these dependencies are unwanted, wrapper methods are more useful than filter methods and embedded and should be recommended. A downside of the wrapper methods however was that they took much more computation time than the filter and embedded methods. Therefore if it does not matter when features have dependencies with each other, a filter method should be recommended.

There was a difference in quality within the wrapper methods. The forward selection and PTA all showed promising results and therefore should be considered for the framework. Floating search showed the best results, however took significantly longer in computation time. Within the wrapper methods, an ordering beforehand showed improved results. Backward elimination wrapper algorithms, stochastic search algorithms and embedded backwards elimination algorithms all were significantly worse in computation time than the other wrapper algorithms and therefore should not be considered in these cases.

A recommendation for the threshold in wrapper methods is not trivial. A higher threshold of $\alpha = 0.01$ gave a smaller feature subset at the cost of a lower classification score. If a smaller subset is desired or more influential features are needed, a bigger threshold should be chosen, whereas it

will be smaller when a higher quality of the feature subset is desired.

Both additions to *TPOT* were beneficial according to the experiment. This was also expected, as H3 states. When the initial number of features was very high (50.000), a combination of initial bias towards feature selection and using feature selection algorithms focusing on preservation of only 200 features performed best. For lower number of features only adding one of the improvements was better, as computation time became the limiting factor in case of using both improvements.

Bibliography

- [1] N. Gehlenborg, S. I. O'donoghue, N. S. Baliga, A. Goesmann, M. A. Hibbs, H. Kitano, O. Kohlbacher, H. Neuweber, R. Schneider, D. Tenenbaum, *et al.*, "Visualization of omics data for systems biology," *Nature methods*, vol. 7, no. 3s, p. S56, 2010.
- [2] A. Brazma, P. Hingamp, J. Quackenbush, G. Sherlock, P. Spellman, C. Stoeckert, J. Aach, W. Ansorge, C. A. Ball, H. C. Causton, *et al.*, "Minimum information about a microarray experiment (miame)—toward standards for microarray data," *Nature genetics*, vol. 29, no. 4, p. 365, 2001.
- [3] J. S. Cottrell and U. London, "Probability-based protein identification by searching sequence databases using mass spectrometry data," *electrophoresis*, vol. 20, no. 18, pp. 3551–3567, 1999.
- [4] K. Dettmer, P. A. Aronov, and B. D. Hammock, "Mass spectrometry-based metabolomics," *Mass spectrometry reviews*, vol. 26, no. 1, pp. 51–78, 2007.
- [5] B. Liu, X. Zhou, Y. Wang, J. Hu, L. He, R. Zhang, S. Chen, and Y. Guo, "Data processing and analysis in real-world traditional chinese medicine clinical data: challenges and approaches," *Statistics in medicine*, vol. 31, no. 7, pp. 653–660, 2012.
- [6] D. F. Sittig, A. Wright, J. A. Osheroff, B. Middleton, J. M. Teich, J. S. Ash, E. Campbell, and D. W. Bates, "Grand challenges in clinical decision support," *Journal of biomedical informatics*, vol. 41, no. 2, pp. 387–392, 2008.
- [7] G. Magni, C. Caldieron, S. Rigatti-Luchini, and H. Merskey, "Chronic musculoskeletal pain and depressive symptoms in the general population. an analysis of the 1st national health and nutrition examination survey data," *Pain*, vol. 43, no. 3, pp. 299–307, 1990.
- [8] P. Bertolazzi, G. Felici, P. Festa, and G. Lancia, "Logic classification and feature selection for biomedical data," *Computers & Mathematics with Applications*, vol. 55, no. 5, pp. 889–899, 2008.
- [9] G. Piatetsky-Shapiro and P. Tamayo, "Microarray data mining: facing the challenges," *ACM SIGKDD Explorations Newsletter*, vol. 5, no. 2, pp. 1–5, 2003.
- [10] A. Lommen, "Metalign: interface-driven, versatile metabolomics tool for hyphenated full-scan mass spectrometry data preprocessing," *Analytical chemistry*, vol. 81, no. 8, pp. 3079–3086, 2009.
- [11] A. Holzinger, M. Dehmer, and I. Jurisica, "Knowledge discovery and interactive data mining in bioinformatics-state-of-the-art, future challenges and research directions," *BMC bioinformatics*, vol. 15, no. 6, p. 11, 2014.
- [12] M. Wilkins, "Proteomics data mining," *Expert review of proteomics*, vol. 6, no. 6, pp. 599–603, 2009.

- [13] D. Teodoro, R. Choquet, E. Pasche, J. Gobeill, C. Daniel, P. Ruch, and C. Lovis, “Biomedical data management: a proposal framework,” in *MIE*, pp. 175–179, Citeseer, 2009.
- [14] M. Y. Galperin, “The molecular biology database collection: 2008 update,” *Nucleic Acids Research*, vol. 36, no. suppl1, pp. D2–D4, 2008.
- [15] A. Sturn, J. Quackenbush, and Z. Trajanoski, “Genesis: cluster analysis of microarray data,” *Bioinformatics*, vol. 18, no. 1, pp. 207–208, 2002.
- [16] A. Karnovsky, T. Weymouth, T. Hull, V. G. Tarcea, G. Scardoni, C. Laudanna, M. A. Sartor, K. A. Stringer, H. Jagadish, C. Burant, *et al.*, “Metscape 2 bioinformatics tool for the analysis and visualization of metabolomics and gene expression data,” *Bioinformatics*, vol. 28, no. 3, pp. 373–380, 2011.
- [17] D. Tabas-Madrid, R. Nogales-Cadenas, and A. Pascual-Montano, “Genecodis3: a non-redundant and modular enrichment analysis tool for functional genomics,” *Nucleic acids research*, vol. 40, no. W1, pp. W478–W483, 2012.
- [18] F. Faul, E. Erdfelder, A.-G. Lang, and A. Buchner, “G* power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences,” *Behavior research methods*, vol. 39, no. 2, pp. 175–191, 2007.
- [19] R. Baumgartner and R. L. Somorjai, “Data complexity assessment in undersampled classification of high-dimensional biomedical data,” *Pattern Recognition Letters*, vol. 27, no. 12, pp. 1383–1389, 2006.
- [20] W. Welthagen, R. A. Shellie, J. Spranger, M. Ristow, R. Zimmermann, and O. Fiehn, “Comprehensive two-dimensional gas chromatography–time-of-flight mass spectrometry (gc \times gc-tof) for high resolution metabolomics: biomarker discovery on spleen tissue extracts of obese nzo compared to lean c57bl/6 mice,” *Metabolomics*, vol. 1, no. 1, pp. 65–73, 2005.
- [21] A. R. T. Donders, G. J. Van Der Heijden, T. Stijnen, and K. G. Moons, “A gentle introduction to imputation of missing values,” *Journal of clinical epidemiology*, vol. 59, no. 10, pp. 1087–1091, 2006.
- [22] M. H. Cartwright, M. J. Shepperd, and Q. Song, “Dealing with missing software project data,” in *Software Metrics Symposium, 2003. Proceedings. Ninth International*, pp. 154–165, IEEE, 2003.
- [23] J. S. Haukoos and C. D. Newgard, “Advanced statistics: missing data in clinical research—part 1: an introduction and conceptual framework,” *Academic Emergency Medicine*, vol. 14, no. 7, pp. 662–668, 2007.
- [24] R. Deneer, “Scoring co-morbidity severity in bariatric patients based on biomarkers: a data mining approach,” 2017.
- [25] R. P. Nair, K. C. Duffin, C. Helms, J. Ding, P. E. Stuart, D. Goldgar, J. E. Gudjonsson, Y. Li, T. Tejasvi, B.-J. Feng, *et al.*, “Genome-wide scan reveals association of psoriasis with il-23 and nf- κ b pathways,” *Nature genetics*, vol. 41, no. 2, pp. 199–204, 2009.
- [26] M. Suárez-Farinas, K. Li, J. Fuentes-Duculan, K. Hayden, C. Brodmerkel, and J. G. Krueger, “Expanding the psoriasis disease profile: interrogation of the skin and serum of patients with moderate-to-severe psoriasis,” *Journal of Investigative Dermatology*, vol. 132, no. 11, pp. 2552–2564, 2012.
- [27] J. Bigler, H. A. Rand, K. Kerkof, M. Timour, and C. B. Russell, “Cross-study homogeneity of psoriasis gene expression in skin across a large expression range,” *PLoS One*, vol. 8, no. 1, p. e52242, 2013.

- [28] J. Kim, R. Bissonnette, J. Lee, J. C. da Rosa, M. Suárez-Fariñas, M. A. Lowes, and J. G. Krueger, “The spectrum of mild to severe psoriasis vulgaris is defined by a common activation of il-17 pathway genes, but with key differences in immune regulatory genes,” *Journal of Investigative Dermatology*, vol. 136, no. 11, pp. 2173–2182, 2016.
- [29] Y. Yao, L. Richman, C. Morehouse, M. De Los Reyes, B. W. Higgs, A. Boutrín, B. White, A. Coyle, J. Krueger, P. A. Kiener, *et al.*, “Type i interferon: potential therapeutic target for psoriasis?,” *PloS one*, vol. 3, no. 7, p. e2737, 2008.
- [30] M. Suárez-Fariñas, S. J. Tintle, A. Shemer, A. Chiricozzi, K. Nogales, I. Cardinale, S. Duan, A. M. Bowcock, J. G. Krueger, and E. Guttman-Yassky, “Nonlesional atopic dermatitis skin is characterized by broad terminal differentiation defects and variable immune abnormalities,” *Journal of Allergy and Clinical Immunology*, vol. 127, no. 4, pp. 954–964, 2011.
- [31] S. Tintle, A. Shemer, M. Suárez-Fariñas, H. Fujita, P. Gilleaudeau, M. Sullivan-Whalen, L. Johnson-Huang, A. Chiricozzi, I. Cardinale, S. Duan, *et al.*, “Reversal of atopic dermatitis with narrow-band uvb phototherapy and biomarkers for therapeutic response,” *Journal of Allergy and Clinical Immunology*, vol. 128, no. 3, pp. 583–593, 2011.
- [32] J. K. Gittler, A. Shemer, M. Suárez-Fariñas, J. Fuentes-Duculan, K. J. Gulewicz, C. Q. Wang, H. Mitsui, I. Cardinale, C. de Guzman Strong, J. G. Krueger, *et al.*, “Progressive activation of t h 2/t h 22 cytokines and selective epidermal proteins characterizes acute and chronic atopic dermatitis,” *Journal of Allergy and Clinical Immunology*, vol. 130, no. 6, pp. 1344–1354, 2012.
- [33] K. J. Cios and G. W. Moore, “Uniqueness of medical data mining,” *Artificial Intelligence in Medicine*, vol. 26, no. 1, pp. 1 – 24, 2002. Medical Data Mining and Knowledge Discovery.
- [34] I. Yoo, P. Alafaireet, M. Marinov, K. Pena-Hernandez, R. Gopidi, J.-F. Chang, and L. Hua, “Data mining in healthcare and biomedicine: A survey of the literature,” *Journal of Medical Systems*, vol. 36, pp. 2431–2448, Aug 2012.
- [35] H. Chen, S. S. Fuller, C. Friedman, and W. Hersh, *Medical informatics: knowledge management and data mining in biomedicine*, vol. 8. Springer Science & Business Media, 2006.
- [36] R. Miotto, F. Wang, S. Wang, X. Jiang, and J. T. Dudley, “Deep learning for healthcare: review, opportunities and challenges,” *Briefings in Bioinformatics*, p. bbx044, 2017.
- [37] D. Blythe, “Rise of the graphics processor,” *Proceedings of the IEEE*, vol. 96, no. 5, pp. 761–778, 2008.
- [38] C. Turkyay, F. Jeanquartier, A. Holzinger, and H. Hauser, *On Computationally-Enhanced Visual Analysis of Heterogeneous Data and Its Application in Biomedical Informatics*, pp. 117–140. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014.
- [39] A. Holzinger and I. Jurisica, *Knowledge Discovery and Data Mining in Biomedical Informatics: The Future Is in Integrative, Interactive Machine Learning Solutions*, pp. 1–18. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014.
- [40] W. Dubitzky, M. Granzow, and D. P. Berrar, *Fundamentals of data mining in genomics and proteomics*. Springer Science & Business Media, 2007.
- [41] Y. Peng, Z. Wu, and J. Jiang, “A novel feature selection approach for biomedical data classification,” *Journal of Biomedical Informatics*, vol. 43, no. 1, pp. 15 – 23, 2010.
- [42] D. Dunbar and G. Humphreys, “A spatial data structure for fast poisson-disk sample generation,” in *ACM Transactions on Graphics (TOG)*, vol. 25, pp. 503–508, ACM, 2006.

- [43] L. Devroye, "Sample-based non-uniform random variate generation," in *Proceedings of the 18th conference on Winter simulation*, pp. 260–265, ACM, 1986.
- [44] L. J. Van't Veer, H. Dai, M. J. Van De Vijver, Y. D. He, A. A. Hart, M. Mao, H. L. Peterse, K. Van Der Kooy, M. J. Marton, A. T. Witteveen, *et al.*, "Gene expression profiling predicts clinical outcome of breast cancer," *nature*, vol. 415, no. 6871, p. 530, 2002.
- [45] D. A. Roff and P. Bentzen, "The statistical analysis of mitochondrial dna polymorphisms: chi 2 and the problem of small samples.," *Molecular biology and evolution*, vol. 6, no. 5, pp. 539–545, 1989.
- [46] R. Bellazzi, M. Diomidous, I. N. Sarkar, K. Takabayashi, A. Ziegler, A. T. McCray, *et al.*, "Data analysis and data mining: current issues in biomedical informatics," *Methods of information in medicine*, vol. 50, no. 6, p. 536, 2011.
- [47] D. Otasek, C. Pastrello, A. Holzinger, and I. Jurisica, *Visual Data Mining: Effective Exploration of the Biological Universe*, pp. 19–33. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014.
- [48] L. Marenco, T.-Y. Wang, G. Shepherd, P. L. Miller, and P. Nadkarni, "Qis: A framework for biomedical database federation," *Journal of the American Medical Informatics Association*, vol. 11, no. 6, pp. 523–534, 2004.
- [49] V. Y. Bichutskiy, R. Colman, R. K. Brachmann, and R. H. Lathrop, "Heterogeneous biomedical database integration using a hybrid strategy: a p53 cancer research database," *Cancer informatics*, vol. 2, p. 277, 2006.
- [50] W. Sperzel, R. Abarbanel, S. Nelson, M. Erlbaum, D. Sherertz, M. Tuttle, N. Olson, and L. Fuller, "Biomedical database inter-connectivity: an experiment linking mim, genbank, and meta-1 via medline.," in *Proceedings of the Annual Symposium on Computer Application in Medical Care*, p. 190, American Medical Informatics Association, 1991.
- [51] F. Aubry, S. Badaoui, H. Kaplan, and R. D. Paola, "Design and implementation of a biomedical image database (bdim)," *Medical Informatics*, vol. 13, no. 4, pp. 241–248, 1988.
- [52] D. Windridge and M. Bober, *A Kernel-Based Framework for Medical Big-Data Analytics*, pp. 197–208. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014.
- [53] S. Selvaraj and J. Natarajan, "Microarray data analysis and mining tools," *Bioinformatics*, vol. 6, no. 3, p. 95, 2011.
- [54] M. Afzal, I. Manzoor, and O. P. Kuipers, "A fast and reliable pipeline for bacterial transcriptome analysis case study: serine-dependent gene regulation in streptococcus pneumoniae," *Journal of visualized experiments: JoVE*, no. 98, 2015.
- [55] M. Wojnarski, A. Janusz, H. S. Nguyen, J. Bazan, C. Luo, Z. Chen, F. Hu, G. Wang, L. Guan, H. Luo, *et al.*, "Rsctc'2010 discovery challenge: Mining dna microarray data for medical diagnosis and treatment," in *International Conference on Rough Sets and Current Trends in Computing*, pp. 4–19, Springer, 2010.
- [56] R. Matthiesen and O. N. Jensen, "Analysis of mass spectrometry data in proteomics," in *Bioinformatics*, pp. 105–122, Springer, 2008.
- [57] J. T. Watson and O. D. Sparkman, *Introduction to mass spectrometry: instrumentation, applications, and strategies for data interpretation*. John Wiley & Sons, 2007.
- [58] A. C. Neves, C. L. Morais, T. P. Mendes, B. G. Vaz, and K. M. Lima, "Mass spectrometry and multivariate analysis to classify cervical intraepithelial neoplasia from blood plasma: an untargeted lipidomic study," *Scientific reports*, vol. 8, no. 1, p. 3954, 2018.

- [59] I. Guyon, S. Gunn, A. Ben-Hur, and G. Dror, “Result analysis of the nips 2003 feature selection challenge,” in *Advances in Neural Information Processing Systems 17* (L. K. Saul, Y. Weiss, and L. Bottou, eds.), pp. 545–552, MIT Press, 2005.
- [60] P. Mahé, M. Arsac, S. Chatellier, V. Monnin, N. Perrot, S. Mailler, V. Girard, M. Ramjeet, J. Surre, B. Lacroix, A. van Belkum, and J.-B. Veyrieras, “Automatic identification of mixed bacterial species fingerprints in a maldi-tof mass-spectrum,” *Bioinformatics*, vol. 30, no. 9, pp. 1280–1286, 2014.
- [61] M. T. Madigan, J. M. Martinko, J. Parker, *et al.*, *Brock biology of microorganisms*, vol. 13. Pearson, 2017.
- [62] S. J. Pocock, *Clinical trials: a practical approach*. John Wiley & Sons, 2013.
- [63] A. Bendele, T. McAbee, G. Sennello, J. Frazier, E. Chlipala, and D. McCabe, “Efficacy of sustained blood levels of interleukin-1 receptor antagonist in animal models of arthritis: comparison of efficacy in animal models with human clinical data,” *Arthritis & Rheumatism: Official Journal of the American College of Rheumatology*, vol. 42, no. 3, pp. 498–506, 1999.
- [64] G. Kan, C. A. Visser, J. J. Koolen, and A. J. Dunning, “Short and long term predictive value of admission wall motion score in acute myocardial infarction. a cross sectional echocardiographic study of 345 patients,” *Heart*, vol. 56, no. 5, pp. 422–427, 1986.
- [65] P. Diaconis and B. Efron, “Computer-intensive methods in statistics,” *Scientific American*, vol. 248, no. 5, pp. 116–131, 1983.
- [66] B. Cestnik, “Kononenkoj., bratkoj.(1987): Assistant-86: A knowledge elicitation tool for sophisticated users,” *Progress in machine learning*.
- [67] P. A. Murtaugh, E. R. Dickson, G. M. Van Dam, M. Malinchoc, P. M. Grambsch, A. L. Langworthy, and C. H. Gips, “Primary biliary cirrhosis: prediction of short-term survival based on repeated patient visits,” *Hepatology*, vol. 20, no. 1, pp. 126–134, 1994.
- [68] L. Misery, V. Sibaud, C. Merial-Kieny, and C. Taieb, “Sensitive skin in the american population: prevalence, clinical data, and role of the dermatologist,” *International journal of dermatology*, vol. 50, no. 8, pp. 961–967, 2011.
- [69] K. Fernandes, J. S. Cardoso, and J. Fernandes, “Transfer learning with partial observability applied to cervical cancer screening,” in *Iberian conference on pattern recognition and image analysis*, pp. 243–250, Springer, 2017.
- [70] W. McKinney *et al.*, “Data structures for statistical computing in python,” in *Proceedings of the 9th Python in Science Conference*, vol. 445, pp. 51–56, SciPy Austin, TX, 2010.
- [71] Y. Yan and J. Yan, “Hands-on data science with anaconda: Utilize the right mix of tools to create high-performance data science applications,” 2018.
- [72] S. v. d. Walt, S. C. Colbert, and G. Varoquaux, “The numpy array: a structure for efficient numerical computation,” *Computing in Science & Engineering*, vol. 13, no. 2, pp. 22–30, 2011.
- [73] E. Jones, T. Oliphant, and P. Peterson, “{SciPy}: open source scientific tools for {Python},” 2014.
- [74] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, *et al.*, “Scikit-learn: Machine learning in python,” *Journal of Machine Learning Research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [75] W. McKinney, *Python for data analysis: Data wrangling with Pandas, NumPy, and IPython*. ” O’Reilly Media, Inc.” , 2012.

- [76] R. S. Olson and J. H. Moore, "Tpot: A tree-based pipeline optimization tool for automating machine learning," in *Workshop on Automatic Machine Learning*, pp. 66–74, 2016.
- [77] L. Yu and H. Liu, "Feature selection for high-dimensional data: A fast correlation-based filter solution," in *Proceedings of the 20th international conference on machine learning (ICML-03)*, pp. 856–863, 2003.
- [78] M. A. Hall, "Correlation-based feature selection of discrete and numeric class machine learning," 2000.
- [79] R. Kohavi *et al.*, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *Ijcai*, vol. 14, pp. 1137–1145, Montreal, Canada, 1995.
- [80] M. Kearns and D. Ron, "Algorithmic stability and sanity-check bounds for leave-one-out cross-validation," *Neural computation*, vol. 11, no. 6, pp. 1427–1453, 1999.
- [81] A. Agresti, *Categorical data analysis*, vol. 482. John Wiley & Sons, 2003.
- [82] A. Edwards, G. Elwyn, and A. Mulley, "Explaining risks: turning numerical data into meaningful pictures," *Bmj*, vol. 324, no. 7341, pp. 827–830, 2002.
- [83] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*. Elsevier, 2011.
- [84] I. S. Lim, P. de Heras Ciechomski, S. Sarni, and D. Thalmann, "Planar arrangement of high-dimensional biomedical data sets by isomap coordinates," in *Computer-Based Medical Systems, 2003. Proceedings. 16th IEEE Symposium*, pp. 50–55, IEEE, 2003.
- [85] J. Biesiada and W. Duch, "Feature selection for high-dimensional data—a pearson redundancy based filter," in *Computer recognition systems 2*, pp. 242–249, Springer, 2007.
- [86] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 27, no. 8, pp. 1226–1238, 2005.
- [87] R. Battiti, "Using mutual information for selecting features in supervised neural net learning," *IEEE Transactions on neural networks*, vol. 5, no. 4, pp. 537–550, 1994.
- [88] N. J.-M. Blackman and J. J. Koval, "Interval estimation for cohen's kappa as a measure of agreement," *Statistics in medicine*, vol. 19, no. 5, pp. 723–741, 2000.
- [89] Y. Peng, Z. Wu, and J. Jiang, "A novel feature selection approach for biomedical data classification," *Journal of Biomedical Informatics*, vol. 43, no. 1, pp. 15–23, 2010.
- [90] C. Ding and H. Peng, "Minimum redundancy feature selection from microarray gene expression data," *Journal of bioinformatics and computational biology*, vol. 3, no. 02, pp. 185–205, 2005.
- [91] C. Catal and B. Diri, "Investigating the effect of dataset size, metrics sets, and feature selection techniques on software fault prediction problem," *Information Sciences*, vol. 179, no. 8, pp. 1040–1058, 2009.
- [92] H. Liu, J. Li, and L. Wong, "A comparative study on feature selection and classification methods using gene expression profiles and proteomic patterns," *Genome informatics*, vol. 13, pp. 51–60, 2002.
- [93] I. Guyon and A. Elisseeff, *An Introduction to Feature Extraction*, pp. 1–25. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- [94] C. Catal and B. Diri, "Investigating the effect of dataset size, metrics sets, and feature selection techniques on software fault prediction problem," *Information Sciences*, vol. 179, no. 8, pp. 1040 – 1058, 2009.

- [95] L. C. Molina, L. Belanche, and À. Nebot, “Feature selection algorithms: A survey and experimental evaluation,” in *Data Mining, 2002. ICDM 2003. Proceedings. 2002 IEEE International Conference on*, pp. 306–313, IEEE, 2002.
- [96] Y. Saeys, I. Inza, and P. Larrañaga, “A review of feature selection techniques in bioinformatics,” *bioinformatics*, vol. 23, no. 19, pp. 2507–2517, 2007.
- [97] W. Duch, *Filter Methods*, pp. 89–117. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- [98] R. M. Heiberger and B. Holland, *Statistical analysis and data display*. Springer, 2004.
- [99] D. Donoho and J. Jin, “Higher criticism thresholding: Optimal feature selection when useful features are rare and weak,” *Proceedings of the National Academy of Sciences*, vol. 105, no. 39, pp. 14790–14795, 2008.
- [100] J. D. Storey and R. Tibshirani, “Statistical significance for genomewide studies,” *Proceedings of the National Academy of Sciences*, vol. 100, no. 16, pp. 9440–9445, 2003.
- [101] J. P. Higgins, S. G. Thompson, J. J. Deeks, and D. G. Altman, “Measuring inconsistency in meta-analyses,” *BMJ: British Medical Journal*, vol. 327, no. 7414, p. 557, 2003.
- [102] J. Reunanen, *Search Strategies*, pp. 119–136. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- [103] B. Alsallakh and L. Ren, “Powerset: A comprehensive visualization of set intersections,” vol. 23, pp. 1–1, 01 2016.
- [104] I. Tsamardinos, G. Borboudakis, P. Katsogridakis, P. Pratikakis, and V. Christophides, “Massively-parallel feature selection for big data,” *arXiv preprint arXiv:1708.07178*, 2017.
- [105] Y.-J. Huang, D.-Y. Chan, D.-C. Cheng, Y.-J. Ho, P.-P. Tsai, W.-C. Shen, and R.-F. Chen, “Automated feature set selection and its application to mcc identification in digital mammograms for breast cancer detection,” *Sensors*, vol. 13, no. 4, pp. 4855–4875, 2013.
- [106] A. Senawi, H.-L. Wei, and S. A. Billings, “A new maximum relevance-minimum multicollinearity (mrmc) method for feature selection and ranking,” *Pattern Recognition*, vol. 67, pp. 47 – 61, 2017.
- [107] A. El Akadi, A. Amine, A. El Ouardighi, and D. Aboutajdine, “A new gene selection approach based on minimum redundancy-maximum relevance (mrmr) and genetic algorithm (ga),” in *Computer Systems and Applications, 2009. AICCSA 2009. IEEE/ACS International Conference on*, pp. 69–75, IEEE, 2009.
- [108] M. Radovic, M. Ghalwash, N. Filipovic, and Z. Obradovic, “Minimum redundancy maximum relevance feature selection approach for temporal gene expression data,” *BMC bioinformatics*, vol. 18, no. 1, p. 9, 2017.
- [109] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, “Optimization by simulated annealing,” *science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [110] T. Jirapech-Umpai and S. Aitken, “Feature selection and classification for microarray data analysis: Evolutionary methods for identifying predictive genes,” *BMC Bioinformatics*, vol. 6, p. 148, Jun 2005.
- [111] T. N. Lal, O. Chapelle, J. Weston, and A. Elisseeff, *Embedded Methods*, pp. 137–165. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- [112] A. L. Blum and P. Langley, “Selection of relevant features and examples in machine learning,” *Artificial intelligence*, vol. 97, no. 1-2, pp. 245–271, 1997.

- [113] K. Jong, E. Marchiori, M. Sebag, and A. Van Der Vaart, "Feature selection in proteomic pattern data with support vector machines," in *Computational Intelligence in Bioinformatics and Computational Biology, 2004. CIBCB'04. Proceedings of the 2004 IEEE Symposium on*, pp. 41–48, IEEE, 2004.
- [114] J. Prados, A. Kalousis, J.-C. Sanchez, L. Allard, O. Carrette, and M. Hilario, "Mining mass spectra for diagnosis and biomarker discovery of cerebral accidents," *Proteomics*, vol. 4, no. 8, pp. 2320–2332, 2004.
- [115] X. Zhang, X. Lu, Q. Shi, X.-q. Xu, E. L. Hon-chiu, L. N. Harris, J. D. Iglehart, A. Miron, J. S. Liu, and W. H. Wong, "Recursive svm feature selection and sample classification for mass-spectrometry and microarray data," *BMC bioinformatics*, vol. 7, no. 1, p. 197, 2006.
- [116] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Machine learning*, vol. 46, no. 1-3, pp. 389–422, 2002.
- [117] P. Geurts, M. Fillet, D. De Seny, M.-A. Meuwis, M. Malaise, M.-P. Merville, and L. Wehenkel, "Proteomic mass spectra classification using decision tree based ensemble methods," *Bioinformatics*, vol. 21, no. 14, pp. 3138–3145, 2005.
- [118] B. Wu, T. Abbott, D. Fishman, W. McMurray, G. Mor, K. Stone, D. Ward, K. Williams, and H. Zhao, "Comparison of statistical methods for classification of ovarian cancer using mass spectrometry data," *Bioinformatics*, vol. 19, no. 13, pp. 1636–1643, 2003.
- [119] A. Liaw, M. Wiener, *et al.*, "Classification and regression by randomforest," *R news*, vol. 2, no. 3, pp. 18–22, 2002.
- [120] C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown, "Auto-weka: Combined selection and hyperparameter optimization of classification algorithms," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 847–855, ACM, 2013.
- [121] P. Gijsbers, "Automatic construction of machine learning pipelines," 2017.
- [122] L. Kotthoff, C. Thornton, H. H. Hoos, F. Hutter, and K. Leyton-Brown, "Auto-weka 2.0: Automatic model selection and hyperparameter optimization in weka," *Journal of Machine Learning Research*, vol. 17, pp. 1–5, 2016.
- [123] J. R. Koza, "Genetic programming," 1997.
- [124] P. Brazdil, J. Gama, and B. Henery, "Characterizing the applicability of classification algorithms using meta-level learning," in *European conference on machine learning*, pp. 83–102, Springer, 1994.
- [125] R. Vilalta, C. G. Giraud-Carrier, P. Brazdil, and C. Soares, "Using meta-learning to support data mining," *IJCSA*, vol. 1, no. 1, pp. 31–45, 2004.
- [126] P. Brazdil, C. Giraud-Carrier, C. Soares, and R. Vilalta, "Development of metalearning systems for algorithm recommendation," *Metalearning: Applications to Data Mining*, pp. 31–59, 2009.
- [127] M. Pazzani and D. Billsus, "Learning and revising user profiles: The identification of interesting web sites," *Machine learning*, vol. 27, no. 3, pp. 313–331, 1997.
- [128] Y.-W. Chen and C.-J. Lin, "Combining svms with various feature selection strategies," in *Feature extraction*, pp. 315–324, Springer, 2006.
- [129] M. A. Hall and L. A. Smith, "Practical feature subset selection for machine learning," 1998.