

Machine Learning for Fish Oil Analysis [★]

Jesse Wood¹[0000–1111–2222–3333], Bach Hoai Nguyen¹[1111–2222–3333–4444],
Bing Xue¹[2222–3333–4444–5555], Mengjie Zhang¹[2222–3333–4444–5555], and
Daniel Killeen²[2222–3333–4444–5555]

¹ Victoria University of Wellington, Te Herenge Waka, PO Box 600, Wellington 6140,
New Zealand

`jesse.wood@ecs.vuw.ac.nz`

`bach.nguyen@ecs.vuw.ac.nz`

`bing.xue@ecs.vuw.ac.nz`

`mengjie.zhang@ecs.vuw.ac.nz`

² Plant and Food Research, Port Nelson, Nelson 7010, New Zealand

`Daniel.Killeen@plantandfood.co.nz`

Abstract. Gas chromatography (GC) can be used to identify chemical compounds present within tissue samples for quality assurance in food science. Existing analytical chemistry techniques for processing GC data are manual and time-consuming. Here, we explore classification algorithms for fish oil data that automate and significantly reduce the time required to process GC data. We find the Linear SVC model can predict the fish species with near-perfect accuracy. Visualisation is used to explore the interpretability of the models such that their efficacy can be verified for use in a factory setting. The fish oil data is high-dimensional and low sample size. We compare state-of-the-art feature selection methods to reduce the dimensionality of the data. High accuracy is possible with very few features for the MRMR and ReliefF feature selection methods. The exploration reveals there are many feature subsets all capable of producing high-accuracy predictions. No clear superset of important features emerges, which indicates there are many important features to choose from.

Keywords: Feature Selection · Gas Chromatography · Support Vector Machines · Food Science

1 Introduction

Gas chromatography [?] is an chemistry technique we use to analyze fish oils. It can determine the structure of chemical compounds present in a given sample [?]. This is important for quality assurance in a factory setting, especially in food science. We want to be confident that our food labels are accurate and reduce/eliminate cross contamination between different food products. To identify cross-contamination we use fish classification. Given a fish oil sample, we

[★] Supported by organization Plant and Food Research.

can identify the fish species (i.e. Bluecod, Tarakihi), and part (Head, Fins). The existing techniques for performing fish classification are time consuming and laborious. Chemists compare a given sample to reference samples to determine which class it likely belongs to. Previous work on gas chromatography [?,?], has shown machine learning can be used to automate classification.

In this paper we explore machine learning techniques to automate the process of identifying fish species and part on GC data. Firstly, classification algorithms are evaluated for their ability to determine the fish species and part. Visualisation is used to explore the interpretability of successful models. It is important to verify their efficacy with domain knowledge before these algorithms can be deployed in a real-world setting. Secondly, feature selection is used to eliminate redundant features, whilst maintaining high-accuracy predictions.

Specifically, our work is divided into two main sections:

1. Classification Algorithms
2. Feature Selection
3. Visualisation

2 Background

This paper is a multi-disciplinary effort. Domain expertise in chemistry and machine learning is required to extract knowledge from the fish oil data. Before we explore the data, here we provide domain knowledge required to understand this paper.

Specially, the background covers:

1. Chromatography methods: how the raw fish oil data is collected.
2. Classification algorithms: introduce classification algorithms used in the paper.
3. Feature selection: main concepts.

2.1 Gas Chromatography

Gas chromatography (GC) is a technique for the analysis of chemical compounds [?,?,?]. The process separates compounds based on their boiling point and molecular weight. A compound is injected as a liquid, then heat is applied to vaporize it into a gas. A boiling point is a temperature which it changes phase, from liquid to gas. This process is referred to as a phase transition. The speed at which a compound is vaporized depends on its boiling point. The vaporized gases travel through a long coiled tube. That tube has a detector at the end, this detects the rate and intensity which compounds reach the tube's end.

Chemists use the chromatograms of known compounds as a reference when classifying new ones. Take a known example, say *methyl eladiate*. We compare this reference sample to an unknown sample. Analysis can infer the unknown compound since they share the same peaks as the known one. GC is not a

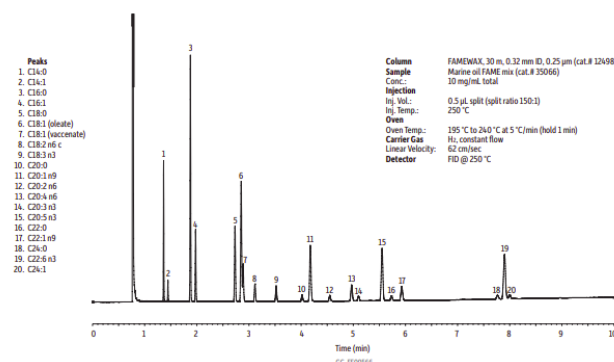


Fig. 1: Gas chromatograph: the artifact of the GC method [?]. The detection is used to visualize intensity (y) and time (x) on a chromatograph.

definitive technique [?], so it is often used in conjunction with other techniques. Mass spectrometry is one such technique [?].

The existing task of classifying chemical compounds based on a chromatograph is laborious [?,?]. The spikes on the graph represent peaks. Each peak represents a resolved chemical compound. Chemists integrate the area under each peak, and compare this to a reference sample, to classify the compound. GC must be performed slowly to ensure that the peaks are not too broad. This ensures each peak resolves and represents a single compound. Once we know what compounds are present in a sample it becomes possible to identify what the sample is. For this fish oil data, we classify a sample into two categories:

1. Species
2. Part

Using machine learning techniques it may be possible to speed the process up. Machine learning algorithms identify patterns in the data. These patterns can be used to classify the sample efficiently. An interpretable and accurate model has the potential to be deployed in a factory setting. It would eliminate the need for manual work. Additionally, an algorithm that can classify unresolved peaks would have an impact on the chemistry field. This increases the speed at which GC is performed, increasing the volumetric efficiency of the production line [?].

2.2 Classification Algorithms

The classification task is to identify the class label for an instance from the dataset. We perform two classification tasks and measure their performance:

1. Species: Identify the species of fish (i.e. Bluecod).
2. Part: Identify the part of the fish (i.e. Head).

A supervised learning method creates a model from a labelled dataset - the train set. We measure the ability of the model to generalize on unseen data - the test set. We give the model a tissue sample from a fish. Based on what it has learnt, it predicts the species and part for that fish. The existing analytical chemistry techniques for performing this task are laborious and time-consuming. We desire a model that is interpretable and has high predictive accuracy. We can then use the model to automate this task. Thus, it has real-world applications for quality assurance in a factory setting.

Support Vector Machines Cortes and Vapnik proposed the Support Vector Machine (SVM) [?]. This model creates a hyperplane that can draw distinct class boundaries between classes. We call these class boundaries the support vectors. We are performing multi-class classification, so it used a one-vs-all approach [?]. This creates a divide between one class and the rest, then repeats for the other classes.

Model The sklearn library provides several SVM models for classification. The default model uses the RBF kernel. Other models use different kernels and parameters [?]. Some models remove trainable parameters. Instead, the user can set the number of support vectors [?].

Kernel The model requires a kernel function. This determines the shape of the support vectors in the hyperplane. Different kernels capture data of varying complexities. The original hyperplane algorithm used a linear kernel [?]. Later, non-linear kernels we introduced. These employ the kernel trick [?]. Figure ?? shows support vectors for each kernel on a 2D plane. This provides an intuition for each kernel.

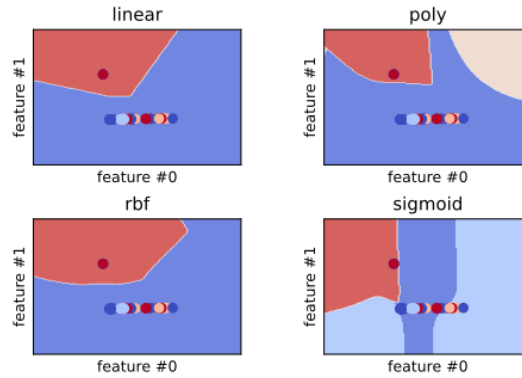


Fig. 2: SVM kernels shapes are shown. Specifically, linear, polynomial, radial basis function (rbf) and sigmoidal kernel are shown.

Hyperplane Coefficients The l1 regularization term leads to sparse models. So, they include fewer features - making them easier to interpret. Eq ?? defines the total hyperplane as

$$\beta_t = \minmax(\sum_{c \in C} |\beta_c|) \quad (1)$$

where there is the number of classes ($c \in C$) sets of hyperplane coefficients. β_t coefficient as the sum of hyperplane coefficients magnitude for each class β_c . We normalize the coefficients with a min-max feature scaling. The total hyperplane for both datasets is given in Figure ??.

2.3 Feature Selection

Feature selection reduces the complexity of the problem space. This helps counteract the curse of dimensionality [?]. Reducing the complexity improves computational efficiency, increases interpretability, and can improve performance. More interpretable models are easier for humans to understand. This means we can verify their efficiency using domain expertise in biochemistry. This is an important factor for real-world applications in a factory setting.

2.4 Visualisation

Two heuristics are optimized for when selecting a suitable model:

1. Interpretability
2. Accuracy

Interpretability is important for verification in a safety-critical environment. We intend to employ the chosen model in a factory setting. Accuracy is preferable, but not at the expense of interpretability. The efficacy of the model must be explainable through domain knowledge. Or else it is difficult to ensure reliability. The focus on interpretability ensures the model can be used in the real world.

Model interpretability is explored through visualisation. We aim to uncover learnt patterns that can be verified with domain knowledge. The desired algorithm should strike a balance between predictive performance and semantically meaningful features.

What constitutes semantic meaning varies from one domain to another. It is easy to build intuition for semantic meaning in computer vision and natural language processes, they correspond to recognisable images and structured text. In the domain of food science, our meaning is derived from performance on the classification task(s) and similarity to underlying chemical compounds.

3 Data processing

- Why the raw data is not applicable to existing classification algorithms?
- Extracting datasets that are ready for classification algorithms:
 - Sum up the intensity.
 - Aligning missing packets.
- Overview of extracted data.

4 Classification Algorithms

We measure the predictive ability of classifiers on both the fish species and part dataset. We are looking for the model with the highest accuracy. As a result, we start broadly by exploring a variety of models from the different families of AI, then we narrow and refine the search.

Specifically, we examine:

1. Ensemble
2. SVM Model
3. SVM Kernel

For each of the following experiments, the same experiment setup is used. We use stratified cross-validation ($k = 10$) to measure the classification accuracy. Each method has its performance recorded on the same cross-folds. Then we average over 30 independent runs. This experimental setup evaluates performance on both the fish species and part datasets.

4.1 Ensemble

This is a broad search for an effective classification model. We explore a classifier from each family of AI. The model with the highest classification accuracy on both datasets is then selected, and explored by later sections in further detail.

We examine 5 classification models:

1. K-Nearest Neighbors [?]
2. Random Forest [?]
3. Naive Bayes [?]
4. Decision Tree [?]
5. Support Vector Machine [?]

Table ?? shows for random forest, decision tree and support vector machine have perfect training accuracy. The decision tree and random forest overfit the training data. Only the SVM achieves similar performance on the test data. The SVM classifier outperforms the other classifiers. It does so for the test set for both the species and part datasets.

Dataset	Method	AvgTrain \pm Std	AveTest \pm Std
Species	KNN	83.57 \pm 1.80	74.88 \pm 12.54
	RF	1.00 \pm 0.00	85.65 \pm 10.76
	DT	1.00 \pm 0.00	76.98 \pm 13.12
	NB	79.54 \pm 1.60	75.27 \pm 4.35
	SVM	1.00 \pm 0.00	98.33 \pm 5.00
Part	KNN	68.95 \pm 3.49	43.61 \pm 13.48
	RF	1.00 \pm 0.00	72.60 \pm 16.15
	DT	1.00 \pm 0.00	60.14 \pm 14.57
	NB	65.54 \pm 2.69	48.61 \pm 12.19
	SVM	1.00 \pm 0.00	87.14 \pm 8.52

Table 1: Accuracy for different classification techniques. Accuracy is given as the stratified k-fold cross validation over 30 independent runs. We compare K-nearest neighbours (KNN), random forest (RF), decision tree (DT), naive bayes (NB) and support vector machines (SVM).

4.2 SVM Model

The classification results showed that SVM was the most effective classifier. Now, we explore the variations in models for the SVM classifier. We use the same cross-validation setup as before.

We examine 3 SVM models [?]:

1. Suport Vector Classification [?]
2. Nu-Support Vector Classification [?]
3. Linear Support Vector Classification

Dataset	Method	AvgTrain \pm Std	AveTest \pm Std
Species	svc	88.96 \pm 1.40	80.00 \pm 12.33
	nusvc	88.30 \pm 1.17	81.73 \pm 12.75
	lsvc	1.00 \pm 0.00	98.33 \pm 5.00
Part	svc	73.25 \pm 3.54	49.03 \pm 12.14
	nusvc	90.31 \pm 1.97	62.36 \pm 15.18
	lsvc	1.00 \pm 0.00	87.16 \pm 8.56

Table 2: Accuracy for different SVM models. Accuracy is given as the stratified k-fold cross validation over 30 independent runs. We compare Support-Vector Classification (SVC), Nu-Support Vector Classification (Nu-SVC) and Linear Support-Vector Classification (LSVC).

Table ?? shows for fish species, SVC and Nu-SVC models have similar performance on both train and test. The Nu-SVC outperforms the SVC for both train and test for the part dataset. Yet, the linear SVC outperforms both models. It

achieves perfect training accuracy for both datasets. For the test, near-perfect (98.33%) on species, and reasonable performance (87.16%) on the part.

4.3 SVM Kernel

Now we know that SVM is the most effective classifier, and the LSVC is the most effective model. To provide an exhaustive search, we explore all possible kernels. We use the same cross-validation setup as before.

We examine 4 SVM kernels [?]:

1. Polynomial
2. Radial Basis Function (rbf)
3. Sigmoid
4. Linear [?]

Dataset	Method	AvgTrain \pm Std	AveTest \pm Std
Species	poly	76.83 \pm 1.18	71.37 \pm 15.86
	rbf	88.96 \pm 1.40	80.00 \pm 12.33
	sigmoid	33.19 \pm 2.36	30.18 \pm 6.50
	linear	1.00 \pm 0.00	97.50 \pm 5.34
Part	poly	70.63 \pm 2.27	53.89 \pm 6.94
	rbf	73.25 \pm 3.54	49.03 \pm 12.14
	sigmoid	37.47 \pm 1.78	33.47 \pm 8.59
	linear	1.00 \pm 0.00	87.36 \pm 10.77

Table 3: Accuracy for different SVM kernels. Accuracy is given as the stratified k-fold cross validation over 30 independent runs. We compare polynomial (poly), radial basis function (rbf), sigmoidal (sigmoid) and linear.

Table ?? shows the sigmoid kernel performs very poorly on training and test for both datasets. The polynomial and RBF kernel achieve comparable performance for both datasets. The linear kernel outperforms all other kernels for both datasets. It has near-perfect (97.50%) test accuracy on fish species. And reasonable performance (87.36%) on the fish part.

4.4 Discussion

We evaluated an ensemble of classification techniques. Naive Bayes performed poorly. This is likely due to the assumption of conditional independence between features. KNN also performed poorly. This is likely due to the high dimensionality of the data. Points drawn from high dimensional spaces tend to never be close together. SVM provided the best results. This model can identify fish species from gas chromatography data with near-perfect accuracy. This prompted further investigation into this technique.

Within support vector machines, the Linear SVC model showed the best performance. Naturally, within SVM kernels, the linear kernel also showed the best performance. There is high predictive performance on the linear kernel. This suggests an underlying pattern that is linearly separable in a hyperplane. Non-linear kernels - polynomial, RBF or sigmoidal - produced diminishing returns. These kernels try to fit complex patterns that are not present in the data.

Performance for all models was better for the fish species than the part. This suggests tissue samples for different species may have distinct chemical compositions. Yet, different fish parts may have fewer underlying structural differences. For GC data the intra-class variation between species provides a larger signal than part variation. For example, we expect there to be more difference between a tarakihi and a bluecod, than there is a similarity between two livers from each species.

5 Feature Selection

- Why feature selection on this data?
- Brief the main ideas of the feature selection algorithms that were used.
- Compare the performance of selected features and using all features.
- (Optional): analyse the selected features.

For each method, we measure classification accuracy with an SVM model [?]. It has linear kernel, l1 regularization [?] and 10,000 maximum iterations. We examine 4 feature selection methods [?]:

1. Chi² [?]
2. Minimum Redundancy Maximum Relevance [?]
3. ReliefF [?]
4. Particle Swarm Optimization [?,?]

We first provide a detailed accuracy comparison for a set feature number ($k = 500$). Then we explore the accuracy for the general case (any k).

5.1 Classification Accuracy $k = 500$

We measure the classification accuracy at $k = 500$ for each method. To allow comparison with PSO, we take the top k features suggested by the algorithm and compare this to the others.

Table ?? shows for the training set, MRMR, ReliefF and PSO have comparable accuracy for both datasets. The Chi² method does not, instead it performs very poorly. For the test set, ReliefF performs best for species, PSO performs best for the part.

Dataset	Method	AvgTrain \pm Std	AveTest \pm Std
Species	Chi ²	95.17 \pm 3.52	81.85 \pm 9.65
	MRMR	99.79 \pm 0.41	95.09 \pm 6.90
	ReliefF	99.71 \pm 0.44	95.12 \pm 6.26
	PSO	99.71 \pm 4.30	93.30 \pm 8.16
Part	Chi ²	96.32 \pm 0.88	64.86 \pm 19.01
	MRMR	97.44 \pm 0.97	78.79 \pm 13.21
	ReliefF	97.82 \pm 1.04	80.28 \pm 5.58
	PSO	97.62 \pm 0.91	82.36 \pm 10.72

Table 4: Accuracy for different feature selection methods. Accuracy is given as the stratified k-fold cross validation over 30 independent runs. We compare chi² (chi), maximum relevance - minimum redundancy (MRMR), reliefF, particle swarm optimisation (PSO).

5.2 Classification Accuracy (all k)

We measure classification accuracy as a function of feature number. We compared this for several FS methods. Due to limitations, PSO optimizes feature number k automatically. So, to compare its performance, we plot the results of 30 independent runs.

Specifically, we compare the following methods:

- reliefF [?]
- MRMR [?]
- PSO [?,?]
- χ^2 [?]

Figure ?? shows accuracy for fish species. We show accuracy on the training set for each feature selection method. At $k = 1050$ all feature selection methods achieve 100% accuracy on the training set. The SVM fits the training data for each method using a fraction of the full feature set. Figure ?? shows accuracy for fish species. We show test set accuracy for each feature selection method. The accuracy reaches a plateau (96% accuracy) at around $k = 1050$ features for all methods. The test performance is less than the train performance, yet the test accuracy is still very high. This suggests the model can generalize well on unseen data for the fish species.

Figure ?? shows accuracy for part dataset. We show train accuracy for each feature selection method. All feature selection methods struggle to fit the training set for the fish part. Even with the full feature set, a perfect train accuracy is never reached. Figure ?? shows accuracy for part dataset. We show the test accuracy for each feature selection method. The classification accuracy fluctuates for all feature selection methods. At around $k = 1050$ features, it begins to decrease. The training accuracy improves, as the test does not from this point onwards. The SVM is overfitting to noise (redundant features) in the training set.

5.3 Disucssion

Feature selection methods helped reduce dimensionality. We evaluated performance with an SVM classifier. In which, ReliefF and PSO were best for fish species and part, respectively. ReliefF can identify conditional dependencies between features when providing feature rankings. ReliefF algorithms are robust and noise-tolerant, which explains their superior performance. PSO provides a combination of global and local searches. A search through a near-infinite combinatorial space of possible feature subsets. This stochastic method is computationally expensive but can offer effective solutions.

For both general and specific cases, and across all methods, the fish species have lower variance than the fish part in classification accuracy. The classification results (§ ??) support this, they also show higher test accuracy for fish species, than fish part. They suggest different fish parts may have fewer underlying structural differences.

For the general case and both datasets, a lot of interesting behaviour happens at $k = 1050$. The fish species reach a plateau, but the fish part accuracy begins to decrease. An accuracy comparable or better than the full dataset is possible with 21.8% of its features.

For the fish species dataset, we see high accuracy with very few features. ReliefF and MRMR can achieve above 90% classification accuracy with $k = 50$. Chi² is not able to mirror this performance. This shows ReliefF and MRMR are very effective feature selection methods for this task.

The PSO may not have a hyperparameter for feature number k . Instead, it automates the selection of this parameter. Yet, it achieves comparable results to other state-of-the-art methods. This automation may prove useful for automating the classification task for online learning. In a factory, we may want to train a model as new data arrives. PSO requires less human intervention, yet still, provides competitive performance.

MRMR and ReliefF both have high accuracy with very few features on the fish species dataset. This suggests that few features are required to construct a reasonable representation of a fish tissue sample. This is a good indication that the fish species dataset contains less noise. This also warrants further investigation into which features are considered important for low k values. This motivates the following section on visualisation.

5.4 Visualisation

In this section, we explore the interpretability of the feature selection methods. Firstly, We the features and their rankings for low k values. Then investigate if there is a superset of important features - these would be universally recognized by all methods - that is, the overlap or common features. Thirdly, the classification accuracy of the overlap is compared to its parts to see if the superset even exists. Finally, we compare the most effective methods (reliefF and MRMR) to the most effective classifier (SVM).

Specifically, we examine:

1. Feature Rankings ($k = 50$)
2. Overlap: MRMR \cap ReliefF
3. Accuracy: Overlap, MRMR, ReliefF
4. Overlap: FS Method \cap Hyperplane Coefficients

Experimental Setup When evaluating classification accuracy we use a stratified cross-validation ($k = 10$). We give the average balanced accuracy over those k folds. When experiments are examining feature selections/rankings, we fit the entire dataset. We repeat each experiment for both fish species and part datasets.

The number of features k is the independent variable for many of these methods. We exclude PSO [?,?] from the analysis. Since k is not a hyper-parameter for this algorithm. We analyse feature importance for three feature selection methods.

Specifically, we compare the following methods:

- reliefF [?]
- MRMR [?]
- χ^2 [?]

Feature Rankings ($k=50$) We compare the feature rankings for the top 50 ($k = 50$) features for each FS method. We normalize scores ($x \in [0, 1]$) to allow comparison between methods. We see if each FS method selects similar features for a low number of features ($k = 50$). If true, this suggests that there are important features. All techniques would recognize these. Previous work [?,?] has shown worse performance for χ^2 . The visualisation also helps gain an intuition for why this could be the case.

These results relate to the fish species dataset. Figure ?? shows reliefF favors higher feature indexes ($i \in [4000, 4800]$). The feature indexes tend to bunch towards this range. Figure ?? shows MRMR samples from the entire set of feature indexes ($i \in [0, 4800]$). This method samples feature indexes with a more uniform distribution than reliefF. Figure ?? shows χ^2 selects chains of features adjacent to each other. These chains hint the method is not able to make a good selection for low values of k . We would expect poor performance for the χ^2 method for low values of k - the feature selection (§ ??) results confirm that claim.

These results relate to the fish part dataset. Figure ?? shows reliefF favours feature indexes bunching towards the higher range. Figure ?? shows MRMR gives a more uniform distribution of features. Figure ?? shows χ^2 tends to cluster of features around index $i = 3000$. The features cluster more for the fish part than the species.

Overlap: MRMR \cap ReliefF We give overlap - the intersection of features selected by MRMR and ReliefF for a given k - as MRMR \cap ReliefF. The overlap measures the similarities between the feature subsets chosen by each method.

If the overlap value is high, this suggests there are a distinct set of important features. Here, we would expect both methods to recognize these early, or, there may be low overlap. A low overlap would suggest that there are many important features, each method can find unique important features on its own. We exclude χ^2 from this analysis because there was little overlap for lower k values with other methods.

Figure ?? shows this overlap. Intuitively, the limit of the overlap, when k approaches all features k_{max} , is the number of features k_{max} . Due to this limit, we are more interested in the left-hand side of the graph, the low k values. For low k values the overlap is very small. There are a few features important enough to be selected early by both methods, this result is consistent for both datasets.

Accuracy: Overlap, MRMR, ReliefF We compare the classification accuracy of the MRMR, ReliefF and their overlap ($\text{MRMR} \cap \text{ReliefF}$) for a given k . χ^2 had poor performance in the classification task [?]. Also, this method has little to no overlap. So, we no longer consider χ^2 in our analysis. Now, we determine if the first k features selected by each FS method are the most important. Say classification accuracy for the overlap exceeds that of its parts (MRMR, ReliefF). This would mean the most important features are the first k features selected by each method. Or, the overlap may share the same classification accuracy as its parts. This would suggest many features are all important. Feature subsets with the same accuracy are as important as each other.

These results relate to fish species. Figure ?? suggests all subsets of features - MRMR, reliefF and overlap - can fit the data for any k . Figure ?? shows different feature subsets achieve similar test accuracy - so, there are many important features. This suggests that each feature subset are equally important. The feature selection results (§ ??) support this, they showed high accuracy with few features selected ($k = 50$). This result supports that conclusion.

These results relate to the fish part. Figure ?? suggests all subsets of features - MRMR, reliefF and overlap - can fit the data for any k . Figure ?? fluctuates wildly. This is consistent with previous findings [?,?]. This supports that classifying fish parts is more difficult than fish species. Due to the noise in the results, it is difficult to distinguish a clear winner for the part dataset.

Overlap: FS Method \cap Hyperplane Coefficients We compare the hyperplane coefficients β_i to the FS methods method. The SVM classifier has high predictive capabilities on both datasets [?,?]. Thus, we expect the FS method most like SVM is of interest.

Figure ?? shows reliefF with the most consistent overlap with SVM coefficients, this is true for both datasets. The χ^2 method shows wild fluctuations, it is similar for high k values for the fish species, but then similar for low k for the fish part. Apart from a small region of low k values, MRMR shares little overlap for the fish species. For the fish part, MRMR shares little overlap with SVM coefficients, when compared to others. In terms of features selected, reliefF is the most similar to the Linear SVM.

Discussion The visualisation has shown us why χ^2 is a poor fs method. It tends to bunch features together. MRMR and reliefF are far better. This is true for the distribution of feature indexes, overlap and classification accuracy. MRMR has the most uniform distribution of feature indexes. Yet, reliefF selects features most like SVM coefficients.

When we select a few features, there was little overlap between MRMR and reliefF. There were no clear k best features consistent for both techniques. The overlap of features selected did not perform better than its parts (MRMR and reliefF). This shows there are lots of important features.

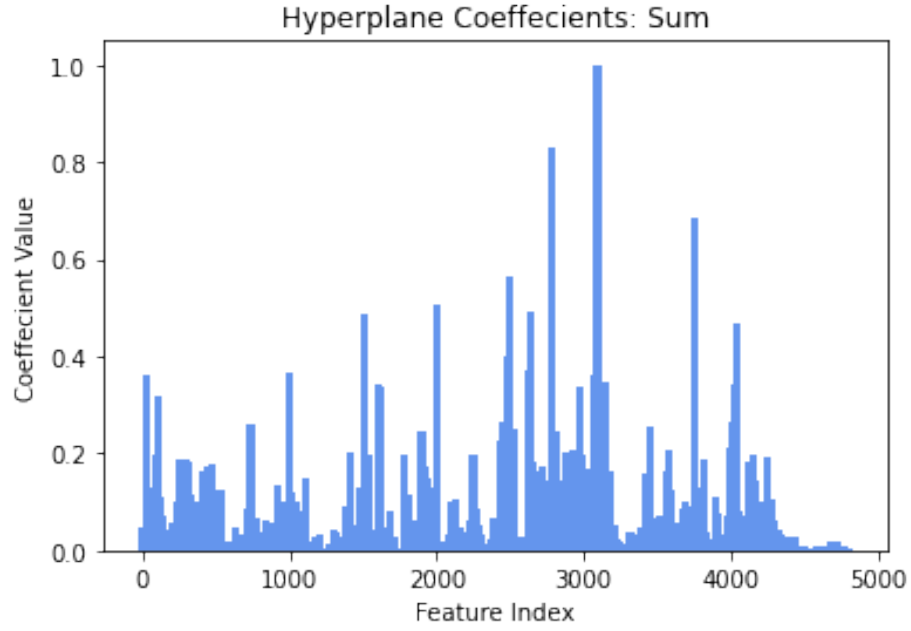
These FS methods were not able to converge on a small superset of important features. That is because there are lots of important features for both datasets. Each feature represents an intensity at a particular timestamp in the gas chromatograph. These results show we can use many chemical combinations to classify fish oil data. It may be easier to match some reference chemicals than others. There is a range of possibilities for combinations of reference chemicals. This research shows we may reduce the cost of this analysis in future.

6 Conclusions and Future Work

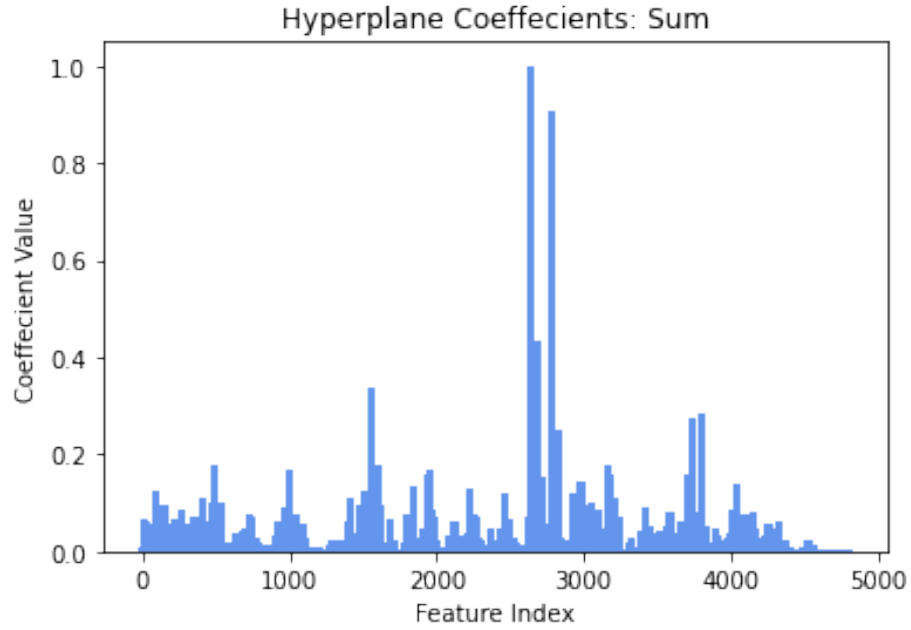
The analysis of fish oil data in this paper has focussed on interpretability. Not only have we found effective classification and feature selection techniques, but we have also tried to explain their performance with visualisation and analytical results. We can draw many conclusions from the analytical results and visualisations, but here we recall the most important:

1. Fish species are easier to predict than fish part - there is more intra-class variation within fish species than there is a similarity between the same part from different fish.
2. The Linear SVM classifier performs better for both classification tasks - the fish oil data is linearly separable on a hyperplane.
3. Near-perfect accuracy can be achieved with very few features for fish species - if this predictive ability exceeds human error this model has real-world applications.
4. Comparison with PSO may be difficult (due to automatic k selection), but this automation may be useful in a factory setting.
5. There are many important features and their subsets that can be used to classify fish oil - thus, there is a range of possibilities for combinations of reference chemicals.

Feature selection is not guaranteed to improve classification accuracy. Yet, it does reduce the complexity, increase interpretability and improve computational efficiency. As with the SVM using an $l1$ regularization, feature selection is also a trade-off between interpretability and performance. Arthur C. Clarke said, "[a]ny sufficiently advanced technology is indistinguishable from magic." [?]. A perfect blackbox model is equivalent to magic, and, it is difficult to have faith in magic, especially if it were to be involved in our food making. Faith isn't needed when our models are interpretable.

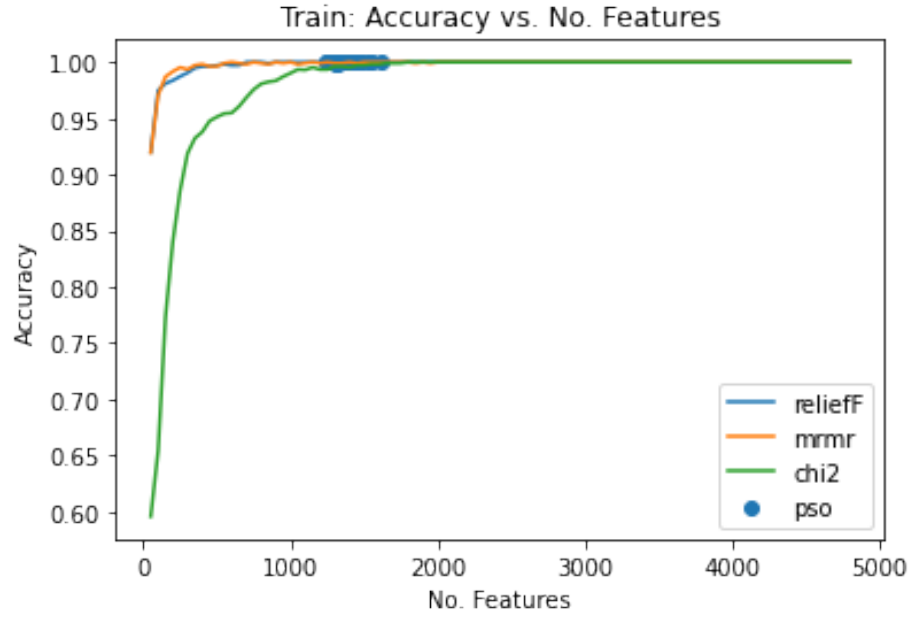


(a) Fish Species Hyperplane Coefficients

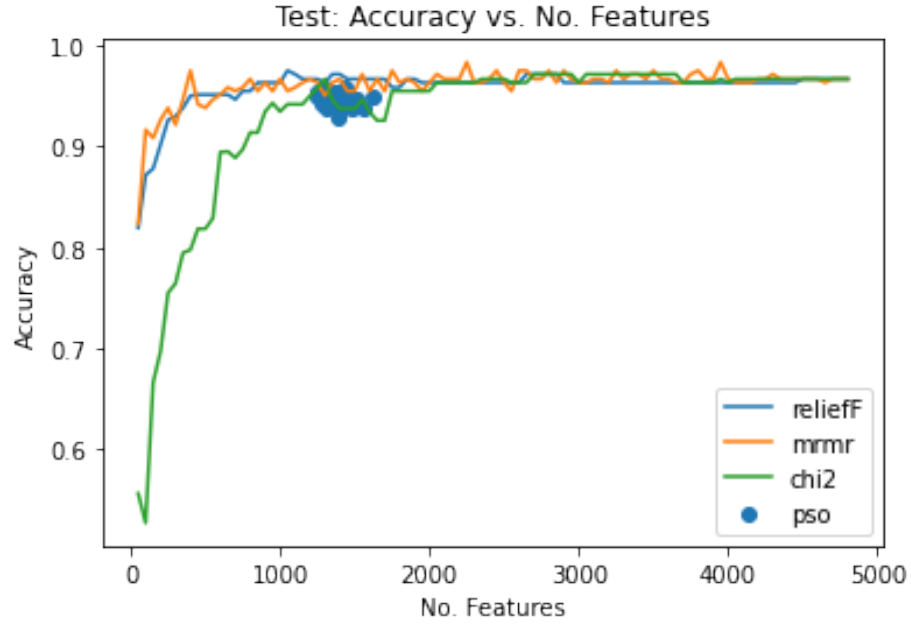


(b) Fish Part Hyperplane Coefficients

Fig. 3: Hyperplane coefficients β_t . The normalized sum of the magnitude of the coefficients for each class is given in Eq ?? . (a) Coefficients for the fish species dataset. (b) Coefficients for the fish part dataset.

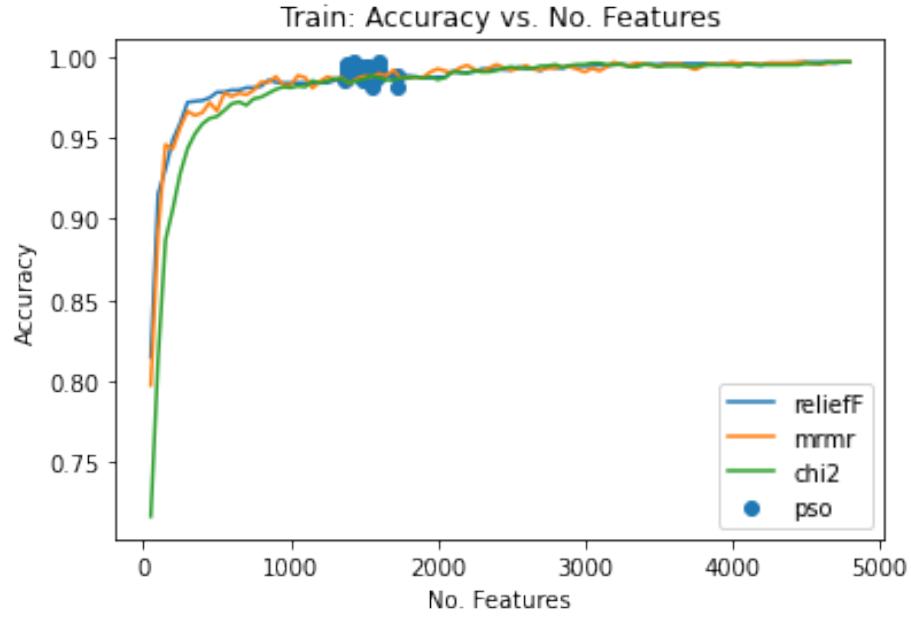


(a) Training set

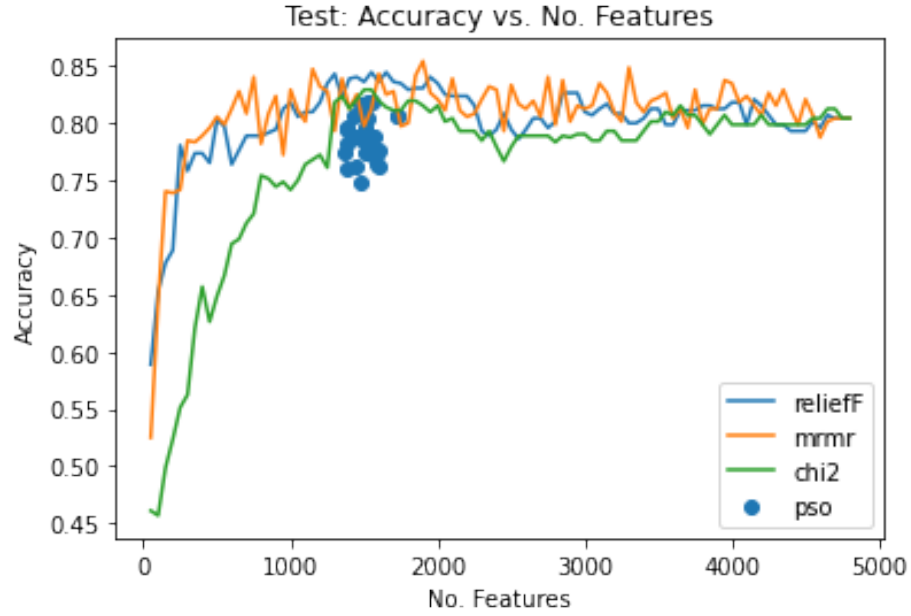


(b) Test set

Fig. 4: Fish species dataset: Classification accuracy for feature selection methods for a given k . We measure the balanced accuracy of k -fold cross-validation. We compare reliefF, maximum relevance - minimum redundancy (MRMR), χ^2 , and particle swarm optimisation (PSO). Fish part dataset. (a) Training set. (b) Test set.

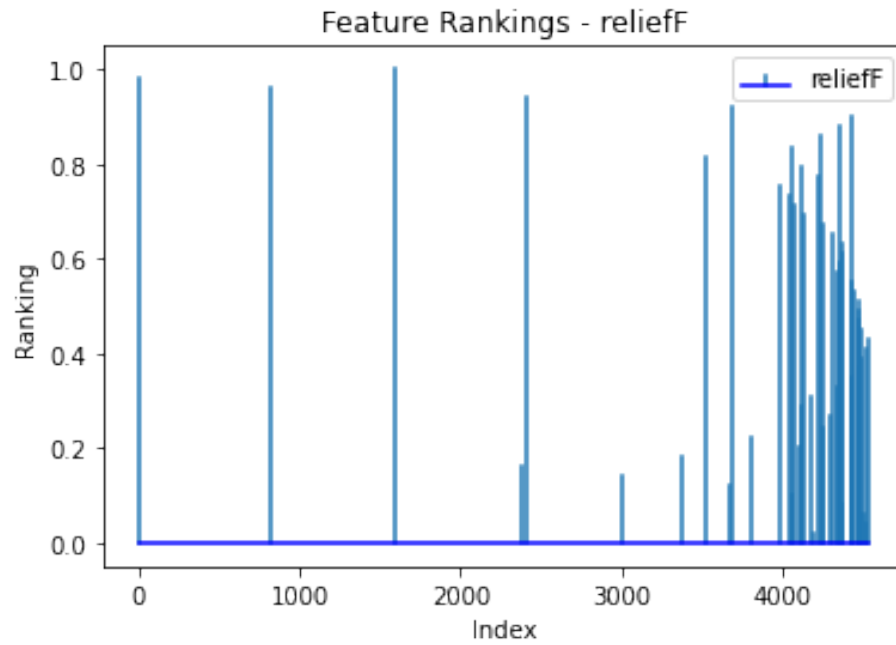


(a) Training set

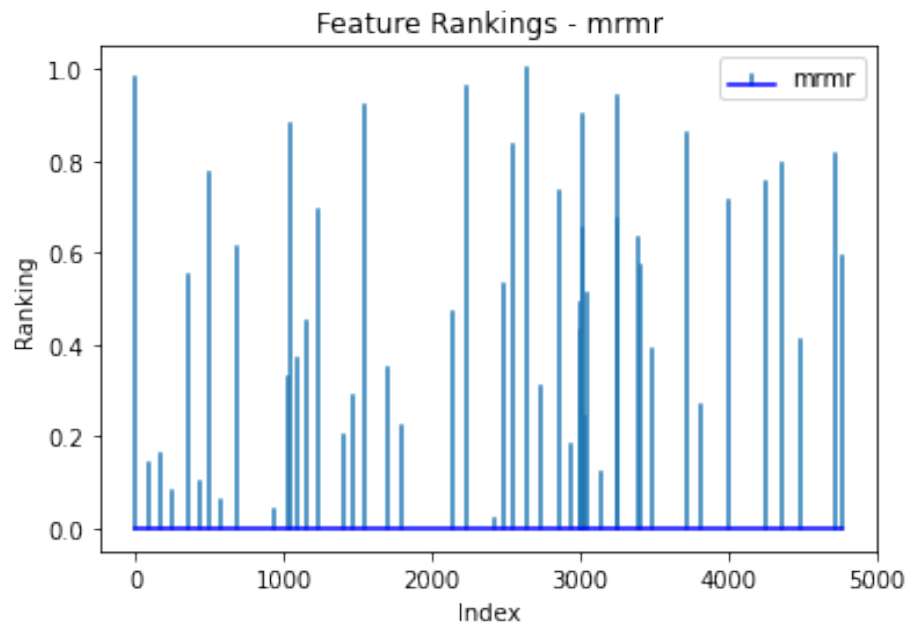


(b) Test set

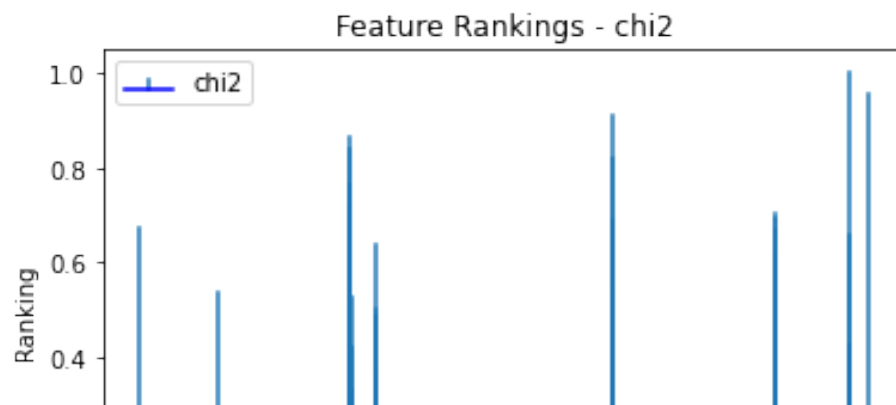
Fig. 5: Fish part dataset: classification accuracy for feature selection methods for a given k . We measure the balanced accuracy of k -fold cross-validation. We compare reliefF, maximum relevance - minimum redundancy (MRMR), χ^2 , and particle swarm optimisation (PSO). (a) Training set. (b) Test set.

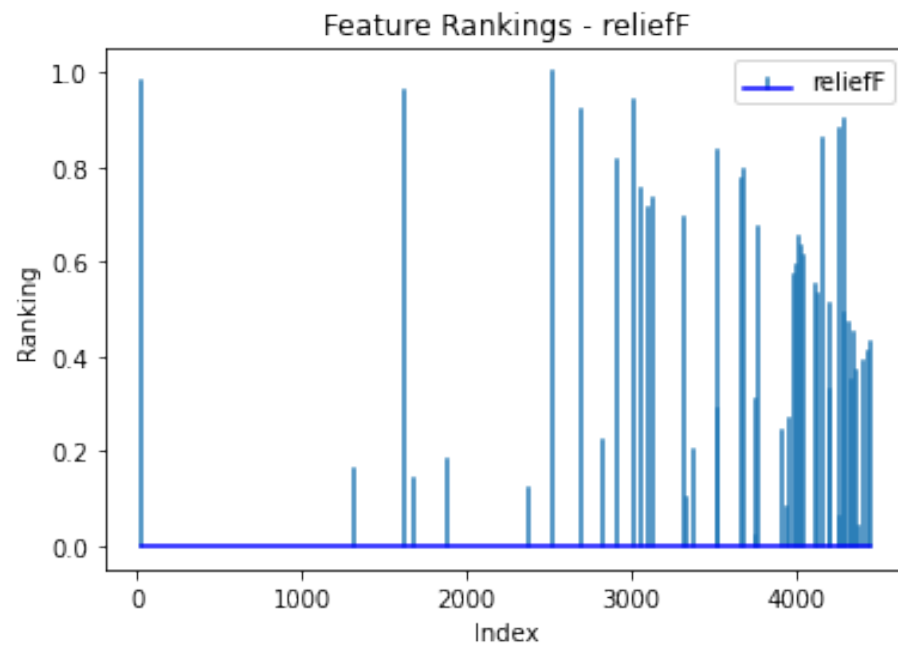


(a) reliefF

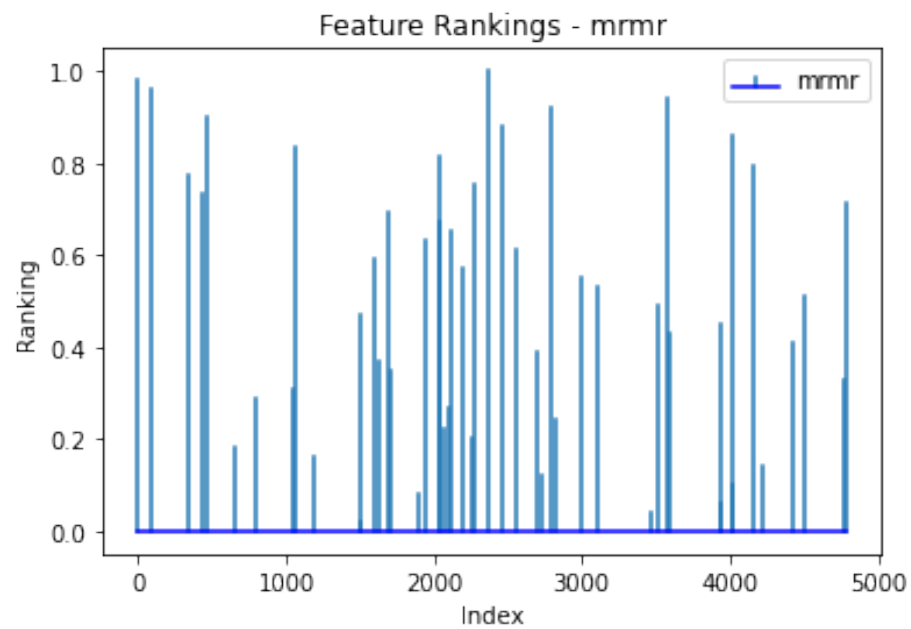


(b) Maximum Relevance — Minimum Redundancy (MRMR)

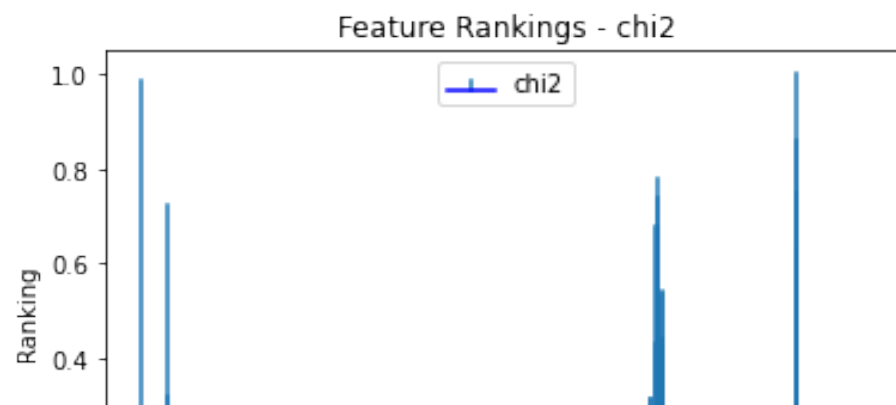


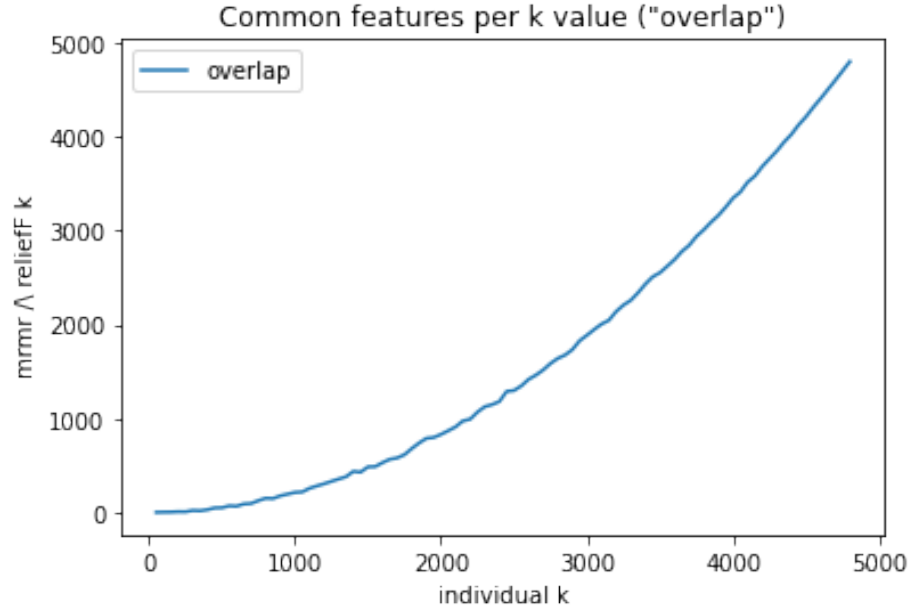


(a) reliefF

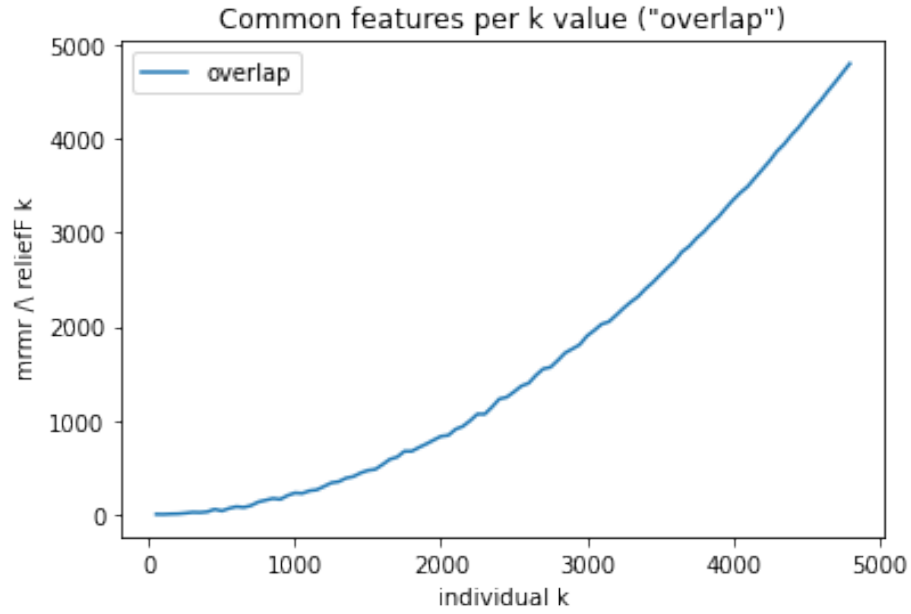


(b) Maximum Relevance — Minimum Redundancy (MRMR)



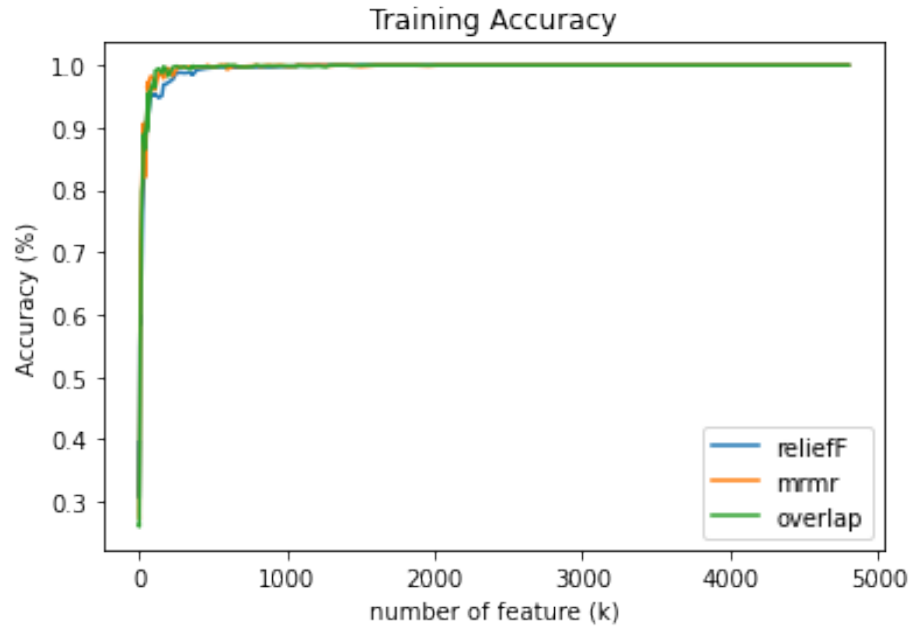


(a) Fish species dataset

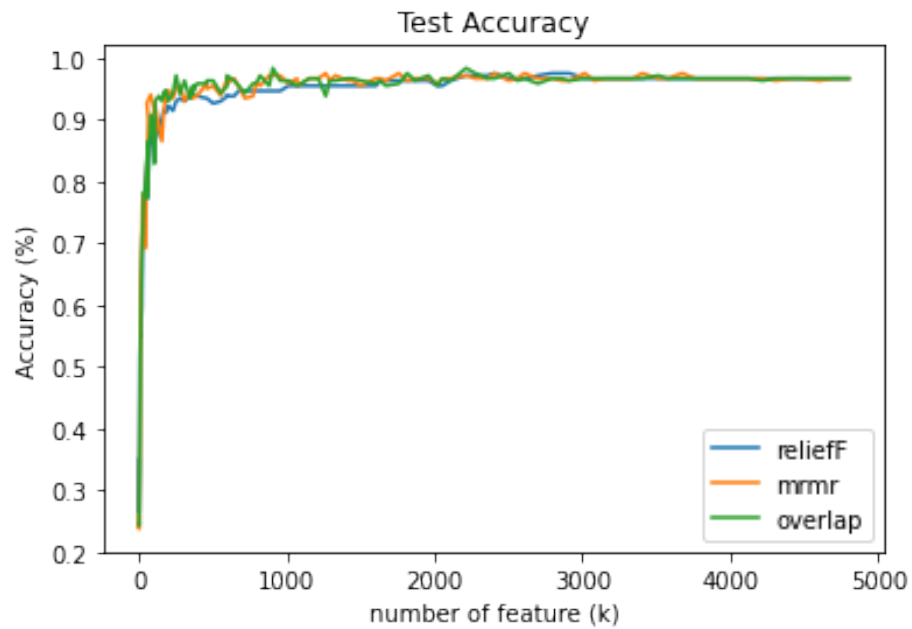


(b) Fish part dataset

Fig. 8: Overlap - common features for MRMR and ReliefF. The overlap demonstrates the number of shared features between the two feature selection methods for a give k value. (a) Fish species. (b) Fish part.

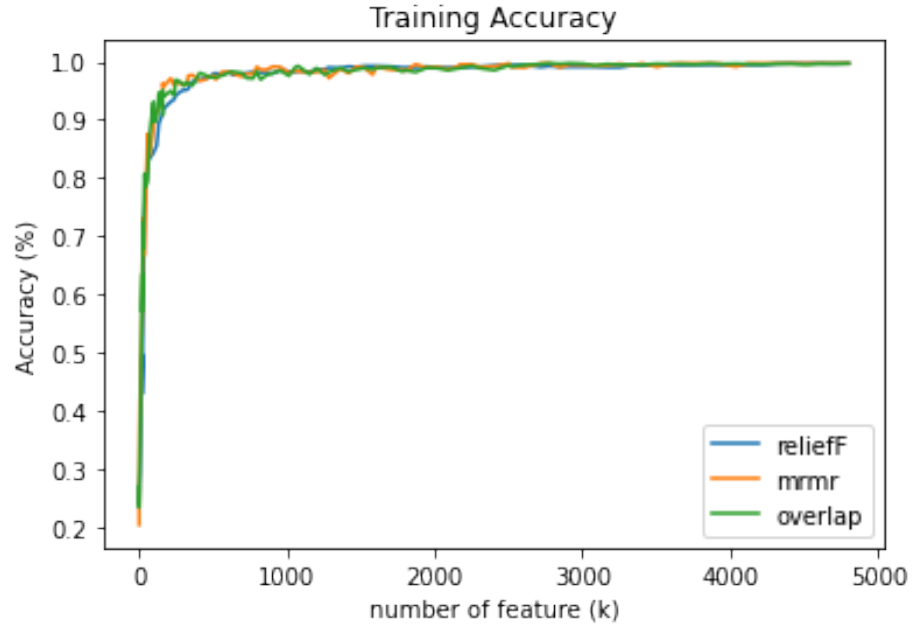


(a) Training set

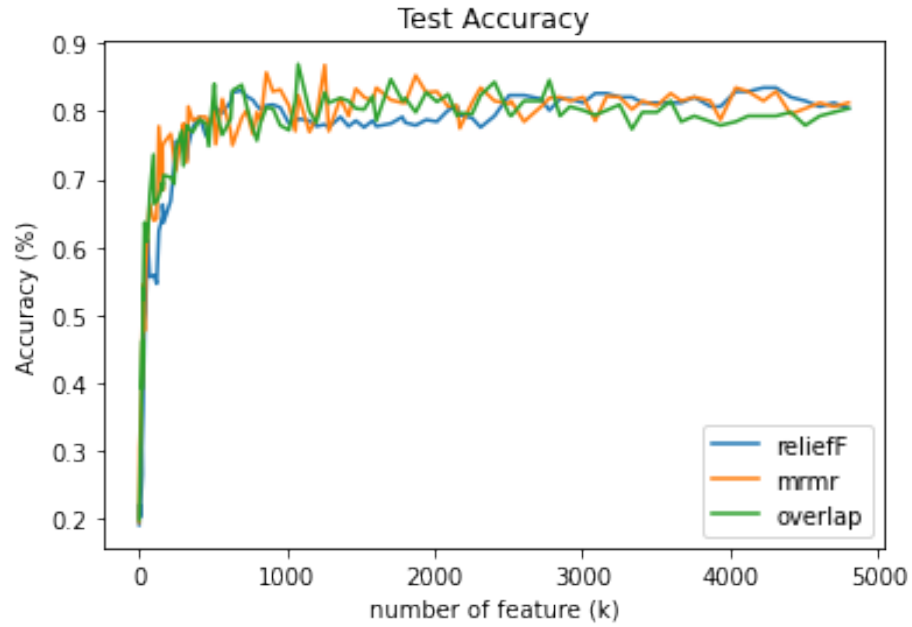


(b) Test set

Fig. 9: Fish species: classification accuracy of the overlap, MRMR and ReliefF for a given k . We measure the balanced accuracy using k -fold cross-validation. The classification accuracy metric shows how important the selected features are for fish part species. (a) Training set. (b) Test set.

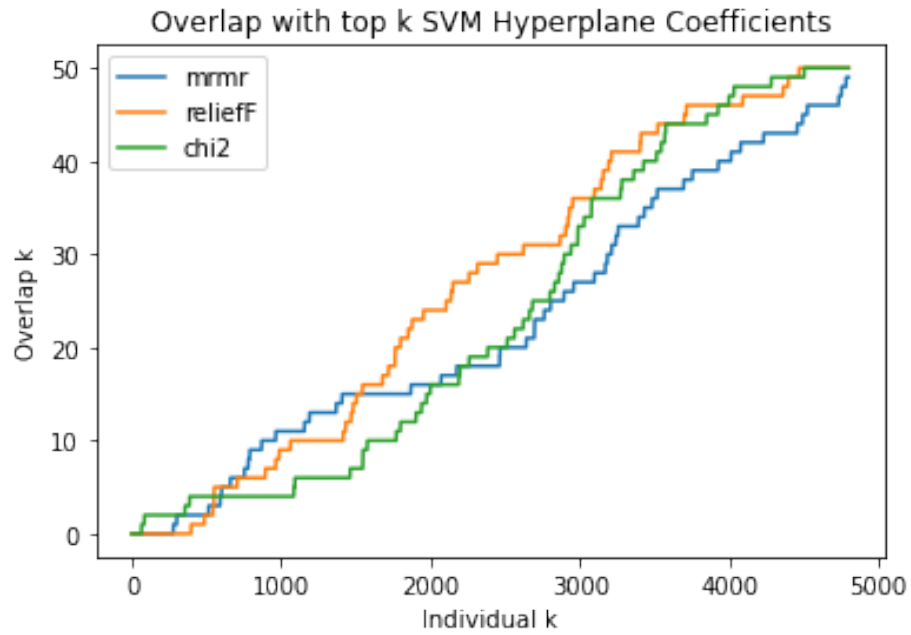


(a) Train set

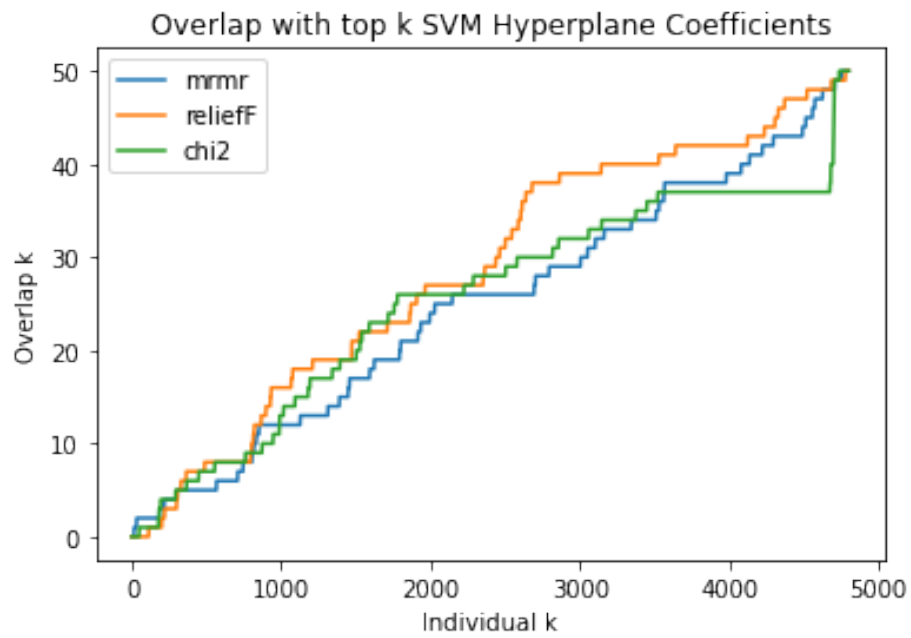


(b) Test set

Fig. 10: Fish part: classification accuracy of the overlap, MRMR and ReliefF for a given k . We measure the balanced accuracy using k -fold cross-validation. The classification accuracy metric shows how important the selected features are for the fish part dataset. (a) Training set. (b) Test set.



(a) Fish species dataset



(b) Fish part dataset

Fig. 11: Overlap - common features between SVM coefficients and feature selection methods for a given k value. We examine the maximum relevance - minimum redundancy (MRMR), relieF, and χ^2 methods. The overlap measure shows how similar each method is to the SVM classification method. (a) Fish species. (b) Fish part.