

# The impact of feature selection on one and two-class classification performance for plant MicroRNAs

Waleed Khalifa, Malik Yousef, Müşerref Duygu Saçar Demirci, Jens Allmer

MicroRNAs (miRNAs) are short nucleotide sequences that form a typical hairpin structure which is recognized by a complex enzyme machinery. It ultimately leads to the incorporation of 18 – 24 nt long mature miRNAs into RISC where they act as recognition keys to aid in regulation of target mRNAs. It is involved to determine miRNAs experimentally and, therefore, machine learning is used to complement such endeavors. The success of machine learning mostly depends on proper input data and appropriate features for parameterization of the data. Although, in general, two-class classification (TCC) is used in the field; because negative examples are hard to come by, one-class classification (OCC) has been tried for pre-miRNA detection. Since both positive and negative examples are currently somewhat limited, feature selection can prove to be vital for furthering the field of pre-miRNA detection. In this study, we compare the performance of OCC and TCC using eight feature selection methods and seven different plant species providing positive pre-miRNA examples. Feature selection was very successful for OCC where the best feature selection method achieved an average accuracy of 95.6%, thereby being ~29% better than the worst method which achieved 66.9% accuracy. While the performance is comparable to TCC, which performs up to 3% better than OCC, TCC is much less affected by feature selection and its largest performance gap is ~13% which only occurs for two of the feature selection methodologies. We conclude that feature selection is crucially important for OCC and that it can perform *on par* with TCC given the proper set of features.

# The impact of feature selection on one- and two-class classification performance for plant microRNAs

Waleed Khalifa<sup>1,2</sup>, Malik Yousef<sup>1,2</sup>, Müşerref Duygu Saçar Demirci<sup>3</sup>, and Jens Allmer<sup>3,4</sup>

<sup>1</sup> Computer Science, The College of Sakhnin, Sakhnin, 30810, Israel

<sup>2</sup> The Institute of Applied Research- the Galilee Society, P.O. Box 437 Shefa Amr, 20200, Israel

<sup>3</sup> Molecular Biology and Genetics, Izmir Institute of Technology, 35430 Urla, Izmir, Turkey

<sup>4</sup> Bionia Incorporated, IZTEKGEB A8, 35430 Urla, Izmir, Turkey

Contact Emails: malik.yousef@gmail.com, duygu.sacar@gmail.com, khwalid@hotmail.com, jens@allmer.de

## Abstract

MicroRNAs (miRNAs) are short nucleotide sequences that form a typical hairpin structure which is recognized by a complex enzyme machinery. It ultimately leads to the incorporation of 18 – 24 nt long mature miRNAs into RISC where they act as recognition keys to aid in regulation of target mRNAs. It is involved to determine miRNAs experimentally and, therefore, machine learning is used to complement such endeavors. The success of machine learning mostly depends on proper input data and appropriate features for parameterization of the data. Although, in general, two-class classification (TCC) is used in the field; because negative examples are hard to come by, one-class classification (OCC) has been tried for pre-miRNA detection. Since both positive and negative examples are currently somewhat limited, feature selection can prove to be vital for furthering the field of pre-miRNA detection. In this study, we compare the performance of OCC and TCC using eight feature selection methods and seven different plant species providing positive pre-miRNA examples. Feature selection was very successful for OCC where the best feature selection method achieved an average accuracy of 95.6%, thereby being ~29% better than the worst method which achieved 66.9% accuracy. While the performance is comparable to TCC, which performs up to 3% better than OCC, TCC is much less affected by feature selection and its largest performance gap is ~13% which only occurs for two of the feature selection methodologies. We conclude that feature selection is crucially important for OCC and that it can perform *on par* with TCC given the proper set of features.

## Introduction

Gene regulation is of prime importance in all living organisms and there are multiple levels at which gene expression can be modulated. MicroRNAs (miRNAs) play a role in post-transcriptional gene regulation [1] and, among other functions, fine-tune the amount of translated protein product [2]. Mature miRNAs are short nucleotide sequences discovered about two decades ago [3]. From databases which host miRNAs like miRBase [4] it can be gleaned that miRNAs exist in a wide range of organisms ranging from viruses [5] to plants [6]. It has also been proposed that the plant miRNA system may have evolved independently [7] and some organisms like yeasts also display differences to the canonical pathway [8]. Any regulatory element itself may be miss-regulated and miRNAs are no exception, and therefore, have been implicated in, for example, human diseases [9], [10] and in plant stress response [11]. While miRNAs may lead to inter-kingdom communication in special cases [12], it is not likely that there is extensive communication among eukaryotes [13]. Experimentally detected and/or validated miRNAs are available in databases such as miRBase [14] and miRTarBase [15]. MicroRNAs' effect can only be established when it is co-expressed with its targets [2], which complicates experimental analysis since only a fraction of the genome is expressed at a given time, in a tissue, or in response to stress conditions; and testing all conditions experimentally is elusive. Additionally, such an analysis needs to be performed on transcript and protein level, concurrently, over multiple time points to establish a causative relationship. Therefore, it seems impossible to experimentally detect all possible miRNAs of any higher eukaryotic organism [16]–[19]. Moreover, it has become clear that even among the experimentally validated miRNAs in miRBase and miRTarBase, there may be dubious

examples [20]. Therefore, carefully designed computational experiments are required to complement experimental approaches for miRNA detection.

Many computational approaches to miRNA detection have been proposed and most of them derive numerical features [21] to describe a pre-miRNA and then use machine learning to establish a model for miRNA identification [22]–[26]. Of these approaches, most, with few exceptions [16], [23], [27], employ two class classification; the latter has been compared previously [24], [28]. Classification in machine learning depends on positive examples for training the classifier in case of one-class classification (OCC) and additionally on negative data in case of two-class classification (TCC). The negative data, however, proves to be difficult to establish (if not impossible), so that all negative datasets currently in use are based on arbitrary selection of examples from parts of a genome deemed not miRNA genic or from randomly generated sequences. While both approaches are questionable, they present the only alternative to using OCC and in the absence of proper benchmark data need to be used for TCC [29]. Since OCC only needs examples for the target class (here positives), it can obliterate the need to define artificial negative examples [30], [31] and can be used to differentiate between target and unknown class. We have recently analyzed the use of OCC for miRNA detection in plants and found that it was competitive in comparison to TCC although the analysis was unduly biased towards TCC [6]. Our previous study also showed that among the hundreds of features proposed for miRNA parameterization [21] some are more discriminative than others. Since feature selection is NP-hard [32], selecting the best subset from more than 1000 features on a per dataset basis is not achievable. Feature selection has been investigated before, but mostly for TCC [33]–[35], while only little has been done for OCC [36]–[38]. In this study, we used different feature selection approaches and compared their effectiveness for OCC and TCC classification performance.

Both machine learning approaches, OCC and TCC, benefit from feature selection. While feature selection is essential for OCC and a difference of about 30% accuracy can be observed, the maximum difference for TCC is ~10%. Moreover, for TCC 7 out of 8 feature selection methods lead to accuracy greater than 90% whereas such high accuracy was only achieved for two methods when using OCC. For the LIG feature selection method, intended as a negative control, both classifiers display lowest performance but TCC is about 20% better than OCC. With increasing accuracy (i.e. better feature selection for OCC), the accuracy for TCC also increases; except for the SFC which is best for OCC but only third best for TCC. While the performance difference for LIG is large, it decreases with the use of better feature selection methods. TCC is only 3% better when the SFC feature selection method is considered, which provided the best performance for OCC. A difference in performance among plant species was observed for both classifiers, but for TCC it was about 5% whereas for OCC it was 15%. In conclusion, feature selection is essential for OCC, but does not affect TCC as much. We propose that due to the lack of true negative data, more focus should be put on the further development of OCC approaches to pre-miRNA detection.

## Materials and methods

### Data

Positive examples for pre-miRNAs from selected plant species were downloaded from miRBase [14] (Releases 20 and 21). *Glycine max* (gma), *Zea mays* (zma), *Sorghum bicolor* (sbi), *Physcomitrella patens* (ppt), *Arabidopsis thaliana* (ath), *Populus trichocarpa* (ptc), and *Oryza sativa* (osa) make up the positive dataset. Negative examples for miRNAs consisted of 980 pseudo pre-miRNAs from the PlantMiRNAPred dataset [39]. For these data, all pre-miRNA features were calculated as described previously [21], [23], [40]. We chose plant pre-miRNAs with large amount of pre-miRNA examples and from different clades for this study. Additionally, plant miRNAs have not been investigated as extensively as metazoan miRNAs which adds to the reason to choose plant pre-miRNAs.

### One class classification

For one-class classification the DDtools [41] implementation of an OCC was utilized. 100 fold Monte Carlo cross validation [42] was performed using randomly sampled 90% of the positive data for training and 10% for testing. Moreover, the pseudo negative sequences were injected as unknown class during testing. We employed *k*-means in this study as previously described [43] since it performed well in respect to OCC although it is a clustering algorithm. During learning, labeled examples are clustered (miRNAs and unknown) and during testing and in prediction, the label of the closest cluster is assigned to the sample.

## Two class classification

Support Vector Machines (SVMs) are used for machine learning and were first proposed by [44]. In bioinformatics and in the field of pre-miRNA detection, SVMs have been used [17], [18], [39], [45]. Here, the WEKA library [46] SVM implementation which is based on LibSVM [47] was utilized. The radial basis function was set to a gamma value of 0.7 and the cost parameter was chosen to be 4.0 and the normalization option was set to true. Any machine learning algorithm needs initial training and we performed a 10 fold Monte Carlo cross validation [42] during learning, by employing random sampling using 90% of the data for training and 10% for testing.

## Feature selection strategies

Feature selection has been shown to be an NP-hard problem and, therefore, other approximate feature selection strategies are being developed. In machine learning for pre-miRNAs, more than 1000 features have been proposed which makes feature selection especially hard. To investigate the impact of feature selection on model performance for OCC and TCC, four negative and four positive feature selection methods were designed. Previously, we found that a set of 50 to 100 features may be sufficient for successful pre-miRNA detection [21]. Using more than 50 features increases the likelihood that the feature set contains some features which may conceal differences among feature selection methods. Therefore a feature set size of 50 was selected for model training in this study.

We have previously performed feature selection for OCC [48] using similar feature selection methods as we propose here, but it is important to compare the impact between OCC and TCC.

Eight feature selection methods were devised and four of them were expected to lead to low performance while the remaining methods were thought to perform well. The former were selecting features with low information gain (LIG), random feature selection (RFS), selecting random feature from feature clusters (RFC), and selecting features from clusters (SFC). The latter were selecting features with high information gain (HIG), selecting the highest information gain from feature clusters (HIC), zero-norm feature selection (ZNF), and Pearson correlation-based feature selection (PCF).

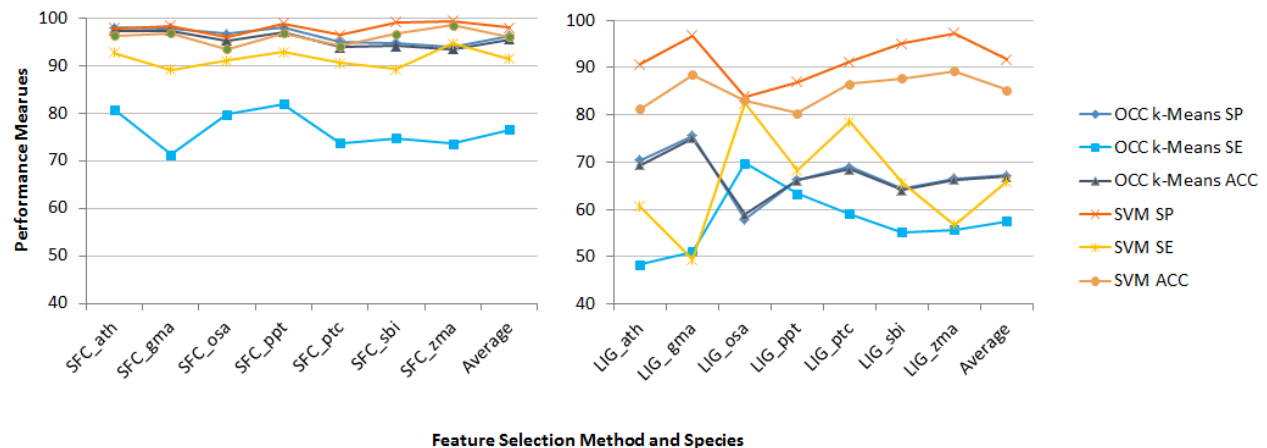
All feature selection methods except for the last two were performed using KNIME [49] and the selected features are available in Supplementary Table 1; information on how to calculate them are provided in Supplementary File 3. The workflows for our feature selection methods, developed in KNIME, are available for download from our website: <http://bioinformatics.iyte.edu.tr/supplements/featsel>.

In order to calculate LIG and HIG, for each dataset, the information gain (IG) among features was established (using KNIME's InformationGainCalculator node) and the 50 features with lowest IG (LIG) or highest IG (HIG) were selected. For RFS, 50 random features were selected using the Row Sampling node in KNIME. To establish RFC, features were clustered using WEKA *k*-Means implementation in KNIME ( $k = 100$ ). From each cluster a random feature was selected and from the 100 random features the final set of 50 features was selected randomly (KNIME's Row Sampling approach). For SFC, clustering was performed as for RFC. Clusters were ordered by number of cluster members (largest to smallest) and the 50 features were chosen from the top. To derive HIC, the same clustering approach as for RFC and SFC was taken, but the features in each cluster were ranked according to IG and the best one was selected. The selected 100 features were again ranked using IG and the best 50 were selected. ZNF is defined to be the non-zero values for all feature vectors of positive examples. Among the non-zero ones, the 50 features with highest sum of values were selected. PCF was established according to Lorena et al. [36] and after Pearson clustering the features with lowest correlation score were retained. Feature selection was performed on a per species basis which led to the selection of different features (Supplementary Table 1). Combined feature selection uses the occurrence of selected features among the 7 selected species and 5 mixed datasets in respect to the top 100 features. The Features were ranked according to their frequency and top 50 were selected (Supplementary Table 1).

## Results and discussion

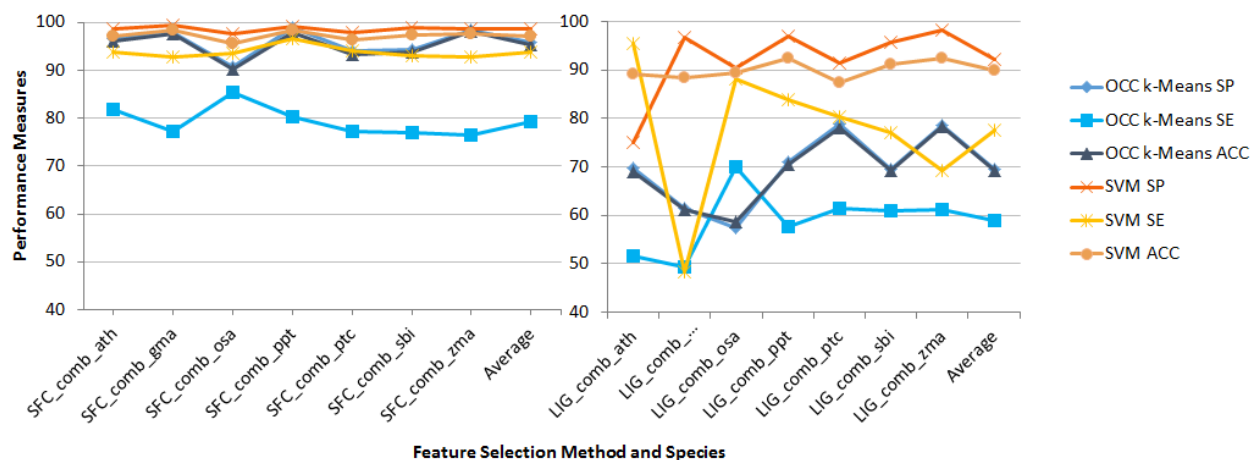
Eight feature selection methods were designed and they were applied to seven plant datasets. For each dataset OCC (100) and TCC (10) models were established using Monte Carlo cross validation (MCCV). Feature selection was performed on a per plant dataset basis. The 50 features selected varied to some extent and, therefore, we defined another feature set (indicated by 'comb') which was created by selecting the features ordered by decreasing incidence based on the individual selections. The selected features are provided in Supplementary Table 1 by their acronyms which are explained in more detailed in our previous studies [21], [24].

We applied the 8 feature selection methods to the 7 plant species' datasets individually and recorded the model performance. Figure 1 shows the average model performance (OCC: 100, TCC: 10 fold cross validation) for the best feature selection method we found (SFC) and the worst one (LIG).



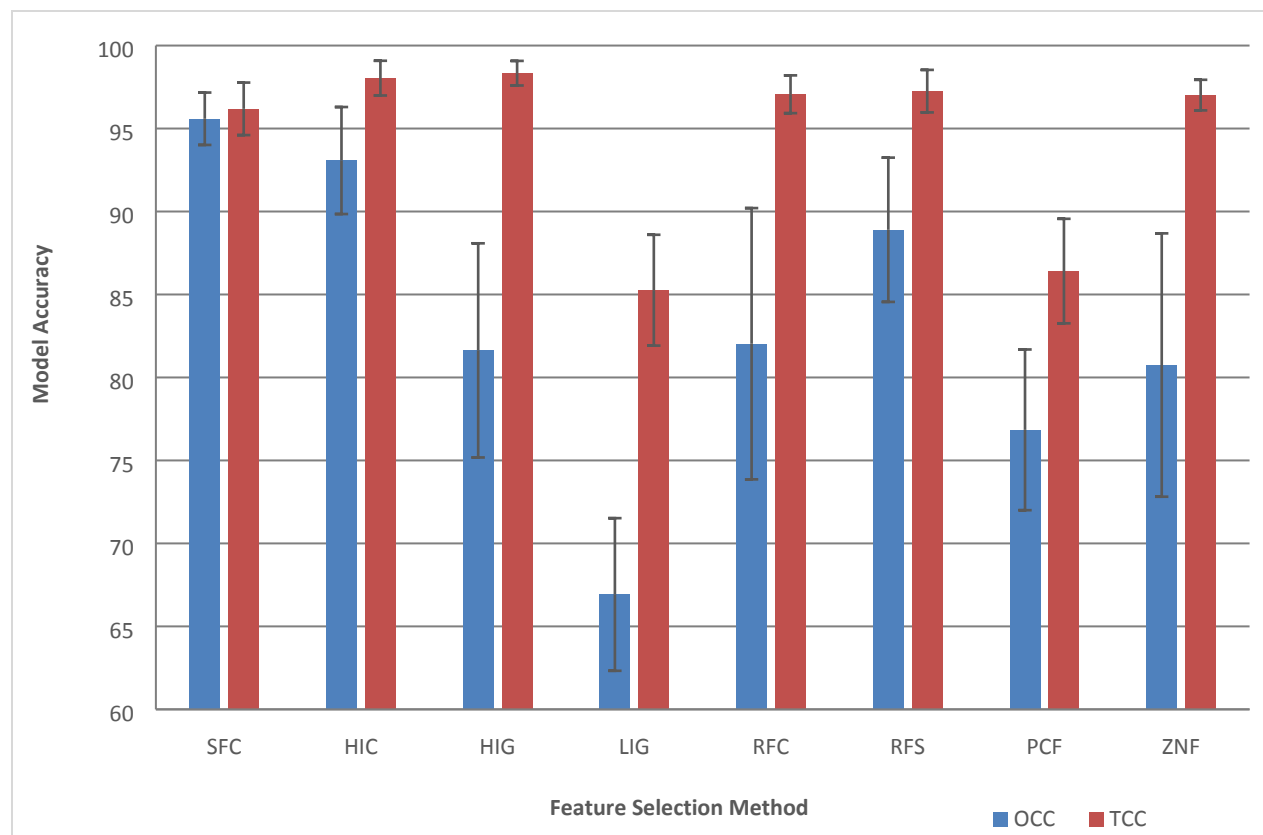
**Figure 1: Average model performance for SFC feature selection method for selected plant species (left) and LIG feature selection model (right). OCC performance is in blue tone and SVM is toned orange. SP: specificity, SE: sensitivity, and ACC: accuracy. Lines between points do not convey meaning, but were used to simplify visual tracking. Supplementary Table 2 contains further information for all feature selection methods as well as standard deviations.**

Sensitivity was the performance measure most affected for both machine learning approaches (Figure 1). For TCC the average accuracy among plant species dropped about 10% between SFC and LIG while it dropped about 30% for OCC. The results for the remaining six feature selection methods are presented in Supplementary Table 2. The impact on using combined feature selection for SFC and LIG is quite similar to individual feature selection (Figure 2). The combined features were not calculated for PCF and ZNF since combination of features was not supported by our workflow in this case. Overall accuracy is slightly reduced for the combined feature selection by on average 1% (OCC) and 2% (TCC) when compared to individual feature selection.



**Figure 2: Average model performance for SFC feature selection method using combined feature set for selected plant species (left) and LIG feature selection model using combined feature set (right). OCC performance is in blue tone and SVM is toned orange. SP: specificity, SE: sensitivity, and ACC: accuracy. Lines between points do not convey meaning, but were used to simplify visual tracking. Supplementary Table 2 contains further information for all feature selection methods as well as standard deviations.**

The performance analysis of the remaining feature selection methods are presented in Supplementary Table 2. In order to compare the performance of all feature selection methods for the two machine learning approaches, the average model accuracy was plotted (Figure 3). It is striking that for most (6 out of 8) TCC performance results the accuracy is above 95% for all plant species.

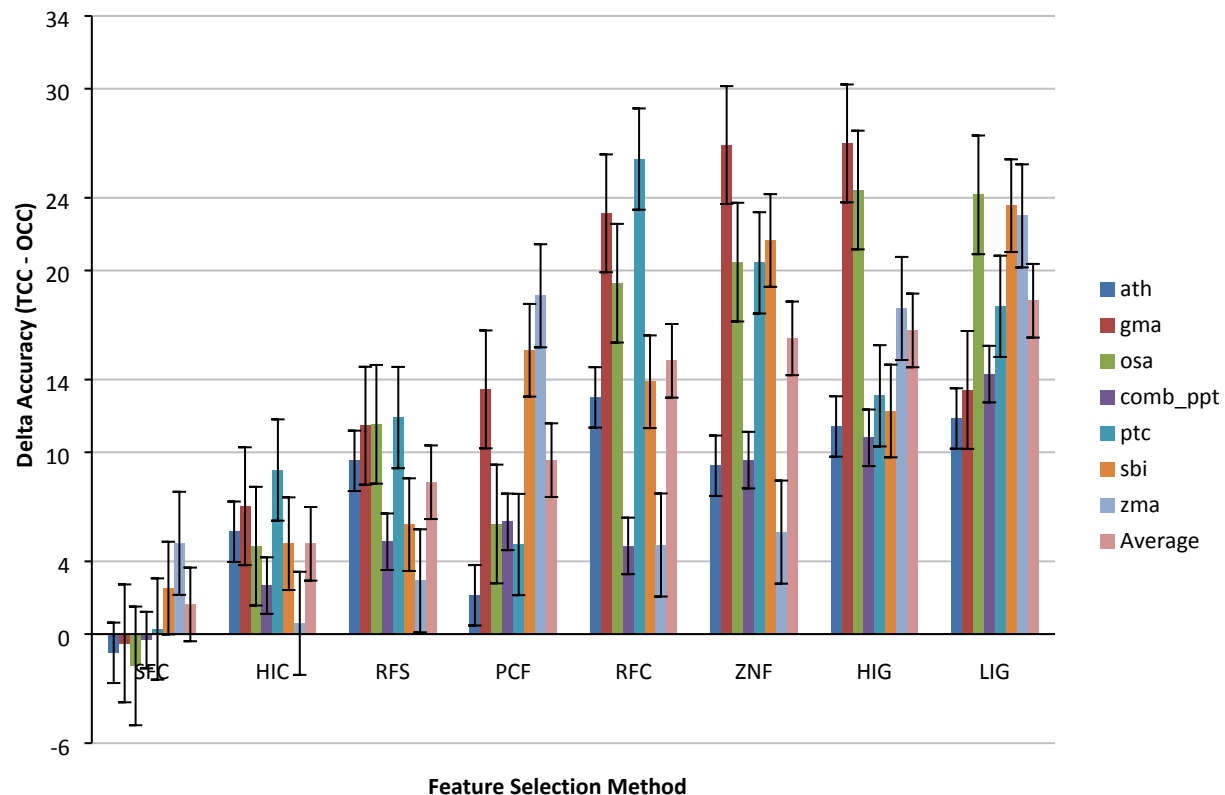


**Figure 3: The average accuracy of OCC and TCC models created for selected plant species using eight different feature selection methods and their standard deviation. Data and figure are available for closer analysis in Supplementary Table 2.**

For OCC, the performance is best for SFC where, on average, for plant species it achieves more than 95% accuracy (Figure 3). All other feature selection methods do not lead to high performing models with HIC being the second best, followed by RFS and PCF. For most feature selection methodologies the accuracy among plant species is quite similar for TCC, but for OCC the differences are much larger.

In order to compare the variance in performance between OCC and TCC, the difference between TCC and OCC accuracy was calculated ( $TCC_{ACC} - OCC_{ACC}$ ) and is presented in Figure 4. Positive values signify better performance of TCC.





**Figure 4: Eight feature selection methodologies applied to OCC and TCC. The difference in accuracy between TCC and OCC is presented. The groups are sorted by increasing average difference. Results are presented on a per species basis. The ‘Average’ averages the OCC or TCC performance among species. Data and figure are available for closer analysis in Supplementary Table 2.**

The most accurate OCC model is on the left (SFC) and it is seen that TCC is outperformed by OCC on several plant species (ath, gma, osa, and ppt). Figure 4 shows that OCC is more affected by feature selection than TCC and, therefore, with increasing effectiveness of the feature selection methodology, the difference between classifiers diminishes. For improper feature selection it can reach up to about 30%; whereas it drops to almost similar performance for the best feature selection method in this study (SFC, ~0.6% on average).

## Conclusions

Many general purpose feature selection methods have been described or used in bioinformatics [50]. For OCC feature selection nothing has been done in the area of pre-miRNA detection while one study investigated feature selection based on OCC for mature miRNA prediction [37]. When considering two class classification of pre-miRNAs SVM recursive feature elimination (RFE) has been used [51]. Meng et al. also used RFE, but modified it and compared to principal component analysis (PCA), correlation-based feature subset selection (CFS), and not using any filtering [52]. They report the best accuracy for SVM using their back SVM-RFE FS with 97.2% closely followed by PCA using SVM with 97.0% accuracy. One approach used genetic algorithm in combination with information gain and also taking into account feature redundancy for FS and achieved almost 99.5% accuracy, alas on a limited dataset [53]. These competing methods using different strategies for FS in pre-miRNA detection do not refer to OCC. However, they clearly show that feature selection has a large impact on model performance. The previous methodologies used correlation among features or feature redundancy for FS but did not put a clear focus on the correlation issue. We, therefore, devised eight feature selection methodologies with a focus on feature correlation and applied them to several plant miRNA datasets. Feature selection was performed on a per plant species basis, but we also investigated the combined feature set using the features shared among species; both of which were not done in previous approaches. Our SFC feature selection methodology was particularly successful

and there was no great difference for feature selection on a per plant basis or when combined (Figures 1, 2; left panels). As expected, the LIG methodology did not perform well at all and was intended as a negative control. However, the SVM learner was not nearly as much affected as the OCC one (Figures 1, 2; right panels) although sensitivity was strongly affected for both learners. Of the eight feature selection methods tested in this study, only 3 show good performance for OCC (SFC, HIC, and RFS; Figure 3) while only two did not seem applicable for SVM (LIG and PCF; Figure 3). For most feature selection methods average SVM performance is above 95% while OCC performance is generally below 90% (Figure 3). It is instructive to analyze the performance difference between SVM and OCC. Figure 4 shows the performance difference and it can be seen that for most feature selection methods SVM performs better than OCC (positive values; Figure 3). However, the SFC feature selection method which is among the best for SVM clearly performs best for OCC and the latter can surpass the SVM performance for several of the selected plant species. From this study it can be concluded that the more successful the feature selection the less difference between OCC and TCC model performance and the better the overall model performance. Thus we conclude, that in the absence of missing negative data OCC should be used and, therefore, additional feature selection strategies should be tried to improve its performance.

# References

- [1] A. E. Erson-Bensan, "Introduction to microRNAs in biological systems.," *Methods Mol. Biol.*, vol. 1107, pp. 1–14, Jan. 2014.
- [2] M. D. Saçar and J. Allmer, "Current Limitations for Computational Analysis of miRNAs in Cancer.," *Pakistan J. Clin. Biomed. Res.*, vol. 1, no. 2, pp. 3–5, 2013.
- [3] R. C. Lee, R. L. Feinbaum, and V. Ambros, "The C. elegans heterochronic gene lin-4 encodes small RNAs with antisense complementarity to lin-14.," *Cell*, vol. 75, no. 5, pp. 843–854, Dec. 1993.
- [4] A. Kozomara and S. Griffiths-Jones, "miRBase: integrating microRNA annotation and deep-sequencing data.," *Nucleic Acids Res.*, vol. 39, no. Database issue, pp. D152–7, Jan. 2011.
- [5] F. Grey, "Role of microRNAs in herpesvirus latency and persistence.," *J. Gen. Virol.*, vol. 96, no. Pt 4, pp. 739–51, Apr. 2015.
- [6] M. Yousef, J. Allmer, and W. Khalifaa, "Plant MicroRNA Prediction employing Sequence Motifs Achieves High Accuracy," 2015.
- [7] E. J. Chapman and J. C. Carrington, "Specialization and evolution of endogenous small RNA pathways.," *Nat. Rev. Genet.*, vol. 8, no. 11, pp. 884–896, Nov. 2007.
- [8] C. Ender and G. Meister, "Argonaute proteins at a glance," *J. Cell Sci.*, vol. 123, no. 11, pp. 1819–1823, Jun. 2010.
- [9] B. Alural, G. A. Duran, K. U. Tufekci, J. Allmer, Z. Onkal, D. Tunali, K. Genc, and S. Genc, "EPO Mediates Neurotrophic, Neuroprotective, Anti-Oxidant, and Anti-Apoptotic Effects via Downregulation of miR-451 and miR-885-5p in SH-SY5Y Neuron-Like Cells.," *Front. Immunol.*, vol. 5, no. September, p. 475, Sep. 2014.
- [10] B. Alural, A. Ozerdem, J. Allmer, K. Genc, and S. Genc, "Lithium protects against paraquat neurotoxicity by NRF2 activation and miR-34a inhibition in SH-SY5Y cells.," *Front. Cell. Neurosci.*, vol. 9, p. 209, May 2015.
- [11] Z. Zhang, J. Yu, D. Li, Z. Zhang, F. Liu, X. Zhou, T. Wang, Y. Ling, and Z. Su, "PMRD: plant microRNA database.," *Nucleic Acids Res.*, vol. 38, no. Database issue, pp. D806–13, Jan. 2010.
- [12] C. Bağcı and J. Allmer, "Removing contamination from genomic sequences based on vector reference libraries," in *2012 7th International Symposium on Health Informatics and Bioinformatics*, 2012, pp. 118–122.
- [13] C. Bağcı and J. Allmer, "One Step Forward, Two Steps Back; Xeno-MicroRNAs Reported in Breast Milk Are Artifacts," *PLoS One*, vol. 11, no. 1, p. e0145065, 2016.
- [14] S. Griffiths-Jones, H. K. Saini, S. van Dongen, and A. J. Enright, "miRBase: tools for microRNA genomics," *Nucleic Acids Res.*, vol. 36, no. Database issue, pp. D154–8, Jan. 2008.
- [15] S.-D. Hsu, F.-M. Lin, W.-Y. Wu, C. Liang, W.-C. Huang, W.-L. Chan, W.-T. Tsai, G.-Z. Chen, C.-J. Lee, C.-M. Chiu, C.-H. Chien, M.-C. Wu, C.-Y. Huang, A.-P. Tsou, and H.-D. Huang, "miRTarBase: a database curates experimentally validated microRNA-target interactions.," *Nucleic Acids Res.*, vol. 39, no. Database issue, pp. D163–9, Jan. 2011.



- [16] M. Yousef, S. Jung, L. C. Showe, and M. K. Showe, "Learning from positive examples when the negative class is undetermined--microRNA gene identification.," *Algorithms Mol. Biol.*, vol. 3, p. 2, Jan. 2008.
- [17] J. Ding, S. Zhou, and J. Guan, "MiRenSVM: towards better prediction of microRNA precursors using an ensemble SVM classifier with multi-loop features.," *BMC Bioinformatics*, vol. 11 Suppl 1, no. Suppl 11, p. S11, Jan. 2010.
- [18] Y. Wu, B. Wei, H. Liu, T. Li, and S. Rayner, "MiRPara: a SVM-based software tool for prediction of most probable microRNA coding regions in genome scale sequences.," *BMC Bioinformatics*, vol. 12, no. 1, p. 107, Jan. 2011.
- [19] W. Ritchie, D. Gao, and J. E. J. Rasko, "Defining and providing robust controls for microRNA prediction.," *Bioinformatics*, vol. 28, no. 8, pp. 1058–61, Apr. 2012.
- [20] M. D. Saçar, H. Hamzeiy, and J. Allmer, "Can MiRBase provide positive data for machine learning for the detection of MiRNA hairpins?," *J. Integr. Bioinform.*, vol. 10, no. 2, p. 215, Jan. 2013.
- [21] M. D. Saçar and J. Allmer, "Data mining for microrna gene prediction: On the impact of class imbalance and feature number for microrna gene prediction," in *2013 8th International Symposium on Health Informatics and Bioinformatics*, 2013, pp. 1–6.
- [22] J. Allmer and M. Yousef, "Computational methods for ab initio detection of microRNAs.," *Front. Genet.*, vol. 3, p. 209, Jan. 2012.
- [23] M. Yousef, J. Allmer, and W. Khalifa, "Sequence Motif-Based One-Class Classifiers Can Achieve Comparable Accuracy to Two-Class Learners for Plant microRNA Detection," *J. Biomed. Sci. Eng.*, vol. 08, no. 10, pp. 684–694, 2015.
- [24] M. D. Saçar and J. Allmer, "Comparison of Four Ab Initio MicroRNA Prediction Tools," in *Proceedings of the International Conference on Bioinformatics Models, Methods and Algorithms*, 2013, pp. 190–195.
- [25] M. D. Saçar and J. Allmer, "Machine learning methods for microRNA gene prediction.," *Methods Mol. Biol.*, vol. 1107, pp. 177–87, Jan. 2014.
- [26] J. Allmer, "Computational and bioinformatics methods for microRNA gene prediction.," *Methods Mol. Biol.*, vol. 1107, pp. 157–75, Jan. 2014.
- [27] L. B. Koski, M. W. Gray, B. F. Lang, and G. Burger, "AutoFACT: an automatic functional annotation and classification tool," *BMC Bioinformatics*, vol. 6, p. 151, 2005.
- [28] I. de On Lopes, A. Schliep, and A. C. de Lf de Carvalho, "The discriminant power of RNA features for pre-miRNA recognition," *BMC Bioinformatics*, vol. 15, no. 1, p. 124, Jan. 2014.
- [29] J. Allmer, "A Call for Benchmark Data in Mass Spectrometry-Based Proteomics," *J. Integr. OMICS*, vol. 2, no. 2, Dec. 2012.
- [30] L. M. Manevitz and M. Yousef, "One-Class SVMs for Document Classification," *J. Mach. Learn. Res.*, vol. 2, pp. 139–154, 2002.
- [31] L. Manevitz and M. Yousef, "One-Class Document Classification via Neural Networks," *Neurocomputing*, vol. 70, no. 7–9, pp. 1466–1481, Mar. 2007.
- [32] E. Amaldi and V. Kann, "On the approximability of minimizing nonzero variables or unsatisfied relations in linear systems," *Theor. Comput. Sci.*, vol. 209, no. 1–2, pp. 237–260, 1998.
- [33] S. Paul, M. Magdon-Ismail, and P. Drineas, "Feature selection for linear SVM with provable guarantees," *J. Mach. Learn. Res.*, vol. 38, pp. 735–743, 2015.
- [34] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene Selection for Cancer Classification using Support Vector Machines," *Mach. Learn.*, no. 46, pp. 389–422, 2002.
- [35] M. E. Ahsen, N. K. Singh, T. Boren, M. Vidyasagar, and M. A. White, "A new feature selection algorithm for two-class classification problems and application to endometrial cancer," in *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, 2012, pp. 2976–2982.
- [36] L. H. N. Lorena, A. C. P. L. F. Carvalho, and A. C. Lorena, "Filter Feature Selection for One-Class Classification," *J. Intell. Robot. Syst.*, pp. 1–17, Sep. 2014.
- [37] P. Xuan, M. Guo, Y. Huang, W. Li, and Y. Huang, "MaturePred: efficient identification of microRNAs within novel plant pre-miRNAs.," *PLoS One*, vol. 6, no. 11, p. e27422, Jan. 2011.
- [38] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA data mining software," *ACM SIGKDD Explor. Newsl.*, vol. 11, no. 1, p. 10, Nov. 2009.
- [39] P. Xuan, M. Guo, X. Liu, Y. Huang, W. Li, and Y. Huang, "PlantMiRNAPred: efficient classification of real and pseudo plant pre-miRNAs.," *Bioinformatics*, vol. 27, no. 10, pp. 1368–76, May 2011.
- [40] M. D. Saçar, C. Bağcı, and J. Allmer, "Computational Prediction of MicroRNAs from *Toxoplasma gondii* Potentially Regulating the Hosts' Gene Expression.," *Genomics. Proteomics Bioinformatics*, vol. 12, no. 5, pp. 228–238, Oct. 2014.

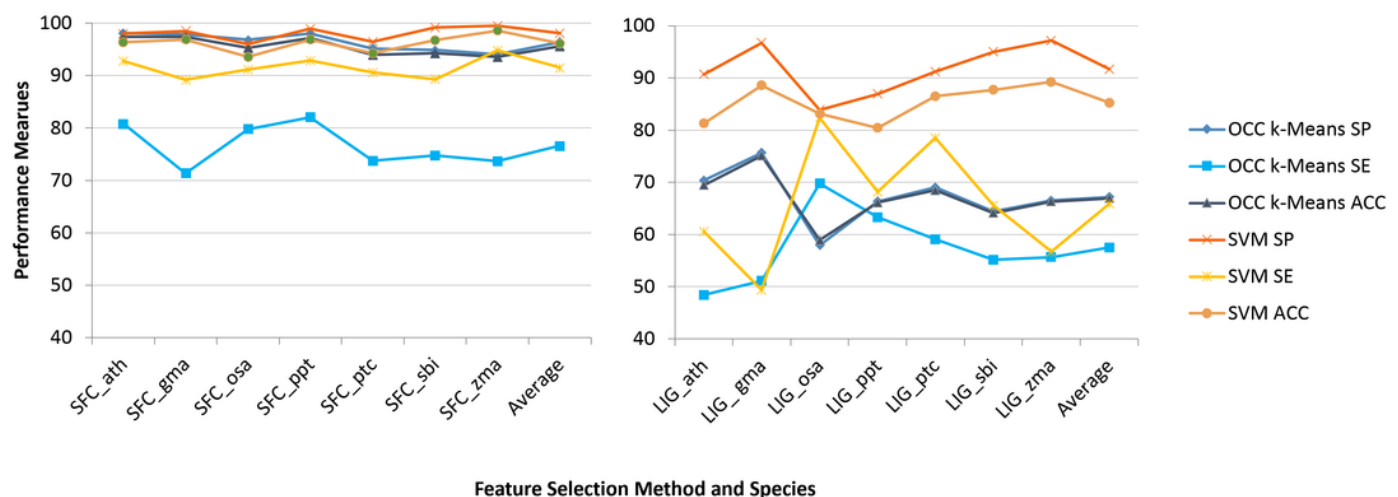
- [41] D. M. J. Tax, "DDtools, the Data Description Toolbox for Matlab." 2015.
- [42] Q.-S. Xu and Y.-Z. Liang, "Monte Carlo cross validation," *Chemom. Intell. Lab. Syst.*, vol. 56, no. 1, pp. 1–11, Apr. 2001.
- [43] M. Yousef, J. Allmer, and W. Khalifa, "Feature Selection for MicroRNA Target Prediction Comparison of One-Class Feature Selection Methodologies," in *Proceedings of the 9th International Joint Conference on Biomedical Engineering Systems and Technologies*, 2016, pp. 219–225.
- [44] V. N. Vapnik, *The nature of statistical learning theory*. New York, New York, USA: Springer-Verlag, 1995.
- [45] K. L. S. Ng and S. K. Mishra, "De novo SVM classification of precursor microRNAs from genomic pseudo hairpins using global and intrinsic folding measures.," *Bioinformatics*, vol. 23, no. 11, pp. 1321–30, Jun. 2007.
- [46] J. E. Gewehr, M. Szugat, and R. Zimmer, "BioWeka--extending the Weka framework for bioinformatics.," *Bioinformatics*, vol. 23, no. 5, pp. 651–3, Mar. 2007.
- [47] C.-C. Chang and C.-J. Lin, "LIBSVM," *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, pp. 1–27, Apr. 2011.
- [48] M. Yousef, M. D. Saçar Demirci, W. Khalifa, and J. Allmer, "Feature Selection Has a Large Impact on One-Class Classification Accuracy for MicroRNAs in Plants," *Adv. Bioinformatics*, vol. 2016, pp. 1–6, 2016.
- [49] M. R. Berthold, N. Cebon, F. Dill, T. R. Gabriel, T. Kötter, T. Meinel, P. Ohl, K. Thiel, and B. Wiswedel, "KNIME - The Konstanz Information Miner," *SIGKDD Explor.*, vol. 11, no. 1, 2009.
- [50] Y. Saeys, I. Inza, and P. Larrañaga, "A review of feature selection techniques in bioinformatics," *Bioinformatics*, vol. 23, no. 19, pp. 2507–2517, 2007.
- [51] J. Shu, K. Chiang, J. Zemleni, and J. Cui, "Computational Characterization of Exogenous MicroRNAs that Can Be Transferred into Human Circulation.," *PLoS One*, vol. 10, no. 11, p. e0140587, 2015.
- [52] J. Meng, D. Liu, C. Sun, and Y. Luan, "Prediction of plant pre-microRNAs and their microRNAs in genome-scale sequences using structure-sequence features and support vector machine.," *BMC Bioinformatics*, vol. 15, no. 1, p. 6595, 2014.
- [53] P. Xuan, M. Z. Guo, J. Wang, C. Y. Wang, X. Y. Liu, and Y. Liu, "Genetic algorithm-based efficient feature selection for classification of pre-miRNAs," *Genet. Mol. Res.*, vol. 10, no. 2, pp. 588–603, 2011.

344

1

Best (SFC) versus worst (LIG) feature selection method on per species feature selection

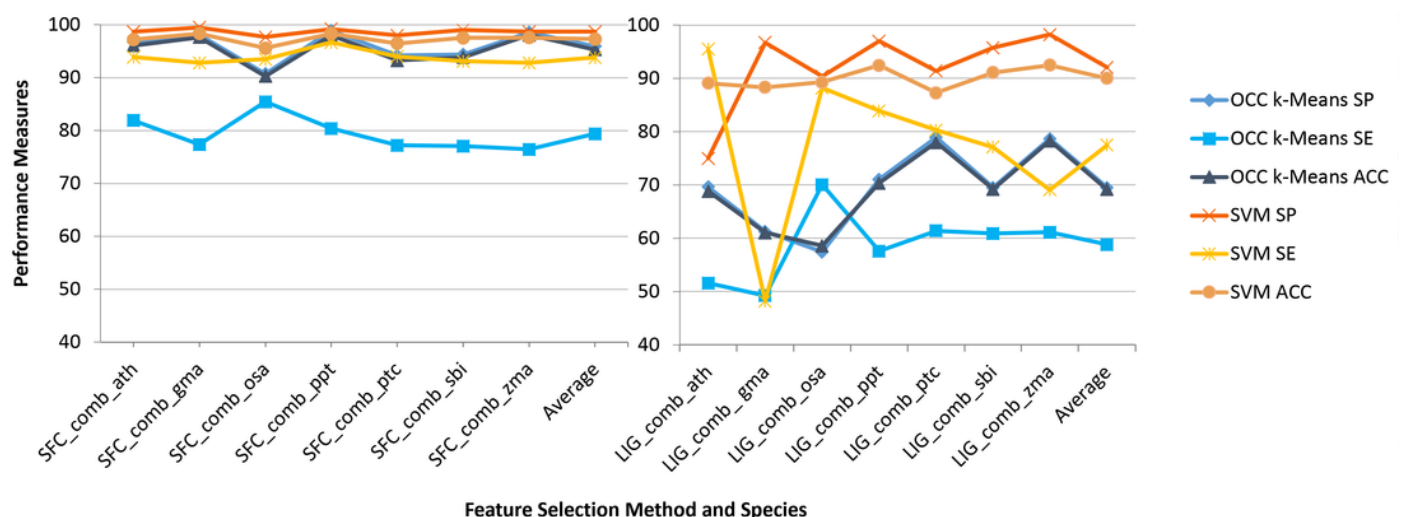
Figure 1: Average model performance for SFC feature selection method for selected plant species (left) and LIG feature selection model (right). OCC performance is in blue tone and SVM is toned orange. SP: specificity, SE: sensitivity, and ACC: accuracy. Lines between points do not convey meaning, but were used to simplify visual tracking. Supplementary Table 2 contains further information for all feature selection methods as well as standard deviations.



# 2

Best (SFC) versus worst (LIG) feature selection method on consensus feature selection

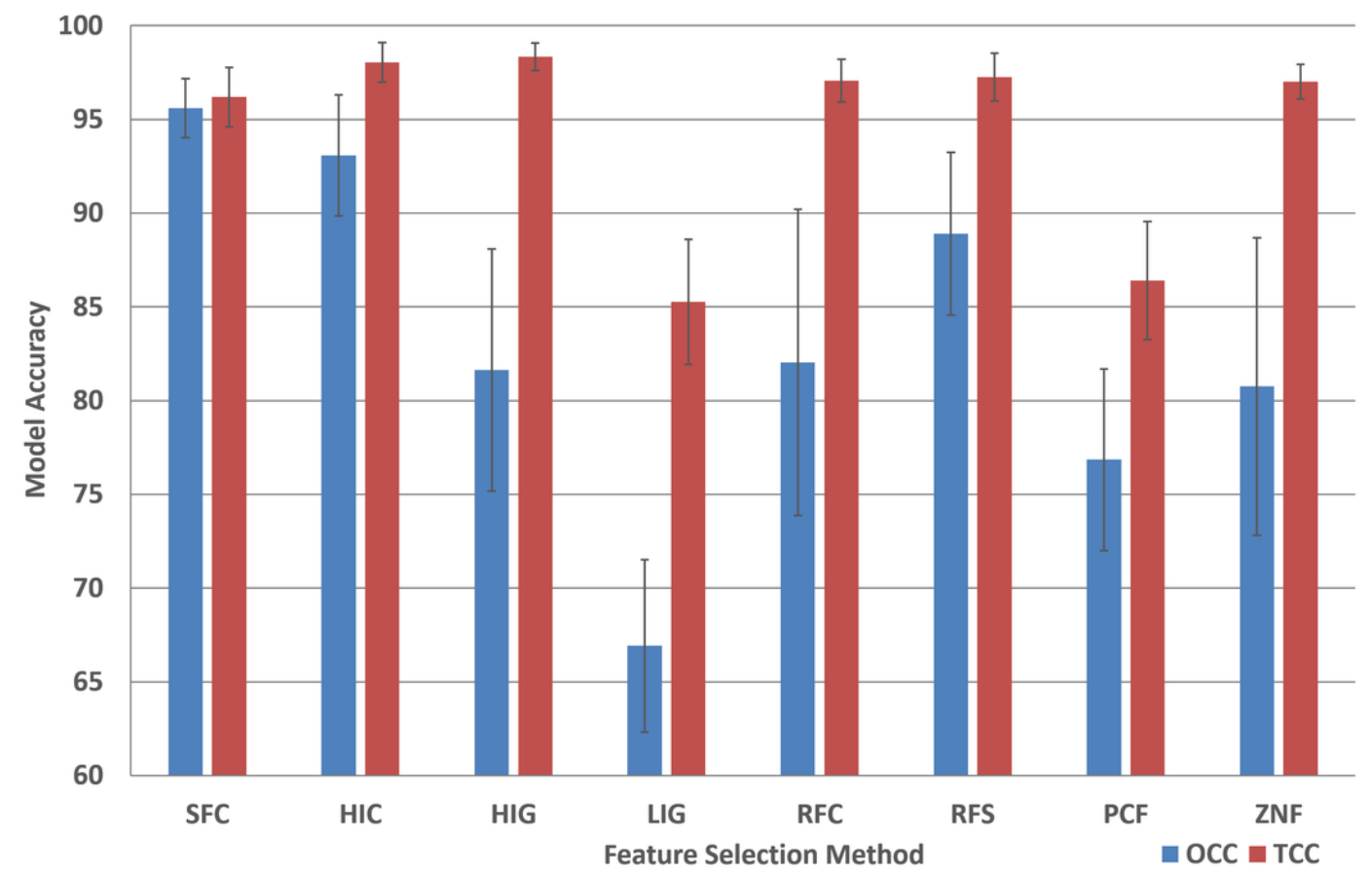
Figure 2: Average model performance for SFC feature selection method using combined feature set for selected plant species (left) and LIG feature selection model using combined feature set (right). OCC performance is in blue tone and SVM is toned orange. SP: specificity, SE: sensitivity, and ACC: accuracy. Lines between points do not convey meaning, but were used to simplify visual tracking. Supplementary Table 2 contains further information for all feature selection methods as well as standard deviations.



# 3

Model accuracy comparison between OCC and TCC in respect to feature selection method

Figure 3: The average accuracy of OCC and TCC models created for selected plant species using eight different feature selection methods and their standard deviation. Data and figure are available for closer analysis in Supplementary Table 2.



# 4

Comparison of the effect of feature selection on two-class versus one-class classification

Figure 4: Eight feature selection methodologies applied to OCC and TCC. The difference in accuracy between TCC and OCC is presented. The groups are sorted by increasing average difference. Results are presented on a per species basis. The 'Average' averages the OCC or TCC performance among species. Data and figure are available for closer analysis in Supplementary Table 2.

