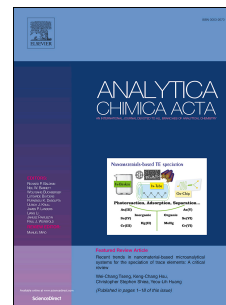# Journal Pre-proof

Open-Source Python Module for Automated Preprocessing of Near Infrared Spectroscopic Data

Jari Torniainen, Isaac O. Afara, Mithilesh Prakash, Jaakko K. Sarin, Lauri Stenroth, Juha Töyräs

Please cite this article as: J. Torniainen, I.O. Afara, M. Prakash, J.K. Sarin, L. Stenroth, J. Töyräs, Open-Source Python Module for Automated Preprocessing of Near Infrared Spectroscopic Data, *Analytica Chimica Acta*, https://doi.org/10.1016/j.aca.2020.02.030.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Credit Author Statement

**Jari Torniainen:** Investigation; Software; Writing - original draft
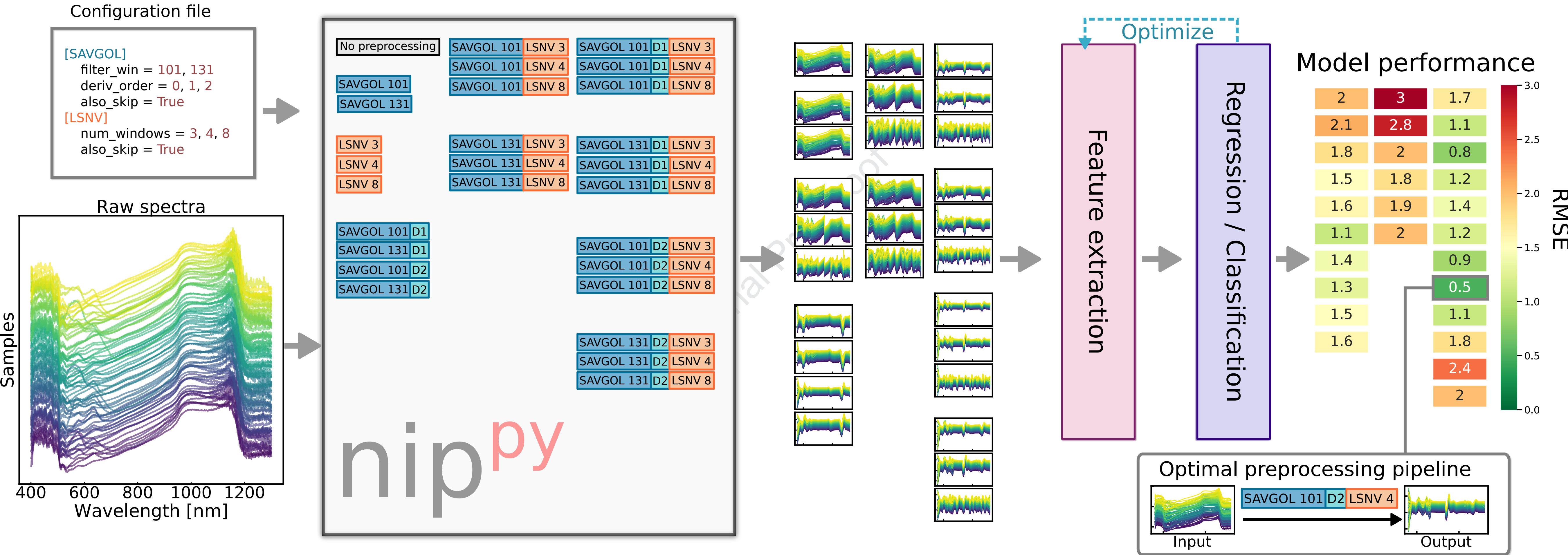**Isaac O. Afara:** Software; Writing - review & editing
**Mithilesh Prakash:** Validation; Writing - review & editing
**Jaakko K. Sarin:** Validation; Writing - review & editing
**Lauri Stenroth:** Supervision; Writing - review & editing
**Juha Töyräs:** Supervision; Project administration; Writing - review & editing

# Open-Source Python Module for Automated Preprocessing of Near Infrared Spectroscopic Data

Jari Torniainen[a,b], Isaac O. Afara[a,b], Mithilesh Prakash[a,b], Jaakko K. Sarin[a,b], Lauri Stenroth[a], Juha Töyräs[a,b,c]

[a]*Department of Applied Physics, University of Eastern Finland, Kuopio, Finland*
[b]*Diagnostic Imaging Center, Kuopio University Hospital, Kuopio, Finland*
[c]*School of Information Technology and Electrical Engineering, The University of Queensland, Brisbane, Australia*

## Abstract

Near infrared spectroscopy (NIRS) is an analytical technique for determining the chemical composition or structure of a given sample. For several decades, NIRS has been a frequently used analysis tool in agriculture, pharmacology, medicine, and petrochemistry. The popularity of NIRS is constantly growing as new application areas are discovered. Contrary to mid infrared spectral region, the absorption bands in near infrared spectral regions are often non-specific, broad, and overlapping. Analysis of NIR spectra requires multivariate methods which are highly subjective to noise arising from instrumentation, scattering effects, and measurement setup. NIRS measurements are also frequently performed outside of a laboratory which further contributes to the presence of noise. Therefore, preprocessing is a critical step in NIRS as it can vastly improve the performance of multivariate models. While extensive research regarding various preprocessing methods exists, selection of the best preprocessing method is often determined through trial-and-error. A more powerful approach for optimizing prepro-

cessing in NIRS models would be to automatically compare a large number of preprocessing techniques (e.g., through grid-search or hyperparameter tuning). To enable this, we present, `nippy`, an open-source Python module for semi-automatic comparison of NIRS preprocessing methods (available at `https://github.com/uef-bbc/nippy`). We provide here a brief overview of the capabilities of `nippy` and demonstrate the typical usage through two examples with public datasets.

## 1. Introduction

Near infrared spectroscopy (NIRS) is a widely used vibrational spectroscopic technique for quantitative evaluation of the composition and structure of a given sample. In NIRS, the target sample is illuminated with near infrared (NIR) light (750 – 2500 nm wavelength range) and the reflected and backscattered light are measured with a spectrometer. NIR-active molecular bonds in the sample absorb the incoming light at different overtone and combination spectral bands, thus producing the NIR absorbance spectrum. Compared to other infrared spectroscopy methods, NIRS has increased penetration depth and less stringent requirements for sample preparation [1]. The robustness and portability of NIRS devices enables their use as diagnostic probes outside laboratory environments [1, 2]. NIRS is a relatively old analytical technique with pertinent research spanning well over three decades. Today, NIRS is utilized heavily in multiple fields, including agriculture [3], food processing [4], pharmaceutical industry [5], and medical research [6].

The objective of a typical NIRS analysis is to relate the measured NIR

spectra to a target property of the sample which is either directly or indirectly linked to its chemical composition or structure. For instance, a common example is the determination of the octane number of gasoline samples from the corresponding NIR spectrum [7]. In this example, the number of CH-bonds in hydrocarbons, like methyl and methylene, correlates with the gasoline octane rating and can, therefore, be detected using NIRS. Compared to traditional mid infrared spectroscopy, the NIR spectrum is more complex and difficult to analyse due to the presence of overlapping overtones and combination bands. The proper analysis of NIR spectral data, therefore, requires multivariate statistical models. On the other hand, the robust-nature of NIRS instrumentation enables integration of sensors for continuous monitoring of chemical processes or as portable scanners for performing measurement and analysis of solid samples on the field (e.g., soil or crops analysis). The best practices in analysing NIRS data have been the topic of active research in the field of chemometrics [1].

The canonical analysis process of regression-based NIRS models is very similar to other chemometric applications. First, a training set of data (i.e., the calibration set) is collected by measuring both the NIR spectrum and the desired target property. The calibration set should contain a relatively large number of samples (e.g., N $= 100 - 200$, estimate derived from recommendations by Burns et al.[8] and ISO 12099:2017[9]) and cover the entire current (and future) range of expected variation of the spectra and the target property. The calibration dataset is used to construct a calibration model, capable of predicting the target property of future samples. The resulting calibration model should be reported with appropriate figures of merit (root

42 mean squared error, RMSE; coefficient of determination, $R^2$; etc.) [10, 11]
43 that approximate the real-world performance.

44 The workflow for generating the calibration model can be divided into
45 three main steps: preprocessing, variable selection (also known as frequency
46 selection or feature extraction), and calibration (i.e., training the classifi-
47 cation or regression model). Preprocessing aims to remove all sources of
48 uninformative variance (e.g., instrumentation and scattering effects) from
49 the measured spectrum. Typical preprocessing steps include normalization,
50 linearization, and smoothing [12]. The variable selection aims to reduce the
51 full feature space (i.e., different wavelengths) to a subset of most important
52 features. Variable selection methods include various decomposition methods,
53 such as principal and independent component analysis, sequential methods
54 like uninformative variable elimination, and optimization techniques like ge-
55 netic algorithms [13, 14]. In the last two cases, the process of variable selec-
56 tion is performed in tandem with the calibration step. Finally, the calibration
57 step trains a classification/regression model which maps the extracted spec-
58 tral features to explain the property of interest. Each of the three steps has
59 multiple viable methods and configurable parameters which can be tuned
60 for achieving the best possible predictive performance. The combination of
61 preprocessing, variable selection, and calibration producing the model with
62 the best performance depends heavily on the instrumentation and intended
63 application.

64 As the workflow for constructing the calibration model dictates how well
65 the technique will perform, a large number of studies have focused on op-
66 timizing this process. A substantial amount of time and effort has been

4

used to compare different preprocessing [15, 16, 17, 12], feature extraction [14, 13, 18], and classification methods [19, 20, 21] in order to build optimal predictive NIRS models. More recently, the effectiveness of different combinations of preprocessing methods (often termed "preprocessing strategies") has been investigated. Engel et al.[22] produced the first critical review of chemometric preprocessing strategies which demonstrated that the model performance between different preprocessing strategies can vary substantially and highlighted the drawbacks in methods currently in use for selecting these strategies. A study by Gerretzen et al.[23] addressed the problem by utilizing *design-of-experiments* approach for selecting the best preprocessing strategy. This approach was later expanded to cover both preprocessing and variable selection, with the aim of improving not only the model performance but also its interpretability[24]. Earlier, Bocklitz et al.[25] presented a similar solution based on genetic algorithms for finding the optimal preprocessing strategy in Raman spectroscopy. Earlier investigations of various methods have resulted in a collective consensus of good methods in the field of NIRS research. However, in order to combine these methods to form the best possible calibration model, each part of the model (i.e., preprocessing, variable selection, and calibration) should be optimized individually for each application.

Recent advances in the field of machine learning have produced powerful methods (such as automated learning and hyperparameter tuning) for seeking further optimized models. The optimization of calibration models, however, has mainly focused on tuning the variable selection and classification/regression parts of the process with little focus on the preprocessing stage. While preprocessing has a critical impact on model outcomes, it is

<sup>92</sup> often not subjected to similar pruning. Likewise, the effect of different pa-
<sup>93</sup> rameter values (e.g., filtering window length) on the model performance is
<sup>94</sup> rarely exhaustively studied or their effects reported.

<sup>95</sup> To help find an effective combination of preprocessing methods, we present
<sup>96</sup> `nippy`: a toolbox to rapidly test different preprocessing combinations in con-
<sup>97</sup> junction with different machine learning methods. `nippy` is a Python 3.6+
<sup>98</sup> module, open source under the MIT license, and fully compatible with the
<sup>99</sup> Python scientific stack (e.g., `scikit-learn` [26] machine learning module
<sup>100</sup> and `scipy` [27]). While nippy is intended for spectral preprocessing, it can
<sup>101</sup> be combined with the powerful outlier detection (e.g., elliptic envelopes, iso-
<sup>102</sup> lation forests, etc.), variable selection (e.g., variance thresholding, recursive
<sup>103</sup> feature elimination, etc.), and model validation (e.g., classification/regression
<sup>104</sup> metrics, dummy estimators, etc.) methods of the `scikit-learn` module in
<sup>105</sup> order to produce state-of-the-art calibration models. It implements all rec-
<sup>106</sup> ommended preprocessing methods (such as scatter correction, smoothing,
<sup>107</sup> and computing derivatives) as well as a framework for testing the effects of
<sup>108</sup> different combinations of preprocessing methods and their parameters. To
<sup>109</sup> be clear, we are not aiming to produce a comprehensive benchmarking of
<sup>110</sup> different preprocessing methods but rather provide a tool for researchers and
<sup>111</sup> practitioners to find an effective preprocessing strategy for their data.

<sup>112</sup> We provide here a brief overview of the module and its capabilities and
<sup>113</sup> demonstrate the typical usage through two example cases using previously
<sup>114</sup> published datasets. Full documentation and the source code of `nippy` are
<sup>115</sup> also available at `https://github.com/uef-bbc/nippy`.

## 2. Materials and methods

### 2.1. Overview of `nippy`

The *near infrared preprocessing in Python toolbox* (abbreviated as `nippy`) is a Python module that enables the user to easily define multiple different preprocessing strategies for NIRS data. The module is written in Python 3.6 and should, therefore, be compatible with all later versions. The module is also fully compatible with the data format used by the `scikit-learn`-module, enabling easy integration with powerful machine learning methods. A standard configuration file provided with the modules can be used to install `nippy` automatically using the Python package management system (i.e., `pip` or `conda`).

The `nippy` module consists of two main components: a preprocessing module and a handler function. The preprocessing module collects all currently accepted preprocessing methods under one file. Inputs to these methods have been standardized and the general naming convention from NIRS literature has been used when possible. To improve performance and reduce the burden of maintenance, the best methods from existing larger modules (such as `sklearn` and `numpy`) have been used. The modular structure of preprocessing methods also enables the user to implement custom preprocessing methods if they so desire. An overview of the preprocessing methods currently provided by `nippy` can be found in the section 2.2. The handler component of `nippy` constructs preprocessing pipelines per users requests. The user provides a list of methods and a list of parameters for each method through a configuration file. The handler parses the configuration and produces all possible permutations which can then be executed as a combination

7

141 of appropriate functions of the preprocessing module. The default usage of
142 the handler function is covered in section 2.3.

143 *2.2. Preprocessing methods*

144 In total, six different preprocessing categories are supported in `nippy` and
145 they are performed in the following order.

146 1. Clipping

147 2. Scatter correction

148 3. Smoothing

149 4. Derivatives

150 5. Trimming

151 6. Resampling

152 It should be noted, that the order in which the preprocessing operations
153 are performed can have an effect on the model performance. The order se-
154 lected for `nippy` represents the consensus among the authors and is based on
155 the order typically found in literature. If needed, the order of the operations
156 can be changed by modifying the source-code (details for this can be found in
157 the online documentation). Some of the above-mentioned operations (such
158 as scatter correction) can contain multiple alternative methods. The han-
159 dler component of `nippy` will make sure that only one operation from each
160 category is included in each of the preprocessing pipes.

161 Clipping operation removes or substitutes data points with values ex-
162 ceeding the user-defined threshold. This method is intended to eliminate
163 short spike-like artifacts that might otherwise distort further analysis. As
164 this method induces short discontinuities in the data, it is advisable to only

165 use it for gross artifacts. Larger noise-contaminated regions can be removed

166 with the trimming operation.

167 Scatter correction methods aim to counter the particle size effect. In

168 `nippy`, the three standard methods (standard normal variate, SNV [28]; mul-

169 tiplicative scatter correction, MSC [29]; and normalization [12]) have been

170 implemented. The non-parametric version of SNV, known as robust normal

171 variate (RNV) is also implemented [30]. In SNV, the correction is performed

172 according to the mean and standard deviation of the spectrum and in RNV

173 (which is intended for more noisy data) according to the median value and

174 user-specified inter-quartile interval. In addition, localized version of SNV

175 (LSNV) [31] is also included, where SNV operation is performed piece-wise

176 inside user-defined spectral windows. In MSC, correction is only possible to

177 the mean of the spectra. The extended version of MSC (EMSC) can also

178 be used, which takes into account both the linear and quadratic terms when

179 performing the correction. [32, 33] Normalization of the spectra can be per-

180 formed to the value ranges specified by the user (e.g., between 0 and 1). If

181 no normalization range is provided, each spectrum is normalized with the re-

182 spective Euclidean norm. Instead of scatter correction, `nippy` also supports

183 baseline correction which only mean-centers the spectra.

184 Smoothing of the NIR spectra can be helpful for removing environmental

185 or instrumentation-related noise. `nippy` provides both convolution-based and

186 Savitzky-Golay filtering [34] as smoothing methods. The Savitzky-Golay

187 filter is also able to return smoothed derivatives of the original spectrum.

188 Sometimes using the entire available range of wavelengths is not de-

189 sired and instead the analysis is constrained to specific optical windows [35].

9

190 Trimming operation enables the extraction of continuous and non-continuous
191 wavelength regions from the full spectral data.

192 As a final step, the spectra can be resampled to a new spectral resolu-
193 tion using the Fourier method. This step can be useful, for instance, when
194 combining spectra acquired with multiple devices having different spectral
195 resolutions.

## 196 2.3. Usage of `nippy`

197 The default format for data entry in `nippy` is `numpy`-matrices. The mod-
198 ule requires one vector specifying the wavelengths and a matrix containing
199 the spectra. The first dimension of the matrix must correspond to the wave-
200 lengths while the second dimension corresponds to samples.

201 Pipelines in `nippy` are defined using a configuration file. Structure of
202 the configuration file follows the INI format (`https://en.wikipedia.org/`
203 `wiki/INI_file`). In the configuration file, preprocessing methods are en-
204 tered as sections and parameters as *key-value* pairs. Most parameters (such
205 as the length of the filtering window) accept multiple variations of the pa-
206 rameter value which can be separated with a comma. An up-to-date list
207 of preprocessing methods and related parameters can be found in the on-
208 line documentation (`https://github.com/UEF-BBC/nippy/blob/master/`
209 `CONFIGURATION.md`). Reading the configuration file will produce a list of
210 nested dictionaries where the dictionaries contain the argument-value pairs
211 of the preprocessing functions.

212 Once the NIR data and the pipeline configuration have been loaded, they
213 can be passed to `nippy`. The handler module of `nippy` then automatically
214 applies all combinations of preprocessing methods defined in the configura-

10

tion file. The operation returns a list of copies of the original data, each modified according to the methods in the pipeline list. The preprocessed data is returned as `numpy`-arrays, meaning it can be fed directly in most machine learning methods. The module also supports iteration through a *Preprocessor*-class, which returns one preprocessed version of the original data at a time.

`nippy` also supports data export to other platforms (such as MATLAB or R). The export formats currently supported are MAT, CSV, and Pickle. An in-depth description of the data structure in exported files can be found from the online documentation.

## 3. Results

The basic operation of `nippy` (fig. 1) is demonstrated through two real-world examples. In both cases, the construction of a calibration model is optimized by finding the best NIRS preprocessing pipeline through grid search. Source code for both examples is available at `https://github.com/uef-bbc/nippy-aca-2019`.

### 3.1. Example 1: Classification of Ethiopian barley variants using NIR

The first example case is a classification task with a publicly available NIR dataset. The dataset (originally published by Kosmowski et al. [36]) consists of NIR measurements of 1200 samples of Ethiopian barley from 24 different barley cultivar variants. The objective of the example is to classify the barley cultivar variant (Ardu 1260 B, Bahati, Bekoji-1, etc.) of a given sample based on the NIR spectrum. This example represents a typical agricultural

11

application where data collected in a field with a portable device is utilized for quality control.

As in Kosmowski et al., the dataset was first split into training (70%) and testing (30%) sets. The test set was stratified over the 24 classes resulting in 15 samples per class. Preprocessing of the data was performed using `nippy`. Similar to the original analysis by Kosmowski et al., classification of the samples was performed using support vector machines (SVMs). A nu-regularized SVM with a polynomial kernel function was trained separately for each individual preprocessing pipeline and the hyperparameters of each model were optimized using a five-fold cross-validation. As the initial inspection of the NIR spectra looked noise-free and relatively flat (fig. 2a), the preprocessing techniques were restricted to Savitzky-Golay filtering and two scatter correction methods (SNV and RNV). The parameter combinations used to build these preprocessing pipelines are listed in table 1.

In total, 38 different preprocessing combinations (table 1) and one baseline model without preprocessing were tested (see fig. 2a for original data and 2b for different preprocessing methods). The effect of each preprocessing method was quantified by comparing the accuracies of hold-out test and training datasets over different pipelines (fig. 3). The best overall accuracy (82.6% for training and 87.2% for test set) was obtained using a preprocessing pipeline with SNV scatter correction and a first-order derivative Savitzky-Golay filtering (3rd order polynomial, 11-point window length).

The baseline model without any preprocessing yielded an accuracy of 75.0% for training and 80.3% for test set. The confusion matrices of the classification results were computed for the baseline model and the best per-

forming preprocessing method (fig. 4). The classification accuracy of the best performing preprocessing pipeline was identical to the best performing model reported by Kosmowski et al. [36].

3.2. Example 2: Regression model for predicting the instantaneous modulus of equine articular cartilage

The second example focuses on complex and noisy data, where NIR spectra were measured from the articular cartilage surface of equine fetlock joints (fig. 5). Articular cartilage is a layer of viscoleastic connective tissue covering the ends of articulating bones within a joint. Material properties of the cartilage layer (such as instantaneous modulus) are an important indicator of joint health. Ability to determine the material properties of cartilage during an arthroscopic procedure could have substantial diagnostic significance in identifying healthy and degraded regions of the joint.

The dataset was originally published by Sarin et al. [37] and optimal regression and variable selection methods for this data were subsequently investigated by Prakash et al. [38]. NIRS measurements were performed on 869 points from the proximal phalanx and the metacarpal bone of five horses. Instantaneous modulus at each measurement point was determined using a custom material testing device (for details see[37]). The dataset and an in-depth description of different variables can be found in [39].

The objective of this example is to construct a calibration model capable of predicting the instantaneous modulus of cartilage from the NIR spectrum. An earlier investigation of model selection and regression methods for this dataset [38] determined that the best prediction performance ($R^2 = 0.51$, RMSEP = 2.46 MPa) was obtained with a five-component partial least

13

squares regression. The preprocessing used in that analysis consisted of a third order Savitzky-Golay filtering (25 nm window size) and second order spectral derivation. Identical regression technique was used here and the effect of preprocessing was investigated using `nippy`. The model was validated with the same holdout test method (N=70, 9% of the dataset) used in the original study[38]. Baseline test set prediction performance without any kind of pretreatment of the spectra was $R^2 = 0.25$ and RMSEP = 3.06 MPa.

The pipelines (i.e., parameter combinations) used for the preprocessing consisted of 3rd order Savitzky-Golay filtering with (12 window sizes, up to the 2nd derivative) and convolution filtering (12 window sizes). Particle size effects were compensated using MSC, SNV, LSNV (four different window sizes), and RNV (four different interquartile ranges). Preprocessing was performed on the full spectral range, as well as subsets of 700 - 900 nm and 850 - 1050 nm. The absence of each aforementioned preprocessing step was also investigated by leaving them out of the analysis. In total, `nippy` was used to generate 1618 comparable pipelines (table 2) which were then compared using PLSR. Number of components for each PLSR model was determined using five-fold cross-validation. Best test set prediction performance in terms of $R^2$ was obtained using the wavelength range of 700 – 950 nm, RNV scatter correction (85% – 15% interquartile range), and a Savitzky-Golay filtering with 73 nm window size (see figs. 6A and 6B). In comparison to the baseline PLSR model, the preprocessing pipeline increased the coefficient of determination by approximately 38% ($R^2 = 0.63$, RMSEP = 2.15 MPa, see fig. 6C). Preprocessing also improved the residuals of the model in terms of magnitude, homoscedasticity, and normality (fig. 6D). Compared

14

<sub>313</sub> to the performance reported earlier with this dataset, optimization of the

<sub>314</sub> preprocessing step provided a modest improvement of approximately 12% in

<sub>315</sub> terms of $R^2$ and 13% in terms of RMSEP.

Table 1: Configuration parameters for example 1.

| Preprocessing operation | Parameter | Values |
|---|---|---|
| SNV | snv_type: | snv, rnv |
| | also_skip: | True |
| SAVGOL | filter_win: | 11,21,51,101 |
| | deriv_order: | 0, 1, 2 |
| | poly_order: | 3 |
| | also_skip: | True |

## 4. Discussion

<sub>316</sub>

<sub>317</sub>     Preprocessing of NIR spectrum is a fundamental part of any NIRS ap-

<sub>318</sub> plication. An optimized preprocessing protocol can substantially improve

<sub>319</sub> the predictive capabilities of NIRS models. In this paper, we presented an

<sub>320</sub> open-source Python module for semi-automatic exploration and comparison

<sub>321</sub> of different preprocessing strategies. Ideally, the tools introduced in this pa-

<sub>322</sub> per should cut down development time when researching or building new

<sub>323</sub> NIR-based analytical applications.

<sub>324</sub>     While improving the prediction performance is the main use of `nippy`, it

<sub>325</sub> also enhances understanding why some preprocessing improves the result by

<sub>326</sub> revealing additional details about the underlying phenomena. By comparing

Table 2: Configuration parameters for example 2.

| Preprocessing operation | Parameter | Values |
|---|---|---|
| MSC | `also_skip:` | True |
| SNV | `also_skip:` | True |
| RNV | `iqr:` | 75-25, 85-15, 65-35, 70-30 |
| | `also_skip:` | True |
| LSNV | `num_windows:` | 3, 5, 7, 9 |
| | `also_skip:` | True |
| SAVGOL | `filter_win:` | 11,25,41,55,71,85,101,115,131,145,161,175 |
| | `deriv_order:` | 0, 1, 2 |
| | `poly_order:` | 3 |
| | `also_skip:` | True |
| SMOOTH | `filter_win:` | 11,25,41,55,71,85,101,115,131,145,161,175 |
| | `also_skip:` | True |
| TRIM | `bins:` | 700 - 950, 850 - 1050 |
| | `also_skip:` | True |

²⁷ which methods work and which do not can yield more insight into what kind

²⁸ of instrumentation is needed by the application. For example, the substantial

²⁹ increase in accuracy as a result of derivation in the Ethiopian barley classifi-

³⁰ cation (section 3.1) most likely indicates that the differentiating factors are

³¹ minute spectral peaks and not the baseline level of the signal.

³²     Software tools for chemometric analysis of NIRS data have existed for a

³³ long time and range from proprietary analysis solutions, such as The Un-

scrambler X (Camo Analytics, Oslo, Norway), OPUS (Bruker Corporation, Billerica, MA, USA), or Pirouette (Infometrix Inc., Bothell, WA, USA), to open-source libraries. While proprietary software can be useful in industrial applications, research of new NIRS analysis techniques is typically conducted with data science and programming oriented methods using tools, such as R, MATLAB, or Python. The main benefits of open-source tools are transparency, customizability, easier access, and lower cost. We have, therefore, limited our comparison of `nippy` to other comparable free or open-source tools, more specifically, `prospectr` [40] and ParLeS [41]

The `prospectr` is one of the most popular R packages for analysing visible and NIR spectroscopic data. The package implements a set of preprocessing and sampling functions in the R language. Preprocessing methods included in `prospectr` are largely similar to those found in `nippy`. The sampling methods present different techniques for selecting training, testing, and validation sets for constructing the calibration models (e.g., Kennard-Stone sampling, DUPLEX sampling, etc.). Variable selection and calibration methods are not included in the `prospectr` package. ParLes is a shareware software solution for constructing NIRS calibration models. It implements the full chain of operations ranging from preprocessing and variable selection to training calibration models. Like `prospectr`, various sampling methods for training the models are also provided.

In comparison, `nippy` shares more features with the `prospectr` package than ParLeS, in the sense that they are both libraries containing functions for preprocessing NIRS data. Main difference between the two is that the `prospectr` package does not enable rapid iteration of multiple comparable

17

preprocessing pipelines. ParLes is more of a stand-alone tool for data exploration that enables the user to perform multiple different manipulations and analysis on the given data. Again, however, the ability to rapidly and programmatically test the effect of different preprocessing pipelines is missing. Furthermore, integration to other analysis platforms (e.g. Python or MATLAB) can not be done directly. As ParLeS is a shareware application, extending the analysis capabilities of the tool is impossible without the help of the original author.

Several earlier chemometric studies have also indicated that preprocessing has a substantial impact on the performance spectroscopic models and depends on such factors as: the preprocessing operations, parameters used, and the order of operations [42, 17, 25, 22, 23, 24]. In addition, Engel et al. [22] pointed out that the sequential optimization of a preprocessing pipeline might not work, as the synergy between different operations can be hard to predict in advance. These studies have also suggested different approaches for finding a suitable preprocessing strategy from multiple comparable alternatives. For instance, Wold et al. [42] suggested utilizing orthogonal signal correction, where preprocessing methods aim to remove linearly uncorrelated spectral components with respect to the target property. Solution proposed by Xu et al. [17] used Monte Carlo sampling for selecting the best preprocessing pipeline from multiple alternatives. Other approaches have suggested genetic algorithms [25] or design of experiments approach [23] for selecting the preprocessing strategy. Combination of preprocessing and variable selection under the same optimization process has also been investigated [24]. While exhaustive grid-search never fails to find the optimal preprocessing

18

combination for a given dataset (provided that the grid of operations and related parameters is dense enough), the operation can be very time-consuming. The proposed hyperparameter tuning methods can converge faster and can enable larger search space for parameters. As a final step, generalization of the optimal solution outside the training data should be verified in order to avoid overfitting.

While the current implementation of `nippy` enables the user to rapidly sift through various combinations of preprocessing operations, it does not provide actual feedback on which of the combinations is the most effective. This feature was intentionally left out as metaheuristics (also known as automated machine learning or hyperparameter search), a relatively recent and very active field of research, can deal with higher level optimization of machine learning pipelines. Metaheuristics can be utilized to find the best solution to a problem by individually tuning the preprocessing, feature extraction, and classification/regression. Several powerful tools, ranging from bayesian optimization [43] to genetic programming [44], already exist for solving this problem. Instead of competing with existing tools, `nippy` was built with focus on compatibility with modern metaheuristic tools. As machine learning tools are nowadays predominantly written in Python and use `numpy`-based matrix structures as its base, the same approach was adopted here. The end goal was to combine decades worth of domain knowledge gained from various NIRS publications and combine it with powerful machine learning frameworks.

In the future, `nippy` could potentially be extended to cover other spectroscopic techniques as well. Raman and mid infrared spectroscopy are very similar to NIRS and, thus, benefit from the same preprocessing operations.

19

However, as NIRS is often used in industrial and agricultural applications, the most crucial need for optimized preprocessing is with this technique.

## 5. Conclusions

To conclude, a lot is known about the different ways to eliminate scattering effects and external noise from NIR spectra. Due to the great diversity of different NIRS applications, the best preprocessing strategy is often dependent on the intended use. Finding the correct combination of preprocessing, variable selection, and calibration has been the focus of much recent research. To facilitate the optimization of preprocessing for NIRS models, we have developed `nippy`, a tool that enables rapid iteration of different preprocessing combinations. We feel that tools, such as `nippy`, are important to chemometrics because they provide researchers easy access to current state-of-art NIR preprocessing and thus enable them to focus on optimizing models instead of reinventing the wheel.

## 6. Acknowledgements

20

## 7. Contributions

Jari Torniainen: Investigation; Software; Writing - original draft

Isaac O. Afara: Software; Writing - review & editing

Mithilesh Prakash: Validation; Writing - review & editing

Jaakko Sarin: Validation; Writing - review & editing

Lauri Stenroth: Supervision; Writing - review & editing

Juha Töyräs: Supervision; Project administration; Writing - review & editing

[1] C. Pasquini, Near infrared spectroscopy: A mature analytical technique with new perspectives – A review, Analytica Chimica Acta 1026 (2018) 8–36.

[2] C. P. Bacon, Y. Mattley, R. DeFrece, Miniature spectroscopic instrumentation: Applications to biology and chemistry, Review of Scientific Instruments 75 (2004) 1–16.

[3] R. V. Rossel, D. Walvoort, A. McBratney, L. J. Janik, J. Skjemstad, Visible, near infrared, mid infrared or combined diffuse reflectance spectroscopy for simultaneous assessment of various soil properties, Geoderma 131 (2006) 59–75.

[4] N. Prieto, R. Roehe, P. Lavin, G. Batten, S. Andres, Application of near infrared reflectance spectroscopy to predict meat and meat products quality: A review, Meat Science 83 (2009) 175–186.

[5] Y. Roggo, P. Chalus, L. Maurer, C. Lema-Martinez, A. Edmond, N. Jent, A review of near infrared spectroscopy and chemometrics in

453  pharmaceutical technologies, Journal of pharmaceutical and biomedical
454  analysis 44 (2007) 683–700.

455  [6] V. R. Kondepati, H. M. Heise, J. Backhaus, Recent applications of near-
456  infrared spectroscopy in cancer diagnosis and therapy, Analytical and
457  Bioanalytical Chemistry 390 (2008) 125–139.

458  [7] J. H. Kalivas, Two data sets of near infrared spectra, Chemometrics
459  and Intelligent Laboratory Systems 37 (1997) 255–259.

460  [8] D. A. Burns, E. W. Ciurczak, Handbook of near-infrared analysis, CRC
461  press, 2007.

462  [9] ISO 12099:2017, Animal feeding stuffs, cereals and milled cereal prod-
463  ucts —- Guidelines for the application of near infrared spectrometry,
464  Standard, International Organization for Standardization, Geneva, CH,
465  2017.

466  [10] P. Williams, P. Dardenne, P. Flinn, Tutorial: Items to be included in a
467  report on a near infrared spectroscopy project, Journal of Near Infrared
468  Spectroscopy 25 (2017) 85–90.

469  [11] M. Larrechi, M. Callao, Strategy for introducing nir spectroscopy and
470  multivariate calibration techniques in industry, TrAC Trends in Ana-
471  lytical Chemistry 22 (2003) 634–640.

472  [12] Å. Rinnan, F. van den Berg, S. B. Engelsen, Review of the most common
473  pre-processing techniques for near-infrared spectra, TrAC - Trends in
474  Analytical Chemistry 28 (2009) 1201–1222.

[13] Z. Xiaobo, Z. Jiewen, M. J. Povey, M. Holmes, M. Hanpin, Variables selection methods in near-infrared spectroscopy, Analytica chimica acta 667 (2010) 14–32.

[14] R. M. Balabin, S. V. Smirnov, Variable selection in near-infrared spectroscopy: benchmarking of feature selection methods on biodiesel data, Analytica chimica acta 692 (2011) 63–72.

[15] A. Candolfi, R. De Maesschalck, D. Jouan-Rimbaud, P. A. Hailey, D. L. Massart, The influence of data pre-processing in the pattern recognition of excipients near-infrared spectra, Journal of Pharmaceutical and Biomedical Analysis 21 (1999) 115–132.

[16] L. J. Chen, L. Xing, L. J. Han, Influence of Data Preprocessing on the Quantitative Determination of Nutrient Content in Poultry Manure by Near Infrared Spectroscopy, Journal of Environment Quality 39 (2010) 1841.

[17] L. Xu, Y. P. Zhou, L. J. Tang, H. L. Wu, J. H. Jiang, G. L. Shen, R. Q. Yu, Ensemble preprocessing of near-infrared (NIR) spectra for multivariate calibration, Analytica Chimica Acta 616 (2008) 138–143.

[18] L.-L. Wang, Y.-W. Lin, X.-F. Wang, N. Xiao, Y.-D. Xu, H.-D. Li, Q.-S. Xu, A selective review and comparison for interval variable selection in spectroscopic modeling, Chemometrics and Intelligent Laboratory Systems 172 (2018) 229–240.

[19] A. Mouazen, B. Kuang, J. De Baerdemaeker, H. Ramon, Comparison among principal component, partial least squares and back propagation

498     neural network analyses for accuracy of measurement of selected soil

499     properties with visible and near infrared spectroscopy, Geoderma 158

500     (2010) 23–31.

501  [20] R. M. Balabin, R. Z. Safieva, E. I. Lomakina, Comparison of linear and

502     nonlinear calibration models based on near infrared (nir) spectroscopy

503     data for gasoline properties prediction, Chemometrics and intelligent

504     laboratory systems 88 (2007) 183–188.

505  [21] R. M. Balabin, R. Z. Safieva, E. I. Lomakina, Gasoline classification us-

506     ing near infrared (NIR) spectroscopy data: Comparison of multivariate

507     techniques, Analytica Chimica Acta 671 (2010) 27–35.

508  [22] J. Engel, J. Gerretzen, E. Szymańska, J. J. Jansen, G. Downey,

509     L. Blanchet, L. M. Buydens, Breaking with trends in pre-processing?,

510     TrAC - Trends in Analytical Chemistry 50 (2013) 96–106.

511  [23] J. Gerretzen, E. Szymańska, J. J. Jansen, J. Bart, H. J. Van Manen,

512     E. R. Van Den Heuvel, L. M. Buydens, Simple and Effective Way for

513     Data Preprocessing Selection Based on Design of Experiments, Analyt-

514     ical Chemistry 87 (2015) 12096–12103.

515  [24] J. Gerretzen, E. Szymańska, J. Bart, A. N. Davies, H. J. van Manen,

516     E. R. van den Heuvel, J. J. Jansen, L. M. Buydens, Boosting model

517     performance and interpretation by entangling preprocessing selection

518     and variable selection, Analytica Chimica Acta 938 (2016) 44–52.

519  [25] T. Bocklitz, A. Walter, K. Hartmann, P. Rösch, J. Popp, How to pre-

520 process Raman spectra for reliable and stable models?, Analytica Chim-
521 ica Acta 704 (2011) 47–56.

522 [26] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion,
523 O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al.,
524 Scikit-learn: Machine learning in python, Journal of machine learning
525 research 12 (2011) 2825–2830.

526 [27] E. Jones, T. Oliphant, P. Peterson, et al., Scipy: Open source scientific
527 tools for python (2001).

528 [28] R. J. Barnes, M. S. Dhanoa, S. J. Lister, Standard normal variate trans-
529 formation and de-trending of near-infrared diffuse reflectance spectra,
530 Applied Spectroscopy 43 (1989) 772–777.

531 [29] H. Martens, S. Jensen, P. Geladi, Multivariate linearity transformation
532 for near-infrared reflectance spectrometry, in: Proceedings of the Nordic
533 symposium on applied statistics, Stokkand Forlag Publishers Stavanger,
534 Norway, pp. 205–234.

535 [30] Q. Guo, W. Wu, D. L. Massart, The robust normal variate transform
536 for pattern recognition with near-infrared data, Analytica Chimica Acta
537 382 (1999) 87–103.

538 [31] Y. Bi, K. Yuan, W. Xiao, J. Wu, C. Shi, J. Xia, G. Chu, G. Zhang,
539 G. Zhou, A local pre-processing method for near-infrared spectra, com-
540 bined with spectral segmentation and standard normal variate transfor-
541 mation, Analytica Chimica Acta 909 (2016) 30–40.

[32] H. Martens, E. Stark, Extended multiplicative signal correction and spectral interference subtraction: new preprocessing methods for near infrared spectroscopy, Journal of pharmaceutical and biomedical analysis 9 (1991) 625–635.

[33] N. K. Afseth, A. Kohler, Extended multiplicative signal correction in vibrational spectroscopy, a tutorial, Chemometrics and Intelligent Laboratory Systems 117 (2012) 92–99.

[34] A. Savitzky, M. J. Golay, Smoothing and differentiation of data by simplified least squares procedures., Analytical chemistry 36 (1964) 1627–1639.

[35] L. A. Sordillo, Y. Pu, S. Pratavieira, Y. Budansky, R. R. Alfano, Deep optical imaging of tissue using the second and third near-infrared spectral windows, Journal of Biomedical Optics 19 (2014) 056004.

[36] F. Kosmowski, T. Worku, Evaluation of a miniaturized nir spectrometer for cultivar identification: The case of barley, chickpea and sorghum in ethiopia, PloS one 13 (2018) e0193620.

[37] J. K. Sarin, M. Amissah, H. Brommer, D. Argüelles, J. Töyräs, I. O. Afara, Near infrared spectroscopic mapping of functional properties of equine articular cartilage, Annals of biomedical engineering 44 (2016) 3335–3345.

[38] M. Prakash, J. K. Sarin, L. Rieppo, I. O. Afara, J. Töyräs, Optimal regression method for near-infrared spectroscopic evaluation of articular cartilage, Applied spectroscopy 71 (2017) 2253–2262.

26

565 [39] J. K. Sarin, J. Torniainen, M. Prakash, L. Rieppo, I. O. Afara, J. Töyräs,
566      Dataset on equine cartilage near infrared spectra, composition, and func-
567      tional properties, Scientific Data 6 (2019) 164.

568 [40] A. Stevens, L. Ramirez-Lopez, An introduction to the prospectr package
569      (2013) 1–22.

570 [41] R. A. Viscarra Rossel, ParLeS: Software for chemometric analysis of
571      spectroscopic data, Chemometrics and Intelligent Laboratory Systems
572      90 (2008) 72–83.

573 [42] S. Wold, H. Antti, F. Lindgren, J. Öhman, Orthogonal signal correc-
574      tion of near-infrared spectra, Chemometrics and Intelligent Laboratory
575      Systems 44 (1998) 175–185.

576 [43] M. Feurer, A. Klein, K. Eggensperger, J. Springenberg, M. Blum,
577      F. Hutter, Efficient and robust automated machine learning, in: Ad-
578      vances in Neural Information Processing Systems, pp. 2962–2970.

579 [44] R. S. Olson, J. H. Moore, Tpot: A tree-based pipeline optimization tool
580      for automating machine learning, in: Workshop on Automatic Machine
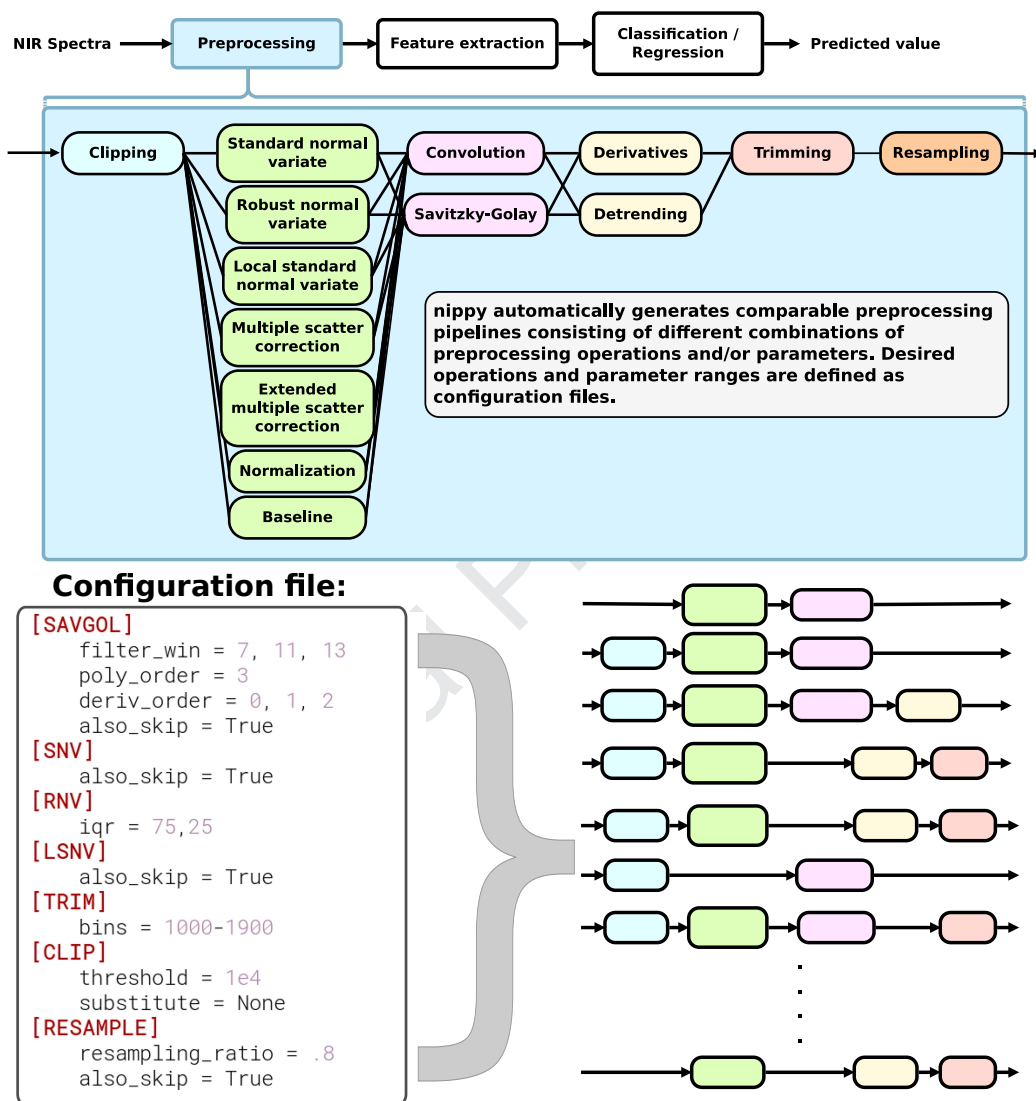581      Learning, pp. 66–74.

27

## Configuration file:

```
[SAVGOL]
    filter_win = 7, 11, 13
    poly_order = 3
    deriv_order = 0, 1, 2
    also_skip = True
[SNV]
    also_skip = True
[RNV]
    iqr = 75,25
[LSNV]
    also_skip = True
[TRIM]
    bins = 1000-1900
[CLIP]
    threshold = 1e4
    substitute = None
[RESAMPLE]
    resampling_ratio = .8
    also_skip = True
```

Figure 1: General operation principle of nippy. Several preprocessing options and parameter values can be combined into multiple competing preprocessing strategies (called pipelines). Comparing the effect of different preprocessors on the prediction performance of the NIRS model will yield the most optimal solution for a given application.

28

Figure 2: A: Overview of the Etiophian barley NIRS data set. Individual lines correspond to the averaged unprocessed spectra of the 24 different barley variants. B: The effect of different parameter combinations of scatter correction and smoothing to the per-class average spectra of the dataset. Rows correspond to different scatter correction treatment while columns represent different filter window lengths used in Savitzky-Golay filtering (3rd polynomial order, no derivation).
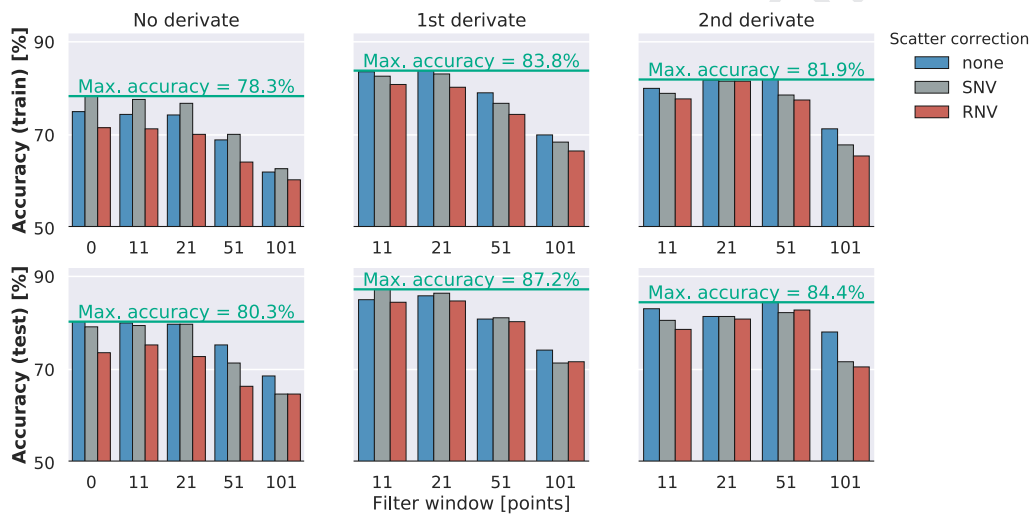
Figure 3: Overall training and testing accuracy with SVM classification for barley cultivars of example 1. Training accuracy was derived from cross-validation while testing accuracy represents the performance of an independent hold-out set. Results have been divided between filtering window length (the effect of omitting filtering was also investigated), scatter correction method, and derivative.
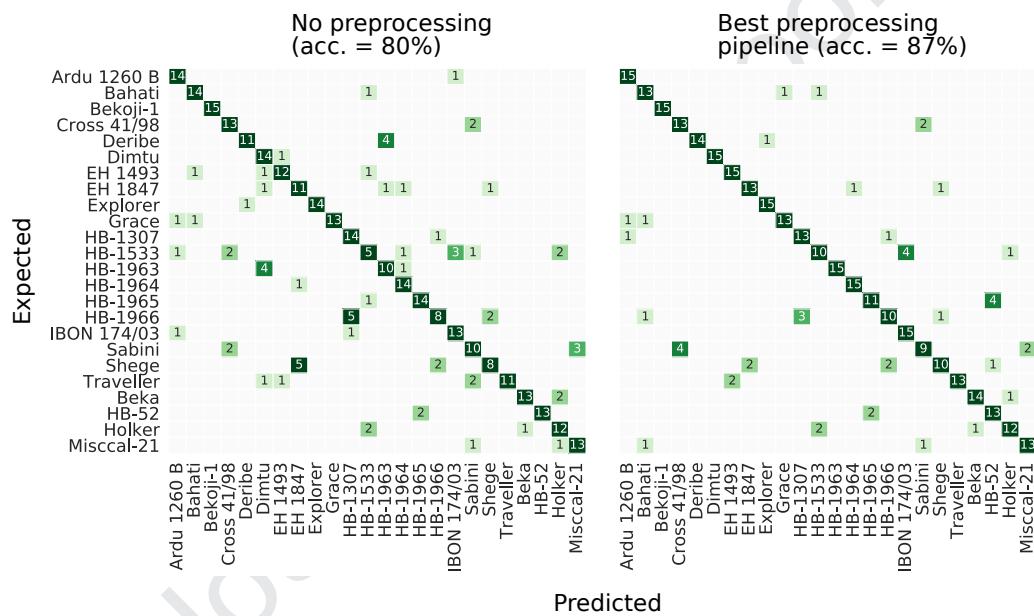
Figure 4: A comparison of confusion matrices between the un-preprocessed data and the data preprocessed with the optimal preprocessing combination of the barley cultivar classification example.
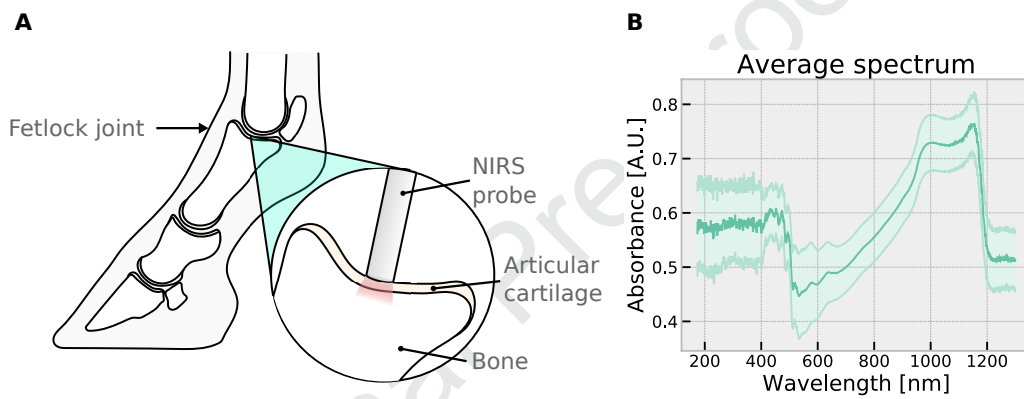
31

Figure 5: A: NIRS spectra was measured from the cartilage surface of the equine fetlock joint. In total, five joints were measured from 44 different areas of interest (AIs). Each AI consisted between 6–25 measurement points (depending on the joint geometry and cartilage condition) resulting in 861 measurement points. Local instantaneous modulus was determined for each measurement site. B: Average of all the collected NIR spectra. Shaded regions represent the standard deviation.
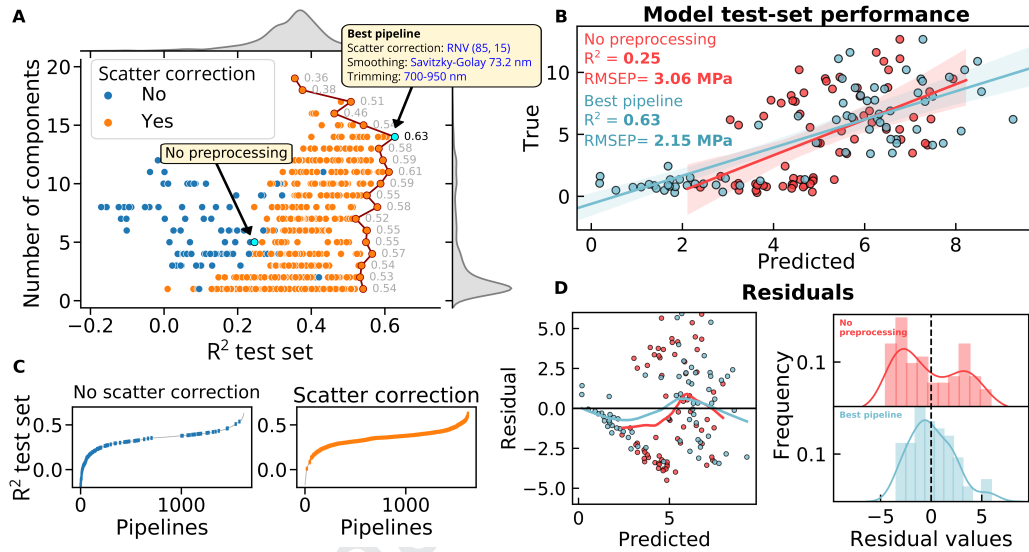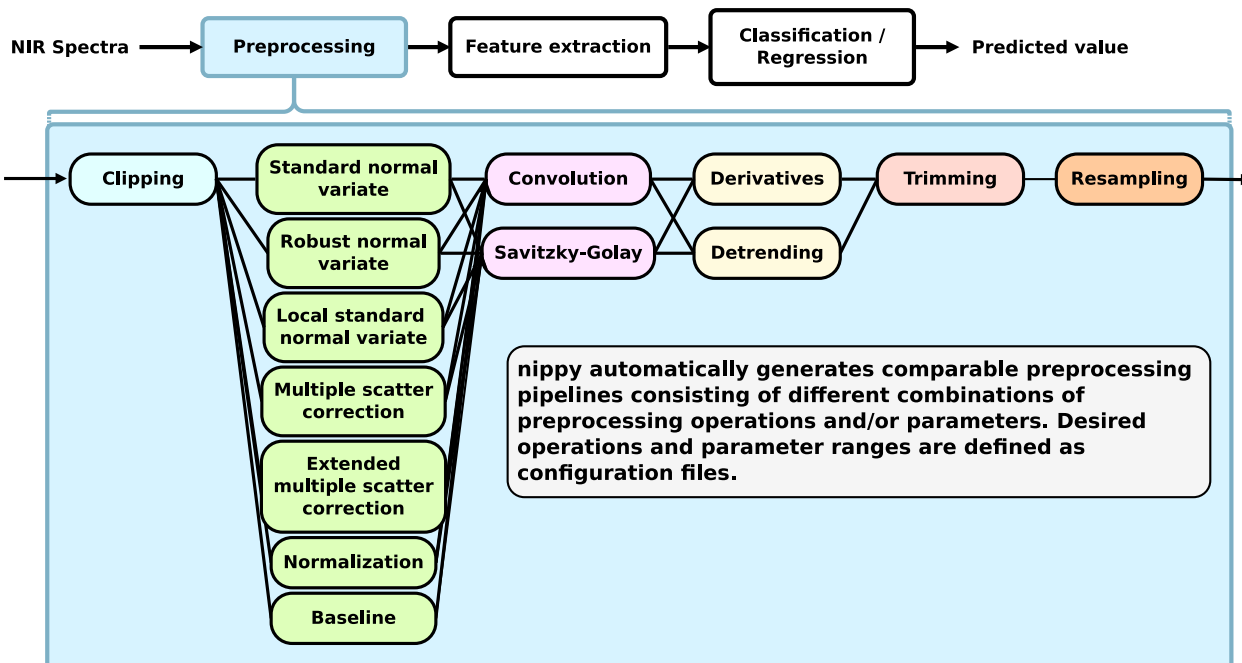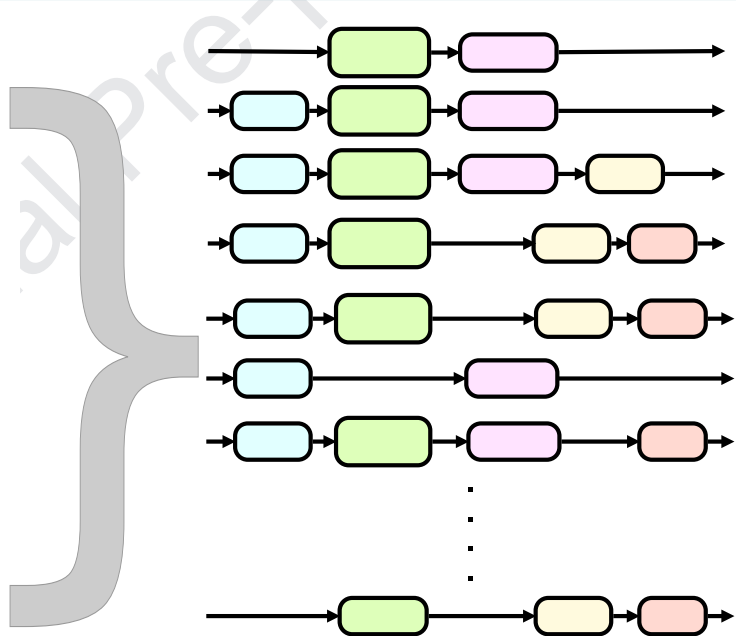
Figure 6: A: All 1618 models generated to predict the instantaneous modulus of articular cartilage presented in terms of model complexity (i.e., the number of latent variables) and test set performance. Color of individual models indicates whether the data was treated with scatter correction. Margin plots indicate the distribution of models in terms of complexity and performance. B: True vs predicted values for the test set between no preprocessing and the best performing pipeline. C: All preprocessing pipelines (scatter correction vs no scatter correction) sorted according to the $R^2$ values of the test set. D: Residual plots and the distributions of residuals between no preprocessing and the best performing pipeline.
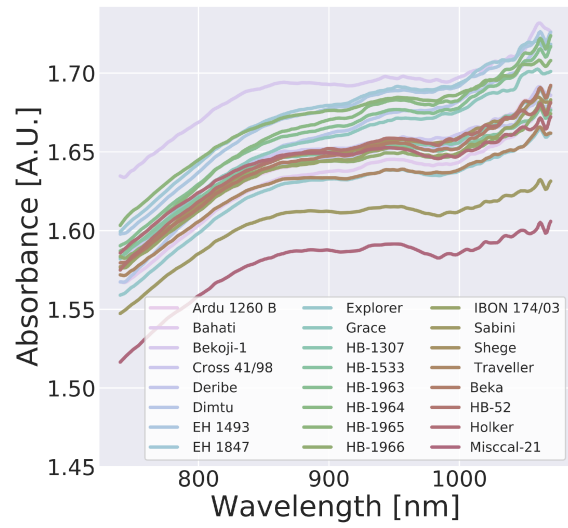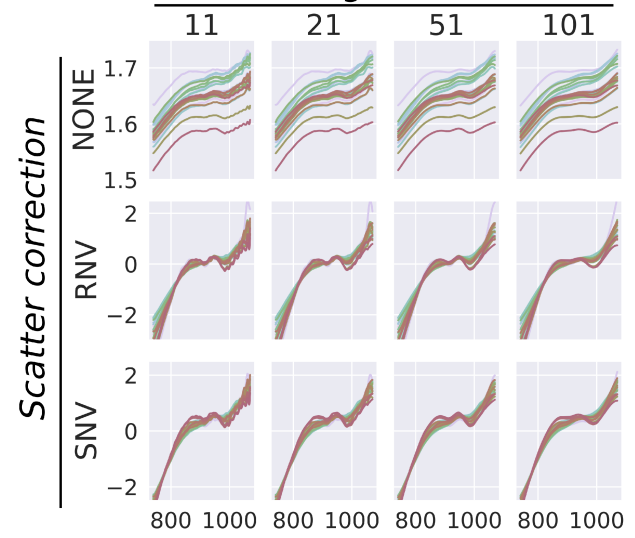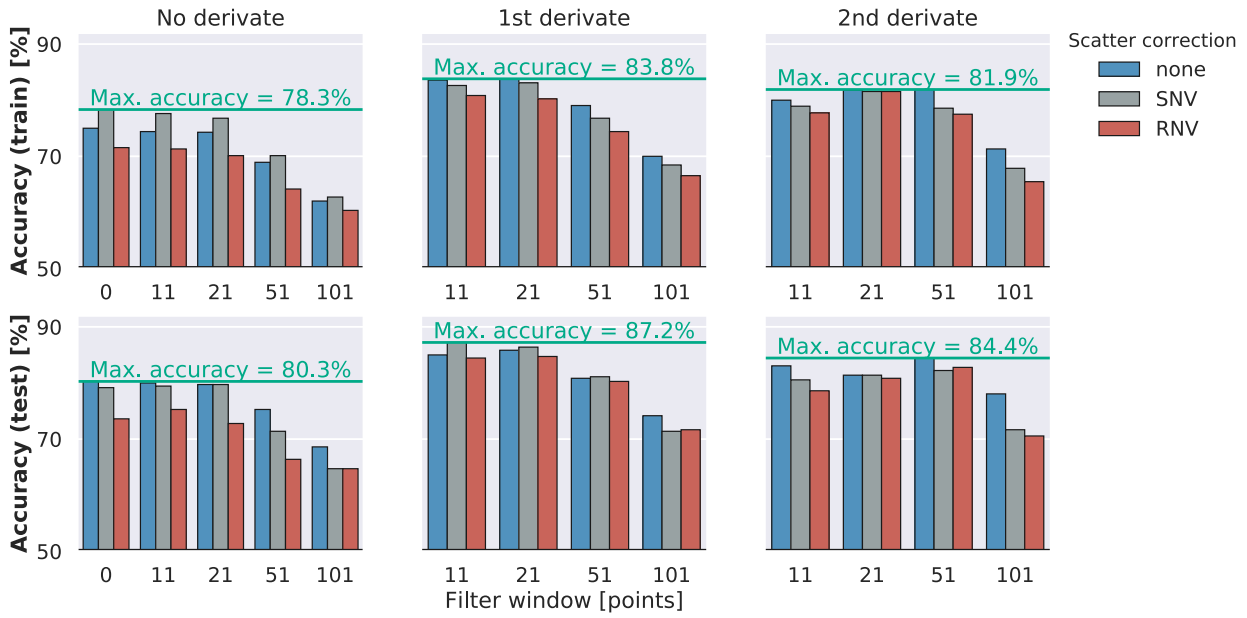
NIR Spectra → **Preprocessing** → **Feature extraction** → **Classification / Regression** → **Predicted value**

**Clipping**

**Standard normal variate**

**Robust normal variate**

**Local standard normal variate**

**Multiple scatter correction**

**Extended multiple scatter correction**

**Normalization**

**Baseline**

**Convolution**

**Savitzky-Golay**

**Derivatives**

**Detrending**

**Trimming**

**Resampling**

nippy automatically generates comparable preprocessing pipelines consisting of different combinations of preprocessing operations and/or parameters. Desired operations and parameter ranges are defined as configuration files.
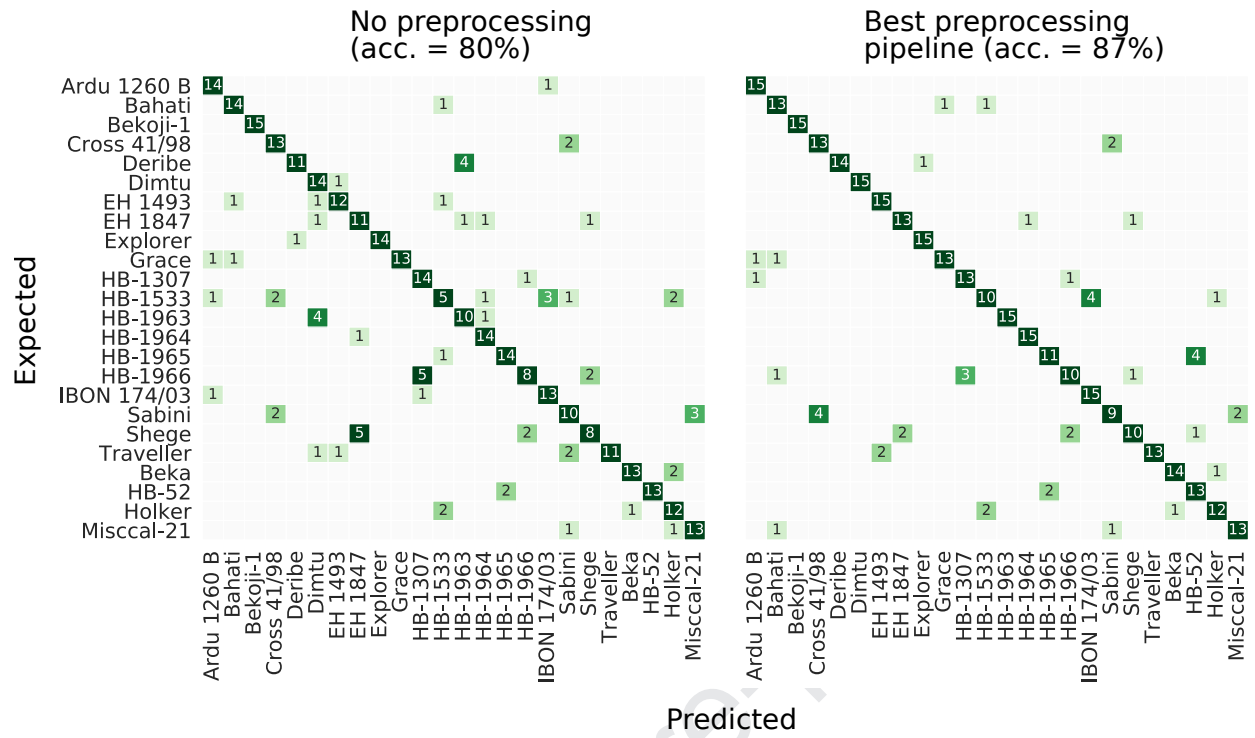
## Configuration file:

```
[SAVGOL]
    filter_win = 7, 11, 13
    poly_order = 3
    deriv_order = 0, 1, 2
    also_skip = True
[SNV]
    also_skip = True
[RNV]
    iqr = 75,25
[LSNV]
    also_skip = True
[TRIM]
    bins = 1000-1900
[CLIP]
    threshold = 1e4
    substitute = None
[RESAMPLE]
    resampling_ratio = .8
    also_skip = True
```
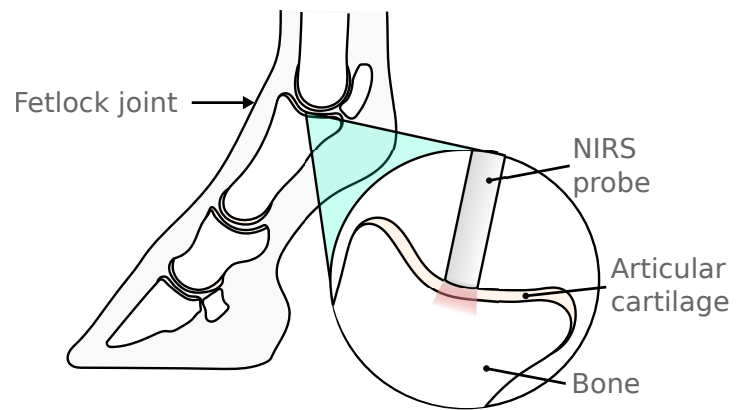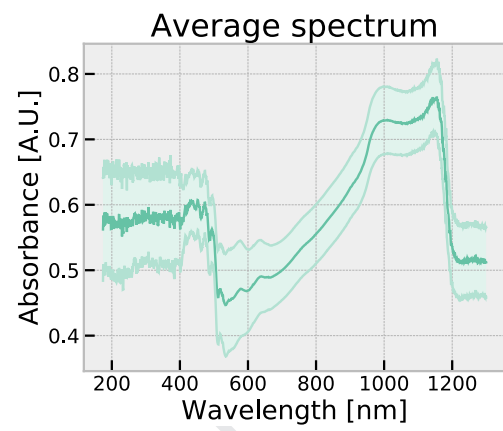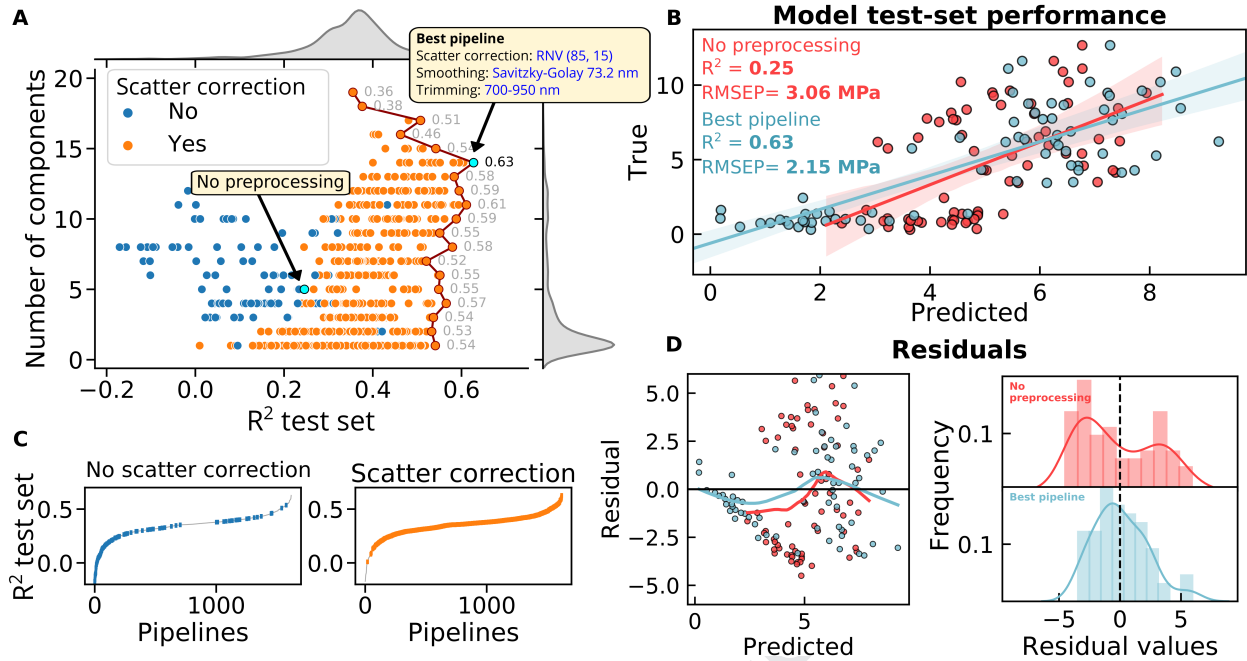
No preprocessing
(acc. = 80%)

Best preprocessing
pipeline (acc. = 87%)

**A**



**B**

**A**



Scatter correction
- No
- Yes

Number of components

Best pipeline
Scatter correction: RNV (85, 15)
Smoothing: Savitzky-Golay 73.2 nm
Trimming: 700-950 nm

No preprocessing

0.36
0.38
0.51
0.46
0.55
0.63
0.58
0.59
0.61
0.59
0.55
0.58
0.52
0.55
0.55
0.57
0.54
0.53
0.54

$R^2$ test set

**B** **Model test-set performance**

No preprocessing
$R^2$ = **0.25**
RMSEP= **3.06 MPa**

Best pipeline
$R^2$ = **0.63**
RMSEP= **2.15 MPa**

True

Predicted

**C**

$R^2$ test set

No scatter correction

Scatter correction

Pipelines

Pipelines

**D** **Residuals**

Residual

Predicted

No preprocessing

Best pipeline

Frequency

Residual values

# Highlights

- Spectral preprocessing affects the performance of chemometric models
- Selection of best preprocessing strategy depends heavily on the intended application
- nippy module can rapidly explore different preprocessing combinations
- Semi-automatic preprocessing tuning for spectroscopic models

**Declaration of interests**

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

On behalf of all the authors;

JARI TORNIAINEN