

Backtransformation: A new representation of data processing chains with a scalar decision function

blinded author names

Received: 07 September 2014, revised: 11 July 2015, second revision: 28 September 2015

Abstract Data processing often transforms a complex signal using a set of different preprocessing algorithms to a single value as the outcome of a final decision function. Still, it is challenging to understand and visualize the interplay between the algorithms performing this transformation. Especially when dimensionality reduction is used, the original data structure (e.g., spatio-temporal information) is hidden from subsequent algorithms. To tackle this problem, we introduce the backtransformation concept suggesting to look at the combination of algorithms as *one* transformation which maps the original input signal to a single value. Therefore, it takes the derivative of the final decision function and transforms it back through the previous processing steps via backward iteration and the chain rule. The resulting derivative of the composed decision function in the sample of interest represents the complete decision process. Using it for visualizations might improve the understanding of the process. Often, it is possible to construct a feasible processing chain with affine mappings which simplifies the calculation for the backtransformation and the interpretation of the result a lot. In this case, the affine backtransformation provides the complete parameterization of the processing chain. This article introduces the theory, provides implementation guidelines, and presents three application examples.

Keywords affine transformations · function composition · processing chain interpretation · processing chain visualization

Mathematics Subject Classification (2000) 68T30 · 68N99 · 68W40

1 Introduction

The basis of machine learning is understanding the data (Chen et al, 2008), and generating descriptive features (Domingos, 2012). Consequently, for numerous data types and processing algorithms, visualization approaches have been developed (Rieger et al, 2004; Rivet et al, 2009; Le et al, 2012; Haufe et al, 2014; Szegedy et al, 2014) as a better abstraction to enhance the understanding of the behaviour of the applied algorithms and of the data. Here, the visualization of an algorithm is often realized in a similar way as for the input data.

To come up with a representation gets way more complicated when algorithms are combined for a more sophisticated preprocessing before applying a final decision algorithm (Verhoeve and de Wulf, 1999; Rivet et al, 2009; Krell et al, 2013a; Kirchner et al, 2013; Feess et al, 2013), i.e., for processing chains. Under these circumstances, understanding and visualization of single algorithms does only explain single steps in the processing chain that are typically not independent from each other. The order of preprocessing algorithms, e.g., influences single processing visualizations, although the value of the final decision function might be not or only weakly influenced. Hence, one is often interested in knowledge about the whole data transformation in the processing chain but a general approach for solving this problem is missing. This situation gets even worse the more complex the data and the associated processing chains become. If dimensionality reduction algorithms are used for example to reduce the complexity of the data and to get rid of the noise, the structure of the output data is usually very different from the original input after the reduction step. In such a case, it is very difficult to understand the connection between decision algorithm, preprocessing, and original data even if single parts can be visualized. Consequently, a concept for representing the complete processing chain in the domain and format of the original input data is required.

Several approaches are described in the literature to visualize the outcome and transformation of classification algorithms, but again, taking the perspective of a single processing step neglecting the processing history (i.e., the preceding algorithms).

When using classifiers with kernels, a direct visualization of the classifier becomes impossible. Baehrens et al (2010) calculate the derivative of the classification function to give information of the classifier dependent on a chosen sample. Unfortunately, this calculation of the derivative is quite complex, difficult to automatize, computationally expensive, and does not consider any preprocessing before the classification. This makes it hard to apply and to generalize for complete data processing chains and high-dimensional data.

Blankertz et al (2011) discuss the visualization of the linear discriminant analysis (LDA) in the context of an electroencephalogram (EEG) based brain-computer interface (BCI) application with different views on the temporal, spatial and spatio-temporal domain. Here, the classifier is applied on spatial features and visualized as a spatial filter together with an interpretation in

46 relation to the original data and other spatial filters. For other visualizations,
47 the classifier weights are not directly used. Furthermore, no complex signal
48 processing chain is used, even though spatial filters are very common for the
49 preprocessing of this type of data. The LDA was applied to the raw data and
50 largely improved with a shrinkage criterion. As a side remark, they mention
51 the possibility to visualize the LDA weights directly, when applied to spatio-
52 temporal features (Blankertz et al, 2011, paragraph before section 6, p. 18).

53 This direct visualization of weights of a linear support vector machine
54 (SVM) has already been suggested by LaConte et al (2005).¹ This approach
55 is intuitive, easy to calculate, and enables a combination with the prepro-
56 cessing. Furthermore, it can be generalized to other data and other classi-
57 fiers (Blankertz et al, 2011).

58 This paper introduces our solution approach denoted as *backtransforma-*
59 *tion*. It incorporates the aforementioned approaches, but with the fundamental
60 difference that it takes all preprocessing steps in the respective chain into ac-
61 count. With this approach, we are able to extract the complete transformation
62 of the data from the chain, so that, e.g., changes in the order of algorithms or
63 the effect of insertions/deletions of single algorithms become immediately vis-
64 ible. Backtransformation also considers processing chains, where the original
65 (e.g., spatio-temporal) structure of the data is hidden. The data processing
66 chain is identified with a (composed) function, mapping the input data to a
67 scalar. In its core, backtransformation is the derivative of this function, cal-
68 culated with the chain rule or numerically. The derivative is either calculated
69 locally for each sample of interest (*general backtransformation*) or globally
70 when the processing chain consists of affine transformations only (affine back-
71 transformation). While the general backtransformation gives information on
72 which components in the data have a large (local) influence on the decision
73 process and which components are rather unimportant, the affine backtrans-
74 formation is independent from the single sample².

75 Numerous established data processing algorithms are affine transforma-
76 tions and it is often possible to combine them to process the data. Hence, we
77 also take a closer look at this type of algorithms and we show that it is possible
78 to retrieve the information on how the data is transformed by the complete
79 decision process, even if a dimensionality reduction algorithm or a temporal
80 filter hide information. The *affine backtransformation* iteratively goes back
81 from the decision algorithm through all processing steps to determine a pa-
82 rameterization of the composed processing function and to enable a semantic
83 interpretation. This results in a helpful representation of the processing chain,
84 where each component in the source domain of the data gets a weight as-
85 signed showing its impact in the decision process. In fact, summing up the
86 products of weights and respective data parts is equivalent to applying the
87 single algorithms on the data step-by-step.

¹ Further methods are presented but they are tailored to functional magnetic resonance imaging (fMRI) data.

² The respective derivatives are constant for every sample and as such not depending on it.

88 In Section 2, the backtransformation concept is introduced. First, we intro-
 89 duce the general backtransformation for differentiable processing chains. This
 90 is followed by the special variant which is obtained when working with affine
 91 transformations. To be even more specific, we discuss the backtransformation
 92 at a processing scheme for segmented time series data in Section 2.3. Here,
 93 we give examples of algorithms for the affine backtransformation, the generic
 94 backtransformation, and also mention cases where it is not applicable. In Sec-
 95 tion 2.4, we describe how the backtransformation is implemented in a generic
 96 way. This is followed by applications of the backtransformation in Section 3
 97 Finally, a conclusion is made in Section 4.

98 2 Methods

99 The requirement to apply the proposed backtransformation as outlined in
 100 the following is that the data processing is a concatenation of differentiable
 101 transformations (e.g., affine mappings) and that the last algorithm in the chain
 102 is a (decision) function which maps the data to a *single scalar*. The final
 103 mapping to the label in case of a classification task is not relevant, here.

104 For each processing stage, the key steps of the backtransformation are to
 105 first choose a mathematical representation of input and output data and then
 106 to determine a parameterization of the algorithm which has to be mapped to
 107 fit to the chosen data representations. Finally, the derivatives of the resulting
 108 transformations have to be calculated and iteratively combined. In its core it is
 109 the application of the chain rule for derivatives (see Section 2.1). For the case
 110 of using only affine mappings, it is just the multiplication of the transformation
 111 matrices, as shown in Section 2.2. Details on the implementation are given in
 112 Section 2.4. For an example of a processing chain of windowed time series data
 113 with a two-dimensional representation of the data see Fig. 1 and Section 2.3.

114 The backward modeling begins with the parametrization of the final deci-
 115 sion function and continues by iteratively combining it backwards with the
 116 preceding algorithms in a processing chain. With each iteration, weights are
 117 calculated, which correspond to the components of the input data of the last
 118 observed algorithm.

119 For the abstract formulation of the backtransformation approach, data
 120 with a one-dimensional representation before and after each processing step
 121 is used. The output of each processing step is fed into the next processing
 122 algorithm.

123 2.1 Backtransformation using the Derivative

This section introduces the general backtransformation. Let the input data be
 denoted with $x^{(0)} = x^{\text{in}} \in \mathbb{R}^{n_0}$ and let the series of processing algorithms be
 represented by differentiable mappings

$$F_0 : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}, \dots, F_k : \mathbb{R}^{n_k} \rightarrow \mathbb{R} \quad (1)$$

which are applied to the data consecutively.³ Then, the application of the processing chain to obtain the output data $x^{\text{out}} := x^{(k+1)}$ from the input data $x^{(0)}$ can be summarized to:

$$x^{\text{out}} = x^{(k+1)} = F(x^{(0)}) = (F_k \circ \dots \circ F_0)(x^{(0)}) . \quad (2)$$

With this notation, the derivative can be calculated with the chain rule:

$$\frac{\partial F}{\partial y} (x^{(0)}) = \frac{\partial F_0}{\partial y^{(0)}} (x^{(0)}) \cdot \frac{\partial F_1}{\partial y^{(1)}} (x^{(1)}) \cdot \dots \cdot \frac{\partial F_{k-1}}{\partial y^{(k-1)}} (x^{(k-1)}) \cdot \frac{\partial F_k}{\partial y^{(k)}} (x^{(k)}) , \quad (3)$$

where $x^{(l)} \in \mathbb{R}^{n_l}$ is the respective input of the l -th algorithm in the processing chain with the mapping F_l and $x^{(l+1)}$ is the output. The terms $\frac{\partial F_l}{\partial y^{(l)}}$ and $\frac{\partial F}{\partial y}$ represent the total differentials of the differentiable mappings and not the partial derivatives. Equation (3) is a matrix product. It can be calculated iteratively using the backtransformation matrices B_l and the derivatives $\frac{\partial F_{l-1}}{\partial y^{(l-1)}} (x^{(l-1)})$:

$$B_k = \frac{\partial F_k}{\partial y^{(k)}} (x^{(k)}) \quad \text{and} \quad B_{l-1} = \frac{\partial F_{l-1}}{\partial y^{(l-1)}} (x^{(l-1)}) \cdot B_l \quad \text{with} \quad l = 1, \dots, k . \quad (4)$$

124 Now each matrix $B_l \in \mathbb{R}^{n_l \times 1}$ has the same dimensions⁴ as the respective
 125 $x^{(l)}$ and tells which change in the components of $x^{(l)}$ will increase (positive
 126 entry in B_l), decrease (negative entry), or will have no effect (zero entry) on
 127 the decision function. The higher the absolute value of an entry (multiplied
 128 with the estimated variance of the respective input), the larger is the influence
 129 of the respective data component on the decision function. Consequently, not
 130 only the backtransformation of the complete processing chain (B_0) but also
 131 the intermediate results (B_l ; $l > 0$) might be used for analyzing the processing
 132 chain. B_k is the matrix used in the existing approaches, which do not consider
 133 the preprocessing (LaConte et al, 2005; Baehrens et al, 2010; Blankertz et al,
 134 2011). Note that the B_l are dependent on the input of the processing chain
 135 and are expected to change with changing input. So the information about
 136 the influence of certain parts in the data is only a *local* information. A *global*
 137 representation is only possible when using affine transformations instead of
 138 arbitrary differentiable mappings F_l .

³ The notation of data and its components differs from the notation in classification tasks. Here, we look at one data sample $x^{(0)}$ with its different processing stages $x^{(l)}$ and the respective changes in each component of the data ($x_{gh}^{(l)}$). The double index notation is applied to account for different axes in the data as in time series (different sensors and time points) or images.

⁴ With $n_{k+1} := 1$ it holds that $\frac{\partial F_l}{\partial y^{(l)}} \in \mathbb{R}^{n_l \times n_{l+1}}$ and the dimensions of B_l are a consequence of the recursion. Another reason for the dimensions of B_l is that B_l corresponds to the mapping of $x^{(l)}$ to the scalar output x^{out} .

139 2.2 Affine Backtransformation

For handling affine transformations like translations, the data vectors are augmented by adding a coordinate with value 1 to have homogenous coordinates. Every affine transformation F can be identified with a tuple (A, T) , where A is a linear mapping matrix and T a translation vector and the corresponding mapping of the processing algorithm applied on data x^{in} reads as

$$x^{\text{out}} = F(x^{\text{in}}) = Ax^{\text{in}} + T = (A \ T) \begin{pmatrix} x^{\text{in}} \\ 1 \end{pmatrix}. \quad (5)$$

So by extending the matrix $(A \ T)$ to a Matrix A' with an additional row of zeros with a 1 at the translational component, the mapping on the augmented data $x'^{\text{in}} = \begin{pmatrix} x^{\text{in}} \\ 1 \end{pmatrix}$ can be written in the simple notation: $x'^{\text{out}} = A'x'^{\text{in}}$. With a processing chain with corresponding matrices A'_0, \dots, A'_k the transformation of the input data x'^{in} can be summarized to

$$x'^{\text{out}} = A'_k \cdot \dots \cdot A'_1 \cdot A'_0 \cdot x'^{\text{in}}. \quad (6)$$

With this notation, the backtransformation concept now boils down to iteratively determine the matrices

$$B_k = A'_k, B_{k-1} = A'_k \cdot A'_{k-1}, \dots, \text{ and } B_0 = A'_k \cdot A'_{k-1} \cdot \dots \cdot A'_1 \cdot A'_0. \quad (7)$$

140 This corresponds to a convolution of affine mappings.⁵ Each $B_l \in \mathbb{R}^{(n_l+1) \times 2}$
 141 defines the mapping of the data from the respective point in the processing
 142 chain (after l previous processing steps) to the final decision value. So *each*
 143 *product* B_l consists of a weighting vector $w^{(l)}$ and an offset $b^{(l)}$ and the artificial
 144 second row with zero entries and 1 in the last column. The term $w^{(l)}$ can now
 145 be used for interpretation and understanding the respective sub-processing
 146 chain or the complete chain with $w^{(0)}$ (see Section 3). The term is equivalent
 147 to the B_l from the backtransformation using the derivative (Section 2.1).

148 The following section renders possible (and impossible) algorithms which
 149 can be used for the affine backtransformation and how the weights from the
 150 backtransformation are determined in detail for a data processing chain ap-
 151 plied on two-dimensional data.

152 2.3 Backtransformation Modeling Example

153 In this section, a more concrete example of applying the backtransformation
 154 principle is given for processing time series epochs of fixed length of several
 155 sensors with the same sampling frequency. We provide examples for affine
 156 transformations to show that there is a large number of available algorithms
 157 to construct a good processing chain. Additionally, cases will be highlighted

⁵ Note that no matrix inversion is required even though one might expect that, because the goal is to find out what the original mapping was doing with the data which sounds like an inverse approach.

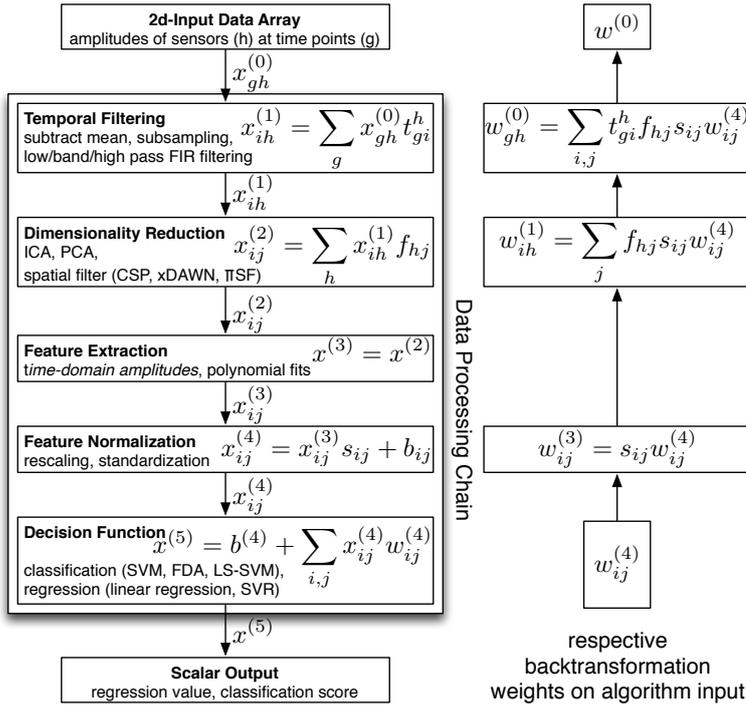


Fig. 1: Illustrative data processing chain scheme with examples of affine algorithms and the formulas for the backtransformation in short. Spatio-temporal data $x_{gh}^{(0)}$ are processed from top to bottom ($x^{(5)}$). Every component of the scheme is optional. Backtransformation takes the classifier parametrization $w^{(4)}$ and projects it iteratively back ($w^{(k)}$) through the processing chain and results in a representation $w^{(0)}$ corresponding to the input domain. For more details refer to Section 2.3.

158 which require the general backtransformation and cases where the backtrans-
 159 formation is not applicable.

160 A possible processing chain with examples of affine mappings and the re-
 161 spective backtransformation weights is depicted in Fig. 1. Note that all compo-
 162 nents of this chain are optional and the presented scheme can be applied to
 163 an arbitrary data processing chain even if dimensions like time and space are
 164 replaced by others or left out (see Sections 2.2 and 3.2).

165 An intuitive way of handling such data is to represent it as two-dimensional
 166 arrays with the time on one axis and space (e.g., sensors) on the other axis,
 167 since important preprocessing steps like temporal and spatial filters just oper-
 168 ate on one axis. So this type of representation eases the use and the parame-
 169 terization of these algorithms compared to the aforementioned mathematically
 170 equivalent one-dimensional representation. Furthermore, a two-dimensional

171 representation of the data helps for its visualization and interpretation. For
 172 parametrization of the two-dimensional arrays, the common double index no-
 173 tation is used, where the data $x^{(0)}$ is represented by its components $x_{gh}^{(0)}$ with
 174 temporal index g and spatial index h . This index scheme will be kept for all
 175 processing stages even if the data could be represented as one-dimensional fea-
 176 ture vectors for some stages. The same indexing scheme can be applied for the
 177 parametrization of the affine data processing algorithms in the chain as will be
 178 shown in the following. The input of the i -th algorithm is denoted with $x^{(i-1)}$
 179 and the output with $x^{(i)}$ respectively. To fit to the concept of backtransfor-
 180 mation, first the parametrization of the decision algorithm will be introduced
 181 and then the preceding algorithms step-by-step. An overview of the process-
 182 ing chain, the chosen parameterizations, and the resulting weights from the
 183 backtransformation is depicted in Fig. 1.

Scalar Decision Function A linear decision function can be parameterized using a decision vector/matrix $w_{ij}^{(4)} \in \mathbb{R}^{m_i \times n_j}$ and an offset $b^{(4)} \in \mathbb{R}$. The transformation of the input $x^{(4)} \in \mathbb{R}^{m_i \times n_j}$ to the decision value $x^{(5)} \in \mathbb{R}$ is then defined as

$$x^{(5)} = b^{(4)} + \sum_{i=1}^{m_i} \sum_{j=1}^{n_j} x_{ij}^{(4)} w_{ij}^{(4)}, \quad (8)$$

184 with m_i time points and n_j sensors. Some examples for machine learning al-
 185 gorithms with linear decision function are SVMs (Vapnik, 1995; Steinwart
 186 and Christmann, 2008; Chang and Lin, 2011), balanced relative margin ma-
 187 chines (Krell et al, 2014a), regularized fishers discriminant analysis classi-
 188 fiers (Mika et al, 2001), passive-aggressive perceptrons (Crammer et al, 2006),
 189 linear regression, support vector regression (Smola and Schölkopf, 2004), ridge
 190 regression, and one-class SVMs (Schölkopf et al, 2001; Krell and Wöhrle, 2015)
 191 and there are many more.

192 Depending on the application, data might be not linearly separable or a
 193 nonlinear separation provides better results. Here, a common approach is to
 194 use nonlinear kernels instead of the linear function. All common kernels are
 195 differentiable, so here the general backtransformation can be still applied in-
 196 stead of the affine backtransformation. As long as the decision function is
 197 differentiable, the general backtransformation can be used, too. When combin-
 198 ing (linear or differentiable) classifiers as an ensemble it depends on the
 199 final gating function, if the resulting scalar comes from an affine/differentiable
 200 function⁶. The same holds for neural networks where different transition func-
 201 tions could be used. Unfortunately, for neural networks the derivative might
 202 not improve the understanding especially when it is showing unexpected local
 203 behavior as explained by Szegedy et al (2014). Nevertheless, most often these
 204 methods are differentiable. If there is no strict step function used but the func-
 205 tion is locally Lipschitz or even locally linear the approximation of a derivative

⁶ A weighted sum of classifiers preserves linearity/differentiability. A majority vote will result in a non-differentiable classifier but when the score is the sum of the voters for the selected class, the resulting function will still be locally linear/differentiable.

could be still used even though some information in the critical points might be hidden. Furthermore, in these cases it is better to use a derivative, which considers the left and the right side for each component.

If there is no scalar output or the function is locally constant it is not possible to derive information from the backtransformation. A decision tree usually produces no useful output function. If the output of a classifier is only ± 1 , no information can be obtained. Another example for a locally constant function could be obtained from a linear decision function $f(x)$ by limiting its values to the interval $[-1, 1]$ with $\min\{1, \max\{-1, f(x)\}\}$. For every x with $|f(x)| > 1$ the resulting new decision function is locally constant and no interpretation of the derivative is possible.

Feature Normalization With a scaling $s \in \mathbb{R}^{m_i \times n_j}$ and transition $b \in \mathbb{R}^{m_i \times n_j}$ and the same indexes as for the linear decision function, an affine feature normalization can be written as

$$x_{ij}^{(4)} = x_{ij}^{(3)} s_{ij} + b_{ij} \text{ with } i \in \{1, \dots, m_i\} \text{ and } j \in \{1, \dots, n_j\} . \quad (9)$$

This covers most standard feature normalization algorithms like rescaling or standardization (Aksoy and Haralick, 2001). Nonlinear scalings, e.g., using absolute values as in $\min\{10, |x_{ij}^{(3)}|\}$, or sample dependent scalings, e.g., division by the Euclidean norm $s_{ij} = \frac{1}{\|x^{(3)}\|_2}$, are not affine mappings and could not be used for the affine backtransformation.

The general backtransformation could still be used for differentiable normalizations like Euclidean normalization if $x^{(3)} \neq 0$. Using min or max results in locally constant behavior which restricts the applicability of the backtransformation.

For the affine backtransformation, the formula of the feature normalization need to be inserted into the formula of the decision function:

$$x^{(5)} = b^{(4)} + \sum_{i,j} \left(x_{ij}^{(3)} s_{ij} + b_{ij} \right) w_{ij}^{(4)} = b^{(3)} + \sum_{i,j} x_{ij}^{(3)} s_{ij} w_{ij}^{(4)} . \quad (10)$$

Here, $b^{(3)} = b^{(4)} + \sum_{i,j} b_{ij} w_{ij}^{(4)}$ summarizes the offset. As denoted in Fig. 1, $s_{ij} w_{ij}^{(4)}$ is the weight to the input data part $x_{ij}^{(3)}$.

Feature Generation For simplicity, the data amplitudes at different sensors have been directly taken as features and nothing needs to be changed in this step ($x^{(3)} = x^{(2)}$). Other linear features like polynomial fits would be possible, too (Straube and Feess, 2013). Nonlinear features (e.g., standard deviation, sum of squares, or sum of absolute values of each sensor) would not work for the affine backtransformation but for the general one. Symbolic features, mapped to natural numbers will be even impossible to analyze with the general backtransformation.

236 *Dimensionality Reduction on the Spatial Component* A spatial filter trans-
 237 forms real sensors to new pseudo sensors by linear combination of the signal
 238 of the original sensors. To use well known dimensionality reduction algorithms
 239 like principal component analysis (Lagerlund et al, 1997; Rivet et al, 2009;
 240 Abdi and Williams, 2010, PCA), and independent component analysis (Jut-
 241 ten and Herault, 1991; Rivet et al, 2009, ICA) for spatial filtering, the space
 242 component of the data is taken as feature component for these algorithms and
 243 the time component for the samples. Examples for typical spatial filters are
 244 common spatial patterns (Blankertz et al, 2008, CSP), xDAWN (Rivet et al,
 245 2009; Woehrle et al, 2015), and π SF (Ghaderi and Straube, 2013).

246 The backtransformation with the spatial filtering is the most important
 247 part of the concept, because spatial filtering hides the spatial information
 248 needed for visualization or getting true spatial information into the classifier.

The number of virtual sensors ranges between the number of real sensors
 and one. The spatial filter for the j -th virtual sensor is a tuple of weights
 $f_{1j}, \dots, f_{n_h j}$ defining the linear weighting of the n_h real sensors. The transfor-
 mation for the i -th time point is written as

$$x_{ij}^{(2)} = \sum_{h=1}^{n_h} x_{ih}^{(1)} f_{hj}, \quad (11)$$

where the time component could be ignored, because the transformation is
 independent of time. The transformation formula can be substituted into for-
 mula (11):

$$x^{(5)} = b^{(3)} + \sum_{i,j} \sum_{h=1}^{n_h} x_{ih}^{(1)} f_{hj} s_{ij} w_{ij}^{(4)} \quad (12)$$

$$= b^{(3)} + \sum_{i,h} x_{ih}^{(1)} \cdot \left(\sum_j f_{hj} s_{ij} w_{ij}^{(4)} \right). \quad (13)$$

249 Equation (13) shows, that the weight $\sum_j f_{hj} s_{ij} w_{ij}^{(4)}$ is assigned to the input
 250 data component $x_{ih}^{(1)}$. If there is no time component, a spatial filter is just a
 251 linear dimensionality reduction algorithm. It is also possible to combine dif-
 252 ferent reduction methods or to do a dimensionality reduction after the feature
 253 generation.

254 For spatial filtering, linear transformations are the common choice. But
 255 for more general dimensionality reduction algorithms like the PCA, it is also
 256 possible to use kernels. Since kernels are usually differentiable, it would be still
 257 possible to apply the generic backtransformation, when such an algorithm is
 258 used in the processing chain.

259 *Detrending, Temporal Filtering, and Decimation* There are numerous discrete-
 260 time signal processing algorithms (Oppenheim and Schaffer, 2009). Detrending
 261 the mean from a time series can be done in several ways. Having a time win-
 262 dow, a direct approach would be to subtract the mean of the time window,

263 or to use some time before the relevant time frame to calculate a guess for
 264 the mean (baseline correction). Often, such algorithms can be seen as finite
 265 impulse response (FIR) filters, which eliminate very low frequencies. Filtering
 266 the variance is a quadratic filter (Krell et al, 2013b) and infinite impulse re-
 267 sponse (IIR) filters have a feedback part. Both filters are not applicable for
 268 the backtransformation, because they have no respective affine transforma-
 269 tions and because they rely on the complete signal which makes it impossible
 270 to obtain a local derivative.

One can either use uniform temporal filtering, which is similar to spatial filtering with changed axis, or introduce different filters for every sensor. As parametrization, t_{gi}^h is chosen for the weight at sensor h for the source g and the resulting time point i with a number of m_g time points in the source domain:

$$x_{ih}^{(1)} = \sum_{g=1}^{m_g} x_{gh}^{(0)} t_{gi}^h. \quad (14)$$

Starting with the more common filter formulation as convolution (filter of length N):

$$x_{ih}^{(1)} = \sum_{l=0}^N a_l \cdot x_{(n-l)h}^{(0)} \stackrel{g:=n-l}{=} \sum_{g=n-N}^n a_{(n-g)} \cdot x_{gh}^{(0)}, \quad (15)$$

271 the filter coefficients a_i can be directly mapped to the t_{gi}^h and the other coef-
 272 ficients can be set to zero.

273 Reducing the sampling frequency of the data by downsampling is a combina-
 274 tion of a low-pass filter and systematically leaving out several time points
 275 after the filtering (decimation). When using a FIR filter, the given parame-
 276 terization of a temporal filter can be used here, too. For leaving out samples,
 277 the matrix t_{gi} for sensor h can be obtained from an identity matrix by only
 278 keeping the rows, where samples are taken from.

The final step is similar to the spatial filtering part:

$$x^{(5)} = b^{(3)} + \sum_{i,h} \left(\sum_{g=1}^{m_g} x_{gh}^{(0)} t_{gi}^h \right) \cdot \left(\sum_j f_{hj} s_{ij} w_{ij}^{(4)} \right) \quad (16)$$

$$= b^{(3)} + \sum_{g,h} x_{gh}^{(0)} \cdot \left(\sum_{i,j} t_{gi}^h f_{hj} s_{ij} w_{ij}^{(4)} \right) \quad (17)$$

$$= b^{(3)} + \sum_{g=1}^{m_g} \sum_{h=1}^{n_h} x_{gh}^{(0)} w_{gh}^{(0)}. \quad (18)$$

279 The input component of the original data $x_{gh}^{(0)}$ finally gets assigned the weight
 280 $w_{gh}^{(0)} = \sum_{i,j} t_{gi}^h f_{hj} s_{ij} w_{ij}^{(4)}$. Note that for some applications it is good to work
 281 on normalized and filtered data for interpreting data and the behavior of the
 282 data processing. In that case, the backtransformation is stopped before the
 283 temporal filtering and the respective weights are used.

284 *Others* The aforementioned algorithms can be combined and repeated (e.g.,
 285 concatenations of FIR filters or PCA and xDAWN). Having a different feature
 286 generator, multiple filters, decimation, or skipping a filter or normalization
 287 the same calculation scheme could be used resulting in different $b^{(3)}$ and $w^{(0)}$.
 288 Nevertheless, $w^{(0)}$ has the same indexes as the original data $x^{(0)}$. After the
 289 final mapping to a scalar by the decision function, a shift of the decision cri-
 290 terion (e.g., using threshold adaptation as suggested by Metzen and Kirchner
 291 (2011)) is possible but has no impact on the backtransformation because it
 292 only requires $w^{(0)}$ and not the offset. If a probability fit (Platt, 1999; Lin et al,
 293 2007; Baehrens et al, 2010) was used, this step has to be either ignored or the
 294 general approach (Section 2.1) has to be applied. Since the probability fit is a
 295 mostly sigmoid function which maps $\mathbb{R} \rightarrow [0, 1]$, it is also possible to visualize
 296 its derivative separately. For the interpretation concerning a sample, the func-
 297 tion value is determined and the respective (positive) derivative is multiplied
 298 with the affine transformation part to get the local importance. Hence, the
 299 relations between the weights remain the same but the absolute values only
 300 change. This approach of mixing the calculations is much easier to implement
 301 and interpret.

302 If nonlinear preprocessing is used to normalize the data (e.g., to have vari-
 303 ance of one), the normalized data can be used as input for the backtransfor-
 304 mation and the respective processing chain. This might be even advantageous
 305 for the interpretation when the visualization of the original data is not help-
 306 ful due to artifacts and outliers. An example for such a case is to work with
 307 normalized image data like the MNIST dataset (LeCun et al, 1998) instead of
 308 the original data, where the size of the images and the position of the digits
 309 varied a lot (see also Section 3.2 and Section 3.3).

310 If any of the algorithms in the observed processing chain is not an affine
 311 mapping, the affine backtransformation cannot be applied. For getting the real
 312 derivatives for the general backtransformation all algorithms need to be differ-
 313 entiable. But if the derivative vanishes at some points due to locally constant
 314 behavior, the backtransformation might be meaningless. On the other hand,
 315 if a generalized derivative can be determined for non-differentiable algorithms
 316 this might still work (Clarke, 1990; Rockafellar and Wets, 2009).

317 2.4 Generic Implementation of the Backtransformation

318 This section gives information on how to apply the backtransformation concept
 319 in practice especially when the aforementioned calculations are difficult or
 320 impossible to perform and a “generic” implementation is required to handle
 321 arbitrary processing chains.

322 The backtransformation has been implemented in a signal processing and
 323 classification environment called pySPACE (Krell et al, 2013a) and can be
 324 directly used⁷. This modular Python software gives simple access to more than

⁷ <http://pyspace.github.io/pyspace/>

200 classification and preprocessing algorithms and so it provides a reasonable
 325 interface for a generic implementation. It provides data visualization tools for
 326 the different processing stages and largely supports the handling of complex
 327 processing chains.
 328

In practice, accessing the single parameterizations for the transformation matrices A_i for the affine backtransformation might be impossible (e.g., because external libraries are used without access to the internal algorithm parameters) or too difficult (e.g., code of numerous algorithms needs to be written to extract these parameters). In this case, the backtransformation approach cannot be applied directly in the way it is described in Section 2.2. Instead, the respective products and weights for the affine backtransformation can be reconstructed with the following trick which only requires the algorithms to be affine. No access to any parameters is needed. First, the offset of the transformation product is obtained by processing a zero data sample with the complete processing chain. The processing function is denoted by F . The resulting scalar output is the offset

$$b^{(0)} = F(0). \quad (19)$$

Second, a basis $\{e_1, \dots, e_n\}$ of the original space (e.g., the canonical basis) needs to be chosen. In the last step, the weights $w_i^{(0)}$, which directly correspond to the base elements, are determined by also processing the respective base element e_i with the processing chain and subtracting the offset $b^{(0)}$ from the scalar output:

$$w_i^{(0)} = F(e_i) - F(0). \quad (20)$$

The calculation of the derivative for the general backtransformation approach is more complicated. Deriving and implementing the derivative function for each algorithm used in a processing chain and combining the derivatives can be very difficult, especially if the goal is to implement it for a large number of relevant algorithms, e.g., as provided in the pySPACE framework. A generic approach would be to use automatic differentiation tools (Griewank and Walther, 2008). These tools generate a program which calculates the derivative directly from the program code. They can also consider the concatenation of algorithms by applying the chain rule. For most standard implementations, open source automatic differentiation tools could be applied. For existing frameworks, it is required to modify each algorithm implementation such that the existing differentiation tools know all derivatives of used elemental functions used in the code, which might be a lot of work. Furthermore, this approach would be impossible if black box algorithms were used. So for simplicity, a different approach, which is similar to the previous one for the affine case can be chosen. This is the numerical calculation of the derivative of the complete decision function via differential quotients for directional derivatives:

$$\frac{\partial F}{\partial e_i}(x_0) \approx \frac{F(x_0 + he_i) - F(x_0)}{h}. \quad (21)$$

Here, e_i is the i -th unit vector, and h is the step size. It is difficult to choose the optimal h for the best approximation, but for the backtransformation

a rough approximation should be sufficient. A good first guess is to choose $h = 1.5 \cdot 10^{-8} \langle x_0, e_i \rangle$ if $\langle x_0, e_i \rangle \neq 0$ and in the other case $h = 1.5 \cdot 10^{-8}$ (Press, 2007). In the backtransformation implementation in pySPACE, the value of $1.5 \cdot 10^{-8}$ can be exchanged easily by the user. It is additionally possible to use more accurate formulas for the differential quotient at the cost of additional function evaluations like

$$\frac{\partial F}{\partial e_i}(x_0) \approx \frac{F(x_0 - he_i) - 8F(x_0 - \frac{h}{2}e_i) + 8F(x_0 + \frac{h}{2}e_i) - F(x_0 + he_i)}{6h}. \quad (22)$$

3 Applications

Having a transformation of the decision algorithm back through different data representation spaces to the original data space might help for the understanding and interpretation of processing chains in several applications (e.g., image detection, classification of neuroscientific data, robot sensor regression) as explained in the following. First, some general remarks will be given on visualization techniques. Afterwards, the affine and the general backtransformation will be applied on handwritten digit classification (Section 3.2 and Section 3.3) because it is a relatively simple problem which can be understood without expert knowledge. Finally, a more complex example is given on EEG data classification (Section 3.4) and an outlook for further applications (Section 3.5).

3.1 Visualization in General

As suggested by LaConte et al (2005) for fMRI data, the backtransformation weights could be visualized in the same way as the respective input data is visualized. This works only if there is a possibility to visualize the data and if this visualization displays the “strength” of the values of the input data. Otherwise, additional effort has to be put into the visualization, or the weights have to be analyzed as raw numbers. For interpreting the weights, it is usually required to also have the original data visualized for comparison (as averaged data or single samples) because higher weights in the backtransformation could be rendered meaningless if the corresponding absolute data values are low or even zero. Additionally to the backtransformation visualization of *one* data processing chain, different chains (with different hyperparameters, training data, or algorithms) can be compared (Krell et al, 2014b). Differences in the weights directly correspond to the differences in the processing. Normally, weights with high absolute values correspond to important components for the processing and weights close to zero are less important and might be even omitted. This very general interpretation scheme does not work for all applications. In some cases, the weights have to be set in relation to the values of the respective data components: If data values are close to zero, high weights

might still be irrelevant, and vice versa. To avoid such problems, it is better to take normalized data, which is very often also a good choice for pure data visualization. Another variant to partially compensate for this issue is to also look at the products of weights and the respective data values.

According to Haufe et al (2014), the backtransformation model is a *backward model* of the original data and as such mixes the reduction of noise with the emphasis of the relevant data pattern. To derive the respective *forward model* they suggest to multiply the respective weighting vector with the covariance matrix of the data. From a different perspective, this approach sounds reasonable, too: If backtransformation reveals that a feature gets a very high weight by the processing chain, but this feature is zero for all except one outlier sample a modified backtransformation would reveal this effect. Furthermore, if a feature is highly correlated with other features, a sparse classifier might just use this one feature and skip the other features which might lead to the wrong assumption, that the other features are useless even though they provide the same information. On the other hand, if features are highly correlated as it holds for EEG data this approach might be also disadvantageous. The processing chain might give a very high weight to the feature, where the best distinction is possible, but the covariance transformation will blur this important information over all sensors and time points. Using such a blurred version for feature selection would be a bad choice. Another current drawback of the method by Haufe et al (2014) is that it puts some assumptions on the data which often do not hold: The expectancy values of noise, data, and signal of interest are assumed to be zero “w.l.o.g.” (without loss of generality). Hence, more realistic assumptions are necessary for better applicability. The effect of the covariance correction by Haufe et al (2014) will be analyzed in Sections 3.2 and 3.3.

Note that in Fig. 1, Section 2.2, and Section 2.3 it has been shown that every iteration step in the backtransformation results in weightings $w^{(i)}$ which correspond to the data $x^{(i)}$. This data is obtained by applying the first i algorithms of the processing chain on the original input data $x^{(0)}$. So depending on the application, it is even possible to visualize data and weights of intermediate processing steps. This can be used to further improve the overall picture of what happens in the processing chain.

3.2 Processing Chain Visualization:

Handwritten Digit Classification: Affine Processing Chain

For a simple application example of the backtransformation approach, the publicly available MNIST dataset is used (LeCun et al, 1998). This dataset contains numerous normalized greyscale images of all digits with a size of 28×28 pixels. They are stored as one-dimensional feature vectors (784 features). For processing, we first applied a PCA on the feature vectors and reduced the dimension of the data to four (or 64). As a second step, the resulting features were normalized to have zero mean and standard deviation of one on the training

403 data. Finally, a linear SVM (Chang and Lin, 2011) with a fixed regulariza-
 404 tion parameter of one is trained on the normalized PCA features. Without
 405 backtransformation, the filter weights for the 4 (or 64) principal components
 406 could be visualized in the domain of the original data and the single weights
 407 assigned by the SVM could be given, but the interplay between SVM and
 408 PCA would remain unknown, especially if all 784 principal components would
 409 be used. This information can only be given with backtransformation and is
 410 displayed in Fig. 2 for the distinction of digit pairs (from 0, 1, and 2). The
 411 generic implementation of the affine backtransformation was used, since only
 412 affine algorithms were used in the processing chain (PCA, feature standardiza-
 413 tion, linear classifier). The forward model to the backtransformation, obtained
 414 by multiplication with the covariance matrix, is also visualized in Fig. 2. Note
 415 that the original data is not normalized (zero mean), although this was an
 416 assumption on the data for the covariance transformation approach by Haufe
 417 et al (2014)⁸.

418 Generally, it can be seen that the classifier focuses on the digit parts, where
 419 there is no overlay between the digits on average. For one class there are high
 420 positive values and for the other there are high negative weights. For the classi-
 421 fication with 64 principal components, the covariance correction smoothes the
 422 weight usage and results in a visualization which is similar to the visualization
 423 of the backtransformation for the classification with 4 principal components.
 424 Hence, the 60 additional components are mainly used for canceling out “noise”.

425 3.3 Processing Chain Visualization: 426 Handwritten Digit Classification: Nonlinear Classifier

427 To show the effect of the generic backtransformation for a nonlinear processing
 428 chain, the evaluation of Section 3.2 is repeated with a radial basis function
 429 kernel for the SVM instead of a linear one. The hyperparameter of the kernel,
 430 γ , has been determined according to Varewyck and Martens (2011). Everything
 431 else remained unchanged. Again the generic implementation was used. Note
 432 that every sample requires its own backtransformation. So for the visualization
 433 of the backtransformation, only the first four single samples were taken.

434 It can be clearly seen in Fig. 3 that there is a different backtransformation
 435 for each sample. Similar to the results in Section 3.2 (Fig. 2), the backtransfor-
 436 mation with covariance correction (when 64 principal components are taken
 437 as features) seems to be more useful in contrast to the raw visualization which
 438 also contains the noise cancellation part. This is surprising because this ap-
 439 proach has been originally developed for linear models and not for nonlinear
 440 ones (Haufe et al, 2014). Using a correction with a “local” covariance would
 441 be more appropriate in this case but more demanding from the computation
 442 and implementation point of view. A large number of principal components
 443 seems to be a bad choice for the nonlinear kernel, because it does not seem to

⁸ Nevertheless, the resulting graphics look reasonable.

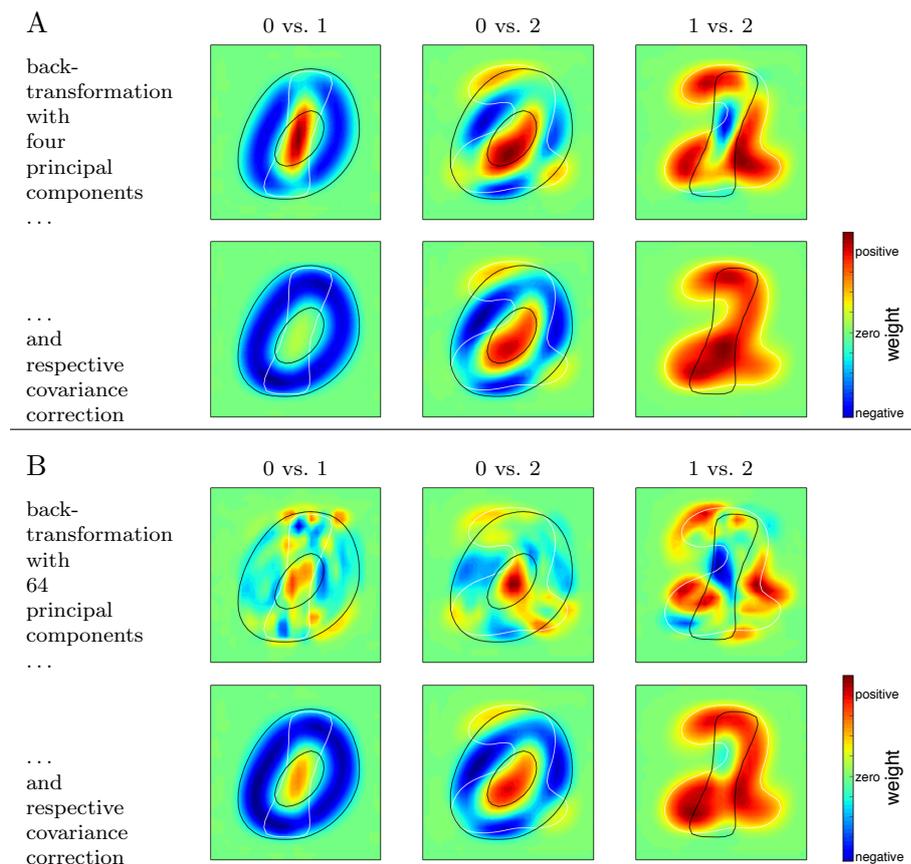


Fig. 2: *Contour plots of backtransformation weights for handwritten digit classification*: The white and black silhouettes display an average contour of the original data (digits 0 vs. 1, 0 vs. 2, and 1 vs. 2). The colored contour plots show the respective weights in the classification process before and after covariance correction with a different number of used principal components (case A and B). Negative weights (blue) are important for the classification of the first class (black silhouette) and positive weights (red) for the second class (white silhouette). Green weights are close to zero and do only contribute weakly to the classification process.

444 generalize that well and is using a lot of small components instead of focusing
 445 on the big shape of the digits.

446 In case of using only 4 principal components, the approach mainly shows
 447 the shape of the digit 2 (or 0 for the first column). In contrast, the visualiza-
 448 tions without covariance correction clearly indicate with a blue color which
 449 parts are relevant for classifying it as the first class and with the red color
 450 which parts are important for the second class. An interesting effect occurs for

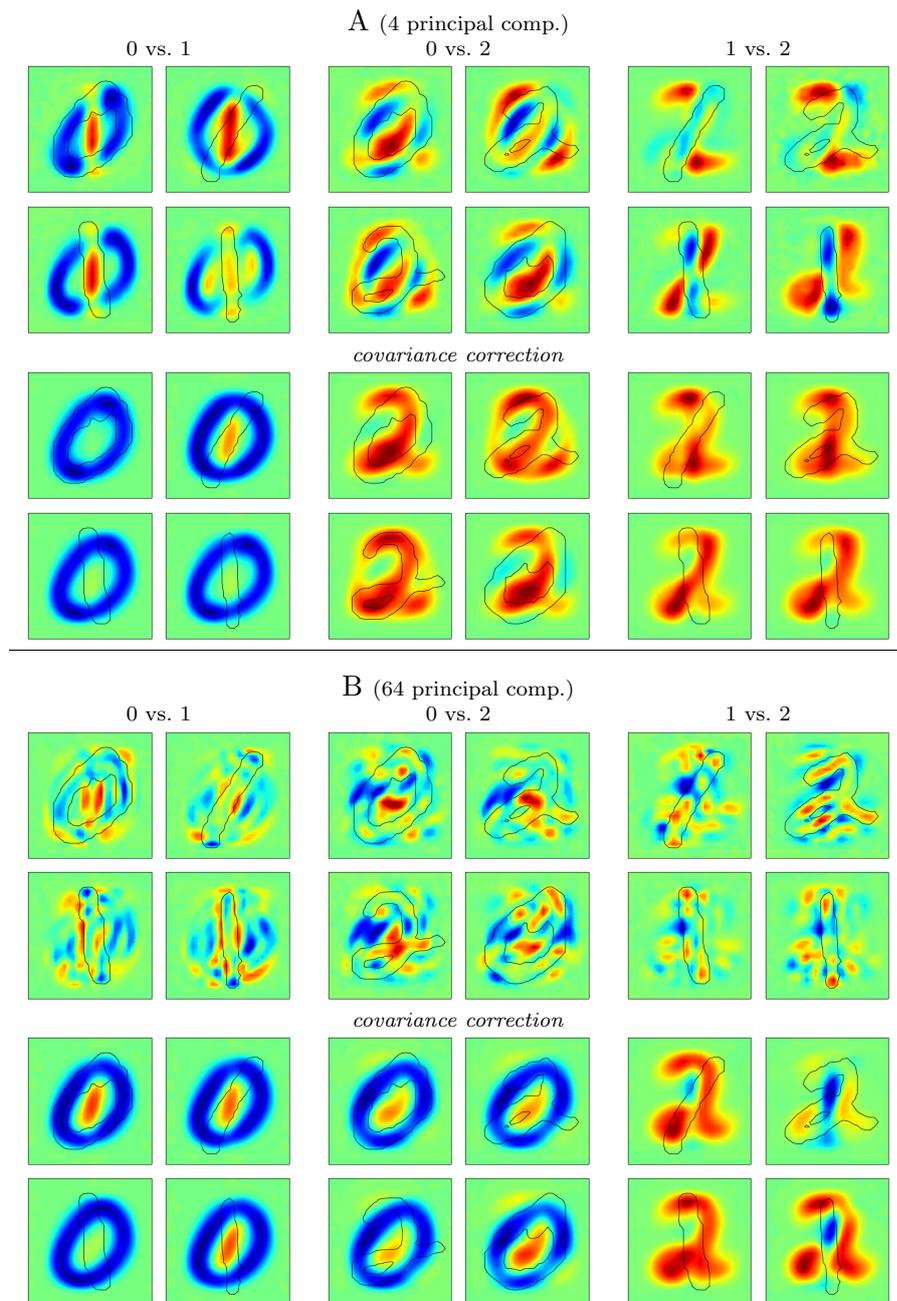


Fig. 3: *Contour plots of backtransformation weights for handwritten digit classification with nonlinear classifier:* The setting is the same as in Fig. 2 except that no average shapes are displayed but the shape of the sample of interest where the backtransformation is calculated for.

451 the first classifier at the fourth digit (1). Here a closer look could be taken at
452 the classifier and the data to find out why there are yellow weights outside the
453 regular shape of the digit 1. This might be the result of some artifacts in the
454 data (e.g., a sample with very bad handwriting near to the observed sample)
455 or an artifact in the processing.

456 In the nonlinear and the linear case with 64 principal components the
457 backtransformation reveals that the decision process is not capable of deriving
458 real shape features for the digits. This might be a reason, why a specially tuned
459 deep neural network performs better in this classification task (Schmidhuber,
460 2012).

461 3.4 Processing Chain Visualization: Movement Prediction from EEG Data

462 The electroencephalogram (EEG) is a very complex signal, measuring elec-
463 trical activity on the scalp with a very high temporal resolution and more
464 than 100 sensors. Several visualization techniques exist for this type of signal,
465 which are used in neuroscience for analysis. When processing EEG data for
466 brain-computer interfaces (BCIs), there is a growing interest in understanding
467 the properties of processing chains and the dynamics of the data, to avoid
468 relying on artifacts and to get information on the original signal back for fur-
469 ther interpretation. Here, very often spatial filtering is used for dimensionality
470 reduction to linearly combine the signals from the numerous electrodes to a
471 largely reduced number of new virtual sensors with much less noise (see Sec-
472 tion 2.3). These spatial filters and much more importantly the data patterns
473 they are enhancing are visualized with similar methods as used for visualizing
474 data. If the spatial filter is the main part of the processing (e.g., only two
475 filters are used), this approach is sufficient to understand the data processing.
476 However, often more filters and other, additional preprocessing algorithms are
477 used. Hence, the original spatial information cannot be determined for the in-
478 put of the classifier. This disables a good visualization of the classifier and an
479 understanding of what the classifier learned from the training data. So here,
480 backtransformation can be very helpful.

481 To illustrate this, a dataset from an EEG experiment was taken (Tabie and
482 Kirchner, 2013). In this experiment, subjects were instructed to move their
483 right arm as fast as possible from a flat board to a buzzer in approximately
484 30 cm distance. The classification task was to predict upcoming movements
485 by detecting movement-related cortical potentials (Johanshahi and Hallett,
486 2003) in the EEG single trials. Before applying the backtransformation and
487 visualizing the data as depicted in Fig. 4, the data has been normalized with a
488 standardization, a decimation, and temporal filtering. Only the last part of the
489 signal, which is close to the movement, was visualized. The processing chain
490 was similar to the one in Section 2.3. The details are described by Seeland
491 et al (2013).

492 The averaged input data in Fig. 4 shows a very strong negative activa-
493 tion at the motor cortex mainly at the left hemisphere around the electrode

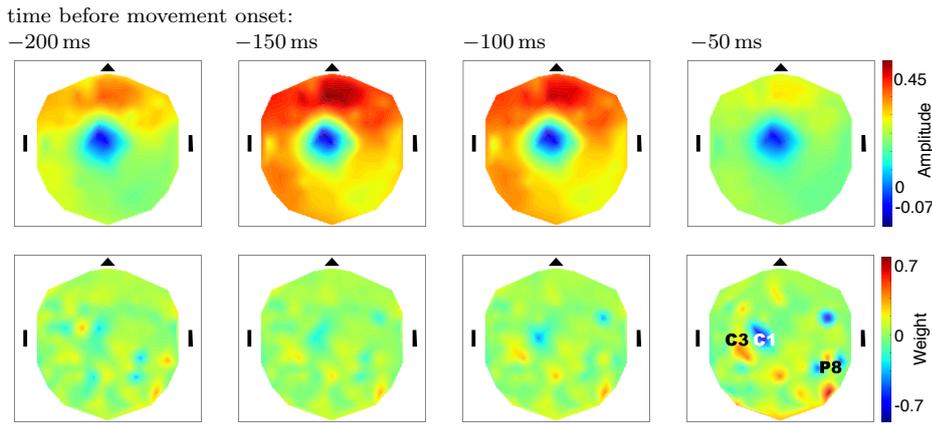


Fig. 4: Visualization of data for movement prediction and the corresponding processing chain: In the first row the average of the data before a movement is displayed as topography plots and in the second row the backtransformation weights are displayed, respectively. The data values from the different sensors were mapped to the respective position on the head, displayed as an ellipse with the nose at the top and the ears on the sides.

494 C1⁹. This activation is consistent with the occurrence of movement related
 495 cortical potentials and is expected from the EEG literature (Johanshahi and
 496 Hallett, 2003). The region of the activation (blue circle on the left hemisphere
 497 at the motor cortex region) is associated with right arm movements, which
 498 the subjects had to perform in the experiment.

499 The backtransformation weights are much more spread over the head com-
 500 pared to the averaged data. There is a major activation at the left motor
 501 cortex at electrodes C1 (negative) and C3 (positive), but also a large spread
 502 activation at the back of the head at the right hemisphere (around the elec-
 503 trode P8). On the time scale, the most important weights can be found at the
 504 last time point, 50 ms before movement onset.

505 This is reasonable, because the most important movement related infor-
 506 mation is expected to be just before the movement starts, although movement
 507 intention can be detected above chance level on average 460 ms before the
 508 movement onset (Lew et al, 2012). Note that the analysis has been performed
 509 on single trials and not on averaged data and that for a good classification
 510 the largest difference is of interest and not the minimal one. The high weights
 511 at C1 and CP3 clearly fit to the high negative activation found in the aver-
 512 aged data and as such highlight the signal of interest. For interpreting the
 513 other weights, two things have to be kept in mind. First, EEG data usually
 514 contains numerous artifacts and second, due to the conductivity of the skin

⁹ A standard extended 10–20 electrode layout has been chosen with 128 electrodes:
http://www.brainproducts.com/filedownload.php?path=downloads/actiCAP-128-channel-Standard-2_1201.pdf.

515 it is possible to measure every electric signal at a certain electrode also on
 516 the other electrodes. Keeping that in mind, the activation around P8 could
 517 be interpreted as a noise filter for the more important class related signal at
 518 C1 and CP3. This required filtering effect on EEG data is closely related to
 519 spatial filtering, which emphasizes a certain spatial pattern (Blankertz et al,
 520 2011, section 4.2). It could be also a relevant signal which cannot be observed
 521 in the plot of the averaged data. These observations are now a good starting
 522 point for domain experts to take a closer look at the raw data to determine
 523 which interpretation fits better.

524 3.5 Applications Beyond Visualization

525 In the following, we shortly describe two further applications of the affine
 526 backtransformation.

527 3.5.1 Group Ranking

528 The formula in Equation (18) has the same structure as a normal linear classi-
 529 fier. Each weight gives an information about the importance of certain signal
 530 components of the input. Summing up the absolute values of one sensor can
 531 now be done in the way as suggested for SVMs (Feess et al, 2013; Lal et al,
 532 2004) to get a sensor ranking:

$$R_h = \sum_{g=1}^{n_g} |w_{gh}^{(0)}|. \quad (23)$$

533 Such a ranking can then be used for sensor selection algorithms to reduce
 534 the number of used electrodes for a BCI and ease comfort and save costs.
 535 It can be used in robotics or other applications too, where the reduction of
 536 input sources can be beneficial. A similar ranking could be also applied to the
 537 time points. The advantage of this ranking method is that it directly operates
 538 on the processing chain and not solely on the input data or feature domain.
 539 Note, that the quality of the ranking also highly depends on the quality of
 540 the processing chain. If a processing chain is worse than an other, chances are
 541 high that also its ranking is worse (Krell, 2015, section 3.4.3).

542 3.5.2 Reinitialization of Linear Classifier with Affine Preprocessing

543 There could be several reasons for exchanging the preprocessing in a signal
 544 processing chain. For example, first some initial preprocessing is loaded but
 545 in parallel a new better fitting data specific processing is trained or tuned
 546 on new incoming data (e.g., a new spatial filter (Woehrle et al, 2015)). If
 547 dimensionality would not be fitting after changing the preprocessing chain, a
 548 new classifier would also be needed. But even if dimensions of old and new
 549 preprocessing were the same it might be good to *adapt* the classifier to that

550 change to have a better initialization. Here, the affine backtransformation can
 551 be used as described in the following.

For this application, a processing chain of affine transformations is assumed which ends with a sample weighting online learning algorithm like the passive aggressive algorithm or a perceptron. Since the classification function is a weighted sum of samples, it enables the following calculation:

$$w = \sum_i \alpha_i y_i \hat{x}_i = \sum_i \alpha_i y_i (Ax_i + T) = A \sum_i \alpha_i y_i x_i + T \sum_i \alpha_i y_i \quad (24)$$

$$= Aw^{(0)} + Tb \text{ with } w^{(0)} = \sum_i \alpha_i y_i x_i \text{ and } b = \sum_i \alpha_i y_i. \quad (25)$$

Here, x_i is the training data with the training labels y_i and \hat{x}_i is the preprocessed training data given to the classifier. The weights α_i are calculated by update formulas of the classifier. During the update step, $w^{(0)}$ must be calculated additionally but neither x_i , y_i , nor α_i are stored. When changing the preprocessing from (A, T) to (A', T')

$$w' = A'w^{(0)} \quad (26)$$

is a straightforward estimate for the new classifier. The advantage of this formula is, that it just requires additionally calculating and storing $w^{(0)}$. So the resulting classifier can be still used for memory efficient online learning. Especially, even if neither (A', T') nor (A, T) is known, w' can be calculated using the new signal processing function $\hat{F}(x) = A'x + T'$:

$$w' = A'w^{(0)} = \hat{F}(w^{(0)}) - T'b = \hat{F}(w^{(0)}) - 0A'w^{(0)}b - T'b = \hat{F}(w^{(0)}) - \hat{F}(0w^{(0)})b. \quad (27)$$

552 So, w' can be computed by processing $w^{(0)}$ and a sample of zero entries in the
 553 signal processing chain. This only requires some minor processing time but
 554 no additional resources. Usually the processing chain is very fast and so the
 555 additional processing time should not be a problem. For giving a proof of concept,
 556 this application of the backtransformation was used in a setting, where
 557 the preprocessing was randomly changed. With the aforementioned approach
 558 the change could be perfectly compensated without any loss in performance
 559 (Krell, 2015, section 2.4.6).

560 4 Conclusion

561 In this paper, a direct approach is given to look at the complete data processing
 562 chain (in contrast to separate handling of its components) and to transform
 563 it to a representation in the same format as the data. This could be used
 564 to improve the understanding of complex processing chains and might enable
 565 several applications in future. It was shown that backtransformation can be
 566 used for visualization of the decision process and a direct comparison with
 567 a visualization of the data is possible and enables an interpretation of the
 568 processing. Our approach extends existing algorithms by also considering the

569 preprocessing, by putting no restrictions on the decision algorithm, by provid-
570 ing the implementation details, and integrating the backtransformation in the
571 pySPACE framework which already comes with a large number of available
572 algorithms.

573 Backtransformation can be used for interpreting the behaviour of the deci-
574 sion process, but it remains an open question on how the further analysis is
575 performed, so that additional investigations and expert knowledge might be
576 required. A related problem occurs when using temporal and spatial filters.
577 Here the solution is to visualize the frequency response and the spatial pat-
578 tern instead of the pure weights of the transformation. The frequency response
579 gives information on how frequencies are filtered out and spatial patterns give
580 information on which signal in space is emphasized by the respective spatial
581 filter. It would be interesting to develop new methods, which improve the inter-
582 pretability of the decision process, e.g., by extending the method of covariance
583 multiplication with a more sophisticated calculation of the covariance matrix
584 or by deriving a different formula for getting the forward model. This might
585 enable the backtransformation to reveal new signals or connections in the data
586 which can then be used to improve the observed data processing chain.

587 In future, it would be interesting to further analyze the application of
588 the backtransformation, e.g., by using other data or processing chains, by
589 analyzing regression problems, or by integrating it into other algorithms and
590 analyzing its benefit.

591 References

- 592 Abdi H, Williams LJ (2010) Principal component analysis. *Wiley Interdisci-*
593 *plinary Reviews: Computational Statistics* 2(4):433–459, DOI 10.1002/wics.
594 101
- 595 Aksoy S, Haralick RM (2001) Feature normalization and likelihood-based sim-
596 ilarity measures for image retrieval. *Pattern Recognition Letters* 22(5):563–
597 582, DOI 10.1016/S0167-8655(00)00112-4
- 598 Baehrens D, Schroeter T, Harmeling S, Kawanabe M, Hansen K, Müller KR
599 (2010) How to Explain Individual Classification Decisions. *Journal of Ma-*
600 *chine Learning Research* 11:1803–1831
- 601 Blankertz B, Tomioka R, Lemm S, Kawanabe M, Müller KR (2008) Optimizing
602 spatial filters for robust EEG single-trial analysis. *IEEE Signal Processing*
603 *Magazine* 25(1):41–56, DOI 10.1109/MSP.2008.4408441
- 604 Blankertz B, Lemm S, Treder M, Haufe S, Müller KR (2011) Single-trial analy-
605 sis and classification of ERP components—a tutorial. *NeuroImage* 56(2):814–
606 825, DOI 10.1016/j.neuroimage.2010.06.048
- 607 Chang CC, Lin CJ (2011) LIBSVM. *ACM Transactions on Intelligent Systems*
608 *and Technology* 2(3):1–27, DOI 10.1145/1961189.1961199
- 609 Chen Ch, Härdle W, Unwin A (2008) *Handbook of Data Visualization*.
610 Springer Handbooks of Computational Statistics, Springer

- 611 Clarke F (1990) Optimization and Nonsmooth Analysis. Society for Industrial
612 and Applied Mathematics, DOI 10.1137/1.9781611971309
- 613 Crammer K, Dekel O, Keshet J, Shalev-Shwartz S, Singer Y (2006) Online
614 Passive-Aggressive Algorithms. *Journal of Machine Learning Research* 7:551
615 – 585
- 616 Domingos P (2012) A few useful things to know about machine learning. *Com-*
617 *munications of the ACM* 55(10):78–87, DOI 10.1145/2347736.2347755
- 618 Feess D, Krell MM, Metzen JH (2013) Comparison of sensor selection mecha-
619 nisms for an ERP-based brain-computer interface. *PLoS ONE* 8(7):e67,543,
620 DOI 10.1371/journal.pone.0067543
- 621 Ghaderi F, Straube S (2013) An adaptive and efficient spatial filter for event-
622 related potentials. In: *Proceedings of the 21st European Signal Processing*
623 *Conference, (EUSIPCO)*
- 624 Griewank A, Walther A (2008) Evaluating Derivatives: Principles and Techn-
625 niques of Algorithmic Differentiation. Society for Industrial and Applied
626 Mathematics
- 627 Haufe S, Meinecke F, Görgen K, Dähne S, Haynes JD, Blankertz B, Bießmann
628 F (2014) On the interpretation of weight vectors of linear models in mul-
629 tivariate neuroimaging. *NeuroImage* 87:96–110, DOI 10.1016/j.neuroimage.
630 2013.10.067
- 631 Johanshahi M, Hallett M (eds) (2003) The Bereitschaftspotential: movement-
632 related cortical potentials. Kluwer Academic/Plenum Publishers, New York,
633 USA
- 634 Jutten C, Herault J (1991) Blind separation of sources, part I: An adaptive
635 algorithm based on neuromimetic architecture. *Signal Processing* 24(1):1–
636 10, DOI 10.1016/0165-1684(91)90079-X
- 637 Kirchner EA, Kim SK, Straube S, Seeland A, Wöhrle H, Krell MM, Tabie
638 M, Fahle M (2013) On the applicability of brain reading for predictive
639 human-machine interfaces in robotics. *PLoS ONE* 8(12):e81,732, DOI
640 10.1371/journal.pone.0081732
- 641 Krell MM (2015) Generalizing, Decoding, and Optimizing Support Vector
642 Machine Classification. PhD thesis, University of Bremen, Bremen, URL
643 <http://nbn-resolving.de/urn:nbn:de:gbv:46-00104380-12>
- 644 Krell MM, Wöhrle H (2015) New one-class classifiers based on the origin
645 separation approach. *Pattern Recognition Letters* 53:93–99, DOI 10.1016/j.
646 patrec.2014.11.008
- 647 Krell MM, Straube S, Seeland A, Wöhrle H, Teiwes J, Metzen JH, Kirchner
648 EA, Kirchner F (2013a) pySPACE - a signal processing and classification
649 environment in Python. *Frontiers in Neuroinformatics* 7(40), DOI 10.3389/
650 fninf.2013.00040, <https://github.com/pyspace>
- 651 Krell MM, Tabie M, Wöhrle H, Kirchner EA (2013b) Memory and Processing
652 Efficient Formula for Moving Variance Calculation in EEG and EMG Sig-
653 nal Processing. In: *Proc. International Congress on Neurotechnology, Elec-*
654 *tronics and Informatics (NEUROTECHNIX 2013)*, ScitePress, Vilamoura,
655 Portugal, pp 41–45, DOI 10.5220/0004633800410045

- 656 Krell MM, Feess D, Straube S (2014a) Balanced Relative Margin Machine The
657 missing piece between FDA and SVM classification. *Pattern Recognition*
658 *Letters* 41:43–52, DOI 10.1016/j.patrec.2013.09.018
- 659 Krell MM, Straube S, Wöhrle H, Kirchner F (2014b) Generalizing, Optimizing,
660 and Decoding Support Vector Machine Classification. In: *ECML/PKDD*
661 *2014 PhD Session Proceedings, Nancy*
- 662 LaConte S, Strother S, Cherkassky V, Anderson J, Hu X (2005) Support vector
663 machines for temporal classification of block design fMRI data. *NeuroImage*
664 *26(2)*:317–329, DOI 10.1016/j.neuroimage.2005.01.048
- 665 Lagerlund TD, Sharbrough FW, Busacker NE (1997) Spatial filtering of mul-
666 tichannel electroencephalographic recordings through principal component
667 analysis by singular value decomposition. *Journal of Clinical Neurophysiol-*
668 *ogy* 14(1):73–82
- 669 Lal TN, Schröder M, Hinterberger T, Weston J, Bogdan M, Birbaumer N,
670 Schölkopf B (2004) Support vector channel selection in BCI. *IEEE Engi-*
671 *neering in Medicine and Biology Society* 51(6):1003–1010, DOI 10.1109/
672 *TBME.2004.827827*
- 673 Le QV, Ranzato M, Monga R, Devin M, Chen K, Corrado GS, Dean J, Ng AY
674 (2012) Building high-level features using large scale unsupervised learning.
675 In: *International Conference on Machine Learning*
- 676 LeCun Y, Bottou L, Bengio Y, Haffner P (1998) Gradient-based learning ap-
677 plied to document recognition. *Proceedings of the IEEE* 86(11):2278–2324,
678 DOI 10.1109/5.726791
- 679 Lew E, Chavarriaga R, Zhang H, Seeck M, del Millan J (2012) Self-paced
680 movement intention detection from human brain signals: Invasive and non-
681 invasive EEG. In: *2012 Annual International Conference of the IEEE Engi-*
682 *neering in Medicine and Biology Society (EMBC)*, pp 3280–3283
- 683 Lin HT, Lin CJ, Weng RC (2007) A note on Platts probabilistic outputs for
684 support vector machines. *Machine Learning* 68(3):267–276, DOI 10.1007/
685 *s10994-007-5018-6*
- 686 Metzen JH, Kirchner EA (2011) Rapid adaptation of brain reading interfaces
687 based on threshold adjustment. In: *Proceedings of the 2011 Conference of*
688 *the German Classification Society, (GfKI-2011), Frankfurt, Germany*, p 138
- 689 Mika S, Rätsch G, Müller KR (2001) A mathematical programming approach
690 to the kernel fisher algorithm. In: *Advances in Neural Information Process-*
691 *ing Systems 13 (NIPS 2000)*, MIT Press, pp 591–597
- 692 Oppenheim AV, Schaffer RW (2009) *Discrete-Time Signal Processing*, 3rd edn.
693 Prentice Hall Press
- 694 Platt JC (1999) Probabilistic Outputs for Support Vector Machines and Com-
695 parisons to Regularized Likelihood Methods
- 696 Press W (2007) *Numerical recipes: the art of scientific computing*, 3rd edn.
697 Cambridge University Press
- 698 Rieger J, Kosar K, Lhotska L, Krajca V (2004) Eeg data and data analysis visu-
699 alization. In: Barreiro J, Martn-Snchez F, Maojo V, Sanz F (eds) *Biological*
700 *and Medical Data Analysis, Lecture Notes in Computer Science*, vol 3337,
701 Springer Berlin Heidelberg, pp 39–48, DOI 10.1007/978-3-540-30547-7_5

- 702 Rivet B, Souloumiac A, Attina V, Gibert G (2009) xDAWN Algorithm to En-
703 hance Evoked Potentials: Application to Brain-Computer Interface. *IEEE*
704 *Transactions on Biomedical Engineering* 56(8):2035–2043, DOI 10.1109/
705 TBME.2009.2012869
- 706 Rockafellar RT, Wets RJB (2009) *Variational analysis*, vol 317. Springer Sci-
707 ence & Business Media
- 708 Schmidhuber J (2012) Multi-column deep neural networks for image classifi-
709 cation. In: *Proceedings of the 2012 IEEE Conference on Computer Vision*
710 *and Pattern Recognition (CVPR)*, IEEE Computer Society, pp 3642–3649
- 711 Schölkopf B, Platt JC, Shawe-Taylor J, Smola AJ, Williamson RC (2001) Esti-
712 mating the support of a high-dimensional distribution. *Neural computation*
713 13(7):1443–1471, DOI 10.1162/089976601750264965
- 714 Seeland A, Wöhrle H, Straube S, Kirchner EA (2013) Online movement pre-
715 diction in a robotic application scenario. In: *6th International IEEE EMBS*
716 *Conference on Neural Engineering (NER)*, San Diego, USA, pp 41–44, DOI
717 10.1109/NER.2013.6695866
- 718 Smola AJ, Schölkopf B (2004) A tutorial on support vector regression. *Statist-
719 ics and Computing* 14(3):199–222, DOI 10.1023/B:STCO.0000035301.
720 49549.88
- 721 Steinwart I, Christmann A (2008) *Support Vector Machines*. Springer
- 722 Straube S, Feess D (2013) Looking at ERPs from Another Perspective: Poly-
723 nomial Feature Analysis. *Perception* 42 *ECVP Abstract Supplement*:220
- 724 Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R
725 (2014) Intriguing properties of neural networks. In: *International Conference*
726 *on Learning Representations*
- 727 Tabie M, Kirchner EA (2013) EMG Onset Detection - Comparison of differ-
728 ent methods for a movement prediction task based on EMG. In: Al-
729 varez S, Solé-Casals J, Fred A, Gamboa H (eds) *In Proceedings of the*
730 *6th International Conference on Bio-inspired Systems and Signal Process-*
731 *ing (BIOSIGNALS-13)*, SciTePress, Barcelona, Spain, pp 242–247, DOI
732 10.5220/0004250102420247
- 733 Vapnik VN (1995) *The nature of statistical learning theory*. Springer
- 734 Varewyck M, Martens JP (2011) A practical approach to model selection
735 for support vector machines with a Gaussian kernel. *IEEE transactions*
736 *on systems, man, and cybernetics Part B, Cybernetics : a publication of*
737 *the IEEE Systems, Man, and Cybernetics Society* 41(2):330–340, DOI
738 10.1109/TSMCB.2010.2053026
- 739 Verhoeve J, de Wulf R (1999) An Image Processing Chain for Land-Cover Clas-
740 sification Using Multitemporal ERS-1 Data. *Photogrammetric Engineering*
741 *& Remote Sensing* 65(10):1179–1186
- 742 Woehrle H, Krell MM, Straube S, Kim SK, Kirchner EA, Kirchner F (2015)
743 An Adaptive Spatial Filter for User-Independent Single Trial Detection of
744 Event-Related Potentials. *IEEE Transactions on Biomedical Engineering*
745 DOI 10.1109/TBME.2015.2402252