

Predictable Feature Analysis

Stefan Richthofer*, Laurenz Wiskott†

*Institute For Neural Computation,
Ruhr-Universität Bochum, Germany*

August 21, 2018

Abstract

Every organism in an environment, whether biological, robotic or virtual, must be able to predict certain aspects of its environment in order to survive or perform whatever task is intended. It needs a model that is capable of estimating the consequences of possible actions, so that planning, control, and decision-making become feasible. For scientific purposes, such models are usually created in a problem specific manner using differential equations and other techniques from control- and system-theory. In contrast to that, we aim for an unsupervised approach that builds up the desired model in a self-organized fashion. Inspired by Slow Feature Analysis (SFA), our approach is to extract sub-signals from the input, that behave as predictable as possible. These “predictable features” are highly relevant for modeling, because predictability is a desired property of the needed consequence-estimating model by definition. In our approach, we measure predictability with respect to a certain prediction model. We focus here on the solution of the arising optimization problem and present a tractable algorithm based on algebraic methods which we call Predictable Feature Analysis (PFA). We prove that the algorithm finds the globally optimal signal, if this signal can be predicted with low error. To deal with cases where the optimal signal has a significant prediction error, we provide a robust, heuristically motivated variant of the algorithm and verify it empirically. Additionally, we give formal criteria a prediction-model must meet to be suitable for measuring predictability in the PFA setting and also provide a suitable default-model along with a formal proof that it meets these criteria.

*Electronic address: stefan.richthofer@ini.rub.de; Corresponding author

†Electronic address: laurenz.wiskott@ini.rub.de

1 Introduction

The motivation for Predictable Feature Analysis (PFA) comes from typical reinforcement-learning settings, where an autonomous agent is placed in an environment and aims to achieve some goal. While many common scenarios are discrete (board-game-like) with rather few states, we consider more natural scenarios, where the input is a continuous signal over time and of high dimension like vision or some other sensory input.¹

PFA is intended as a tool to help the agent make sense of this vast amount of incoming data. Our approach is to look for information in the signal that helps to understand and manipulate the environment in the desired way. To achieve its goal, the agent must be able to plan its actions and thus needs to understand, how the environment behaves – it needs a model that is capable of predicting the outcomes of possible actions. It has been frequently proposed that predictable features are crucial to obtain such a model, see [1] for a review. In contrast to most common approaches from control theory, we attempt to perform the modeling without putting previously known, problem-specific information (usually a representation of the environment in form of differential equations and system theoretic setups) into the model, but look for a truly unsupervised, self-organized approach.

Slow Feature Analysis (SFA) is an algorithm that has most characteristics we are looking for (and as such also served as the name-giving pattern for PFA). It is an algorithm that has proven valuable in several fields and problems concerning signal- and data analysis. The idea is that a drastic, yet reasonable dimensionality reduction can be obtained by focusing on slowly varying sub-signals, the so-called “slow features”. These are considered most relevant, because slowness usually indicates invariance and invariant problem representations are crucial for typical data-analysis and recognition tasks, such as regression and classification. Many of these tasks have proven to become much more feasible on the reduced signal after SFA has been applied. For instance, tasks like the self-organization of complex-cell receptive fields, the recognition of whole objects invariant to spatial transformations, the self-organization of place-cells, extraction of driving forces, or nonlinear blind source separation were successfully performed on the basis of SFA (see [2, 3, 4, 5, 6]).

PFA extracts sub-signals from the input using the same methods like SFA does, but instead of the slowest features, it selects those that are best predictable by a certain prediction model. In section 3 we give the criteria a model must meet to be suitable for this purpose. While there are also model-independent notions of predictability like in the information bottleneck approach [7, 8], focusing on concrete models has the advantage that if PFA finds predictable sub-signals, one is directly able to actually perform the prediction, since the appropriate model is given. The arising optimization problem turned out to be significantly harder than that of SFA, because it is a nested problem: The features extracted must be optimized for predictability, but judging their predictability is an optimization problem by itself. Optimizing these problems in turns usually converges to sub-optimal solutions that depend on the starting point. In this work we discuss the details of the PFA-problem and present a tractable algorithm, thus setting up the basis for future PFA-related work.

There have been other approaches to use notions of predictability. For instance [9] considers scenarios involving embodied agents. They let two sensors predict each other in order to retrieve representation-invariant information. [8] combines notions of predictability with SFA to better understand principles of sensory coding strategies. There also exists an ICA-based approach to

¹ Since it comes from a technical setup, the signal would still be discrete. However, we would not regard its discreteness as states but would conceptually treat it as a continuous signal.

predictability-driven dimensionality reduction, see [10]. ForeCA (Forecastable Component Analysis), an independently developed method, is based on the same paradigm as PFA, but proposes a model-independent approach [11]. A further difference is that PFA (optionally) searches for well predictable Systems, while ForeCA selects best predictable single components. In future work, we are going to compare PFA- and ForeCA-results to get a better understanding for strengths and weaknesses of the two approaches. Finally, there has also been a previous version of our PFA approach [12].

2 Extracting predictable features

Given an input-signal $\mathbf{x}(t)$ with n components, our goal is to extract a certain number (r) of well predictable output-components, referred to as “predictable features”. Since our approach is inspired by SFA, we start with a summary of that algorithm.

2.1 Recall SFA

In the SFA-setting, the optimized property is slow variation. Extraction is in principle performed by linear transformation and projection.² The parameters of these mappings are optimized over a finite training-phase Ω_t consisting of equidistant time points. To make the method more powerful, a non-linear expansion \mathbf{h} can be applied to the signal – usually using monomials of low degree.

In order to avoid the trivial constant solution, the output is constrained to have unit variance and zero mean. Additionally, the output-components must be pairwise uncorrelated. This way the repeated occurrence of the same component is avoided. Mean is defined using $\langle s(t) \rangle_t := \frac{1}{|\Omega_t|} \sum_{t \in \Omega_t} s(t)$ (average of a signal over the training phase). To fulfill the constraints, the expanded signal is sphered over the training-phase, i.e. its mean is shifted to zero and the covariance-matrix is normalized to the identity matrix:

$$\tilde{\mathbf{z}}(t) := \mathbf{h}(\mathbf{x}(t)) - \langle \mathbf{h}(\mathbf{x}(t)) \rangle_t \quad (\text{make mean-free}) \quad (1)$$

$$\mathbf{z}(t) := \mathbf{S}\tilde{\mathbf{z}}(t) \quad \text{with} \quad \mathbf{S} := \langle \tilde{\mathbf{z}}\tilde{\mathbf{z}}^T \rangle^{-\frac{1}{2}} \quad (\text{normalize covariance}) \quad (2)$$

Summing up all SFA-constraints, the following optimization problem is derived:

$$\begin{aligned} & \text{For } i \in \{1, \dots, r\} \\ & \underset{\mathbf{a}_i \in \mathbb{R}^n}{\text{minimize}} \quad \mathbf{a}_i^T \langle \dot{\mathbf{z}}\dot{\mathbf{z}}^T \rangle \mathbf{a}_i \\ & \text{subject to} \quad \mathbf{a}_i^T \langle \mathbf{z} \rangle = 0 \quad (\text{zero mean}) \\ & \quad \mathbf{a}_i^T \langle \mathbf{z}\mathbf{z}^T \rangle \mathbf{a}_i = 1 \quad (\text{unit variance}) \\ & \quad \mathbf{a}_i^T \langle \mathbf{z}\mathbf{z}^T \rangle \mathbf{a}_j = 0 \quad \forall j < i \quad (\text{pairwise uncorrelated}) \end{aligned} \quad (3)$$

To describe the algorithm, we first define the extraction matrix $\mathbf{A}_r := (\mathbf{a}_1, \dots, \mathbf{a}_r) \in \mathbb{R}^{n \times r}$ and the reduced identity $\mathbf{I}_r \in \mathbb{R}^{n \times r}$ consisting of the first r euclidean unit vectors as columns. Now please

²In a strict sense, the transformation is affine because it clears the signal’s mean. Additionally one can count the non-linear expansion as part of the extraction.

note that because of the sphering, it holds that $\langle \mathbf{z} \rangle = 0$ and $\langle \mathbf{z}\mathbf{z}^T \rangle = \mathbf{I}$, thus having the constraints equal to

$$\exists \mathbf{A} \in O(n): \quad \mathbf{A}_r = \mathbf{A}\mathbf{I}_r \quad (4)$$

where $O(n) \subset \mathbb{R}^{n \times n}$ denotes the space of orthogonal transformations, i.e. $\mathbf{A}\mathbf{A}^T = \mathbf{I}$. Choosing \mathbf{a}_i as the eigenvectors of $\langle \mathbf{z}\mathbf{z}^T \rangle$, corresponding to the eigenvalues in descending order, yields an \mathbf{A}_r that solves (3) globally. We denote the extracted signal with $\mathbf{m} := \mathbf{A}_r^T \mathbf{z}$. [13] describes this procedure in detail.

2.2 Modeling the PFA-problem

In order to measure predictability, we focus on a certain prediction-model. Because it is simple and very popular, we use **linear, auto-regressive prediction** as our default model – it is successfully used in many fields for modeling time-related problems. A signal is regarded well predictable, if each value can be approximated by a linear combination of some (p) recent values. Expressing this formally, we face the problem of finding vectors \mathbf{a} and \mathbf{b} such that

$$\begin{aligned} \mathbf{a}^T \mathbf{z}(t) &\stackrel{!}{\approx} b_1 \mathbf{a}^T \mathbf{z}(t-1) + \dots + b_p \mathbf{a}^T \mathbf{z}(t-p) \\ &= \mathbf{a}^T \text{hist}_{\mathbf{z},p}(t) \mathbf{b} \end{aligned} \quad (5)$$

with hist defined as the signal's history over the recent p time-steps:

$$\text{hist}_{\mathbf{z},p,\Delta}(t) := \sum_{i=1}^p \mathbf{z}(t-i\Delta) \mathbf{e}_i^T \quad \text{with} \quad \mathbf{e}_i \in \mathbb{R}^p, \quad (\mathbf{e}_1, \dots, \mathbf{e}_p) = \mathbf{I}_{p,p}. \quad (7)$$

Here $\mathbf{I}_{p,p}$ denotes the p -dimensional identity, thus \mathbf{e}_i denotes the i -th p -dimensional Euclidean unit vector. Δ defaults to 1: $\text{hist}_{\mathbf{z},p} := \text{hist}_{\mathbf{z},p,1}$.

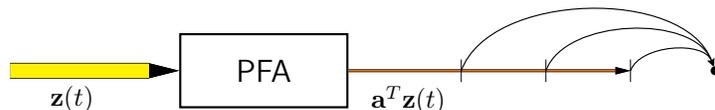


Figure 1: Illustration of PFA

Like in SFA, we optimize the parameters over the training phase Ω_t and also adopt the SFA constraints to avoid trivial or repeated solutions. The first steps of PFA are indeed equal to those in SFA, i.e. we also allow for a non-linear expansion and also start with a sphering-step. As far as possible, we use the notation that was introduced in 2.1. The common way to extend (52) to multiple dimensions can be written as

$$\mathbf{m}(t) \stackrel{!}{\approx} \mathbf{B}_1 \mathbf{m}(t-1) + \dots + \mathbf{B}_p \mathbf{m}(t-p) \quad \text{with} \quad \mathbf{B}_i \in \mathbb{R}^{n \times n}, \text{ diagonal} \quad (8)$$

In this form, it does not fulfill all criteria from section 3 (it is not orthogonal agnostic for $n > 1$). Nevertheless, we mention strategies to solve (54) in the appendix, section A.2. Here we proceed by refining it to be suitable for PFA:

$$\mathbf{m}(t) \stackrel{!}{\approx} \mathbf{B}_1 \mathbf{m}(t-1) + \dots + \mathbf{B}_p \mathbf{m}(t-p) \quad \text{with} \quad \mathbf{B}_i \in \mathbb{R}^{n \times n} \quad . \quad (9)$$

The difference to the first formulation is that each extracted component's prediction may depend on all other extracted components. Note that (54) and (9) are equal for $n = 1$. A massive advantage of model (9) is that we can initially fit it to our data in full dimension and search for the best-fitted components afterwards. For (54), this would not be possible, because the fitting-quality of each component is not invariant under the transformation used for extraction. We formalize the need for such a model-property in section 3.

To formalize fitting, we briefly introduce the notion of a general prediction model and of the fitting-error. We speak of a general prediction-model \mathbf{g} as

$$\mathbf{g} \in \mathcal{G} : \quad \mathbf{z}(t) \stackrel{!}{\approx} \mathbf{g}(\text{hist}_{\mathbf{z},p,\Delta}(t)) \quad (10)$$

where \mathcal{G} is the model-class, i.e. the set of possible realizations of \mathbf{g} . We measure the prediction error in an average least squares sense by

$$\text{err}(\mathbf{g}, \mathbf{z}) := \langle \|\mathbf{z} - \mathbf{g}(\text{hist}_{\mathbf{z},p,\Delta})\|^2 \rangle \quad (11)$$

Now fitting a prediction-model on a given sample in a least squares sense can be expressed as

$$\underset{\mathbf{g} \in \mathcal{G}}{\text{minimize}} \quad \text{err}(\mathbf{g}, \mathbf{z}) \quad (12)$$

By $\mathbf{g}_{\mathbf{z}}^*$, we denote the global solution of (12) and define the following shortcut-notation:

$$\text{err}(\mathbf{z}) := \text{err}(\mathbf{g}_{\mathbf{z}}^*, \mathbf{z}) \quad (13)$$

To formalize (9) as a prediction model in this notation, we combine the coefficient-matrices \mathbf{B}_i to a single broad matrix and define $\mathbf{g}_{\mathbf{B}}$ and $\mathcal{G}_{(9)}$:

$$\mathbf{B} := (\mathbf{B}_1, \dots, \mathbf{B}_p) \in \mathbb{R}^{n \times np} \quad (14)$$

$$\mathbf{g}_{\mathbf{B}}(\text{hist}_{\mathbf{z},p}(t)) := \mathbf{B} \text{vec}(\text{hist}_{\mathbf{z},p}(t)) \quad (15)$$

$$\mathcal{G}_{(9)} := \{ \mathbf{g}_{\mathbf{B}} : \mathbf{B} \in \mathbb{R}^{n \times np} \} \quad (16)$$

By analytic optimization, we obtain the following regression formula to fit $\mathbf{g}_{\mathbf{B}} \in \mathcal{G}_{(9)}$ to $\mathbf{m} = \mathbf{A}_r^T \mathbf{z}$:

$$\mathbf{B}_{\mathbf{z}}(\mathbf{A}_r) = \mathbf{A}_r^T \langle \mathbf{z} \zeta^T \rangle \underline{\mathbf{A}}_r \left(\underline{\mathbf{A}}_r^T \langle \zeta \zeta^T \rangle \underline{\mathbf{A}}_r \right)^{-1} \quad (17)$$

Here we used $\zeta(t) := \text{vec}(\text{hist}_{\mathbf{z},p}(t))$ and the following shortcut notation defined for any matrix \mathbf{M} :

$$\underline{\mathbf{M}} := \mathbf{I}_{p,p} \otimes \mathbf{M} = \begin{pmatrix} \mathbf{M} & \mathbf{0} \\ & \ddots \\ \mathbf{0} & \mathbf{M} \end{pmatrix} \quad (18)$$

p times \mathbf{M}

If $r = n$, $\mathbf{A} = \mathbf{I}$ and thus $\mathbf{A}_r = \mathbf{I}$, we write $\mathbf{W} := \mathbf{B}_{\mathbf{z}}(\mathbf{I}) = \langle \mathbf{z} \zeta^T \rangle \langle \zeta \zeta^T \rangle^{-1}$. (See A.1 for an overview of all notation in this document.) It sometimes happens that $\langle \zeta \zeta^T \rangle$ is not (cleanly) invertible due to some very small or even zero eigenvalues. We regard it best practice to project away the eigenspaces corresponding to eigenvalues below a critical threshold. The intuition behind this is that spaces corresponding to (almost-)zero eigenvalues indicate redundancies in the signal and should not be

used for prediction anyway. To perform this, first compute an eigenvalue decomposition on $\langle \zeta \zeta^T \rangle$. Replace eigenvalues below the threshold by 0 and replace the other ones by their multiplicative inverse. After that undo the decomposition and use the resulting matrix as a proxy for $\langle \zeta \zeta^T \rangle^{-1}$.

For $r = n$ and $\mathbf{A} \in \text{O}(n)$, we have $\mathbf{B}_z(\mathbf{A}) = \mathbf{A}^T \mathbf{W} \underline{\mathbf{A}}$. Since \mathbf{z} is sphered, we can state the following compact notation of the PFA-problem:

$$\underset{\mathbf{A} \in \text{O}(n)}{\text{minimize}} \quad \text{err}(\mathbf{A}_r^T \mathbf{z}) \quad (19)$$

Inserting our default model, we have $\text{err}(\mathbf{A}_r^T \mathbf{z}) = \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{B}_z(\mathbf{A}_r) \underline{\mathbf{A}}_r^T \zeta\|^2 \rangle$. However, because (17) is an involved term, mainly due to the projection under an inversion symbol, (19) appears to be intractable by every method known to us³. Instead of solving it directly, we propose the following tractable relaxation:

$$\underset{\mathbf{A} \in \text{O}(n)}{\text{minimize}} \quad \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{B}_z(\mathbf{A}) \underline{\mathbf{A}}^T \zeta\|^2 \rangle = \langle \|\mathbf{A}_r^T (\mathbf{z} - \mathbf{W} \zeta)\|^2 \rangle \quad (20)$$

Informally speaking, problem (20) asks for components that are optimally predictable, if the prediction may be based on the entire input signal, rather than just on the extracted components themselves. From now on we denote a global optimum of (19) with \mathbf{A}_r^* and of (20) with $\mathbf{A}_r^{(0)}$.

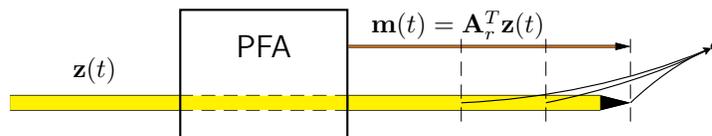


Figure 2: Illustration of relaxation (20)

To solve (20) globally, we write it as

$$\underset{\mathbf{A} \in \text{O}(n)}{\text{minimize}} \quad \text{Tr} \left(\mathbf{A}_r^T \langle (\mathbf{z} - \mathbf{W} \zeta) (\mathbf{z} - \mathbf{W} \zeta)^T \rangle \mathbf{A}_r \right) \quad (21)$$

and choose \mathbf{A} such that it diagonalizes $\langle (\mathbf{z} - \mathbf{W} \zeta) (\mathbf{z} - \mathbf{W} \zeta)^T \rangle$ and sorts the r smallest eigenvalues to the upper left. This method can be described as performing PCA on the residuals of the least squares fit. By some calculus, one can show formal equivalence of this approach to the method proposed in [14]. In section 4 we prove that if $\text{err}((\mathbf{A}_r^*)^T \mathbf{z}) = 0$, then $\mathbf{A}_r^{(0)}$ is also a global solution of (19).⁴ More precisely speaking, the relaxation gap of (20) depends on $\text{err}((\mathbf{A}_r^*)^T \mathbf{z})$ in a continuous manner and is zero, if that error is zero. If the optimal sub-signal has a significant prediction error, the solution obtained as $\mathbf{A}_r^{(0)}$ usually suffers from overfitting and is sub-optimal for (19). In the following, we offer a heuristic method to overcome this overfitting.

³Not counting evolutionary and other inherent local optimization approaches, since we aim for the global solution. Experiments showed us that locally optimal solutions are usually still of high error and of low relevance for the model.

⁴Note that if $\mathbf{A}_r^{(0)}$ is used as solution for (19), the prediction model must be refitted to the reduced signal to get optimal prediction. For this, calculate $\mathbf{B}_z(\mathbf{A}_r^{(0)})$ as defined in (17).

2.3 Avoiding overfitting

To reduce overfitting, we propose the heuristics that signals well predictable in terms of (19) yield a lower error-propagation to subsequent predictions than signals that are well predictable in terms of (20) but not in terms of (19). We ground this on the intuition that the prediction of the latter ones is partly based on noisy data – thus subsequent predictions inherit a higher error. To formalize this idea we define

$$\mathbf{V} := \langle \zeta(t+1)\zeta^T(t) \rangle_t \langle \zeta\zeta^T \rangle^{-1} \quad (22)$$

and can thus perform iterated prediction as follows:

$$\mathbf{z}(t) \approx \mathbf{W}\mathbf{V}^i\zeta(t-i) \quad (23)$$

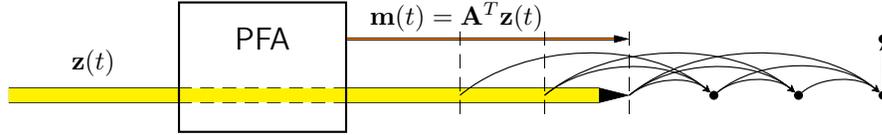


Figure 3: Illustration of iterated prediction

Based on this, we propose the following optimization problem:

$$\underset{\mathbf{A} \in \mathcal{O}(n)}{\text{minimize}} \quad \sum_{i=0}^k \langle \|\mathbf{A}_r^T (\mathbf{z} - \mathbf{W}\mathbf{V}^i\zeta(t-i))\|^2 \rangle_t \quad (24)$$

We can solve it globally in a rather similar way like (20). To do so, we write it as

$$\underset{\mathbf{A} \in \mathcal{O}(n)}{\text{minimize}} \quad \text{Tr} \left(\mathbf{A}_r^T \sum_{i=0}^k \langle (\mathbf{z} - \mathbf{W}\mathbf{V}^i\zeta(t-i)) (\mathbf{z} - \mathbf{W}\mathbf{V}^i\zeta(t-i))^T \rangle_t \mathbf{A}_r \right) \quad (25)$$

and then solve it by diagonalizing $\sum_{i=0}^k \langle (\mathbf{z} - \mathbf{W}\mathbf{V}^i\zeta(t-i)) (\mathbf{z} - \mathbf{W}\mathbf{V}^i\zeta(t-i))^T \rangle_t$ and sorting the lowest r eigenvalues to the upper left. We denote the global solution of (24) by $\mathbf{A}_r^{(k)}$. How to optimally choose k for a certain problem is currently an open question, but we know from experiments that up to some value, increasing k improves $\text{err}((\mathbf{A}_r^{(k)})^T \mathbf{z})$. Beyond that value, increasing k lowers the quality again. Our intuition is that the critical value is related to the maximal time distance, over which the signal holds any auto-correlation – investigating this formally will be subject of future work. To give some impression of the technique and as a proof of concept, we demonstrate it on a synthetic example. We know that basic trigonometric functions are losslessly predictable with our default model for $p = 2$. Because of the theorem in section 4, it would make no sense to work with a losslessly predictable signal. So we add some white noise $\eta(t)$ to it and define an example signal $\mathbf{x}(t) := (\sin(0.1t) + 0.7\eta(t), \sin(0.2t) + \eta(t), \sin(0.4t) + 5.3\eta(t))^T$. We train the algorithm with $\Delta = 1$ once on 1000 samples and once on 2000 samples, always extracting two components, i.e. $r = 2, p = 2$. As a lower bound for any error not involving overfitting, we evaluate (20). For our example we get a lower bound of ≈ 1.206 for 1000 samples and ≈ 1.22 for 2000 samples. These are plotted as red horizontal lines in the following. To measure the amount

of overfitting, we add many dimensions of white random data to our signal and mix everything up by a random, orthogonal transformation. While \mathbf{x} has always the same noise-seed, the added data is generated with different noise in every run. The following results are averaged over about 150 runs, and plot the prediction error against the noise-dimension.⁵

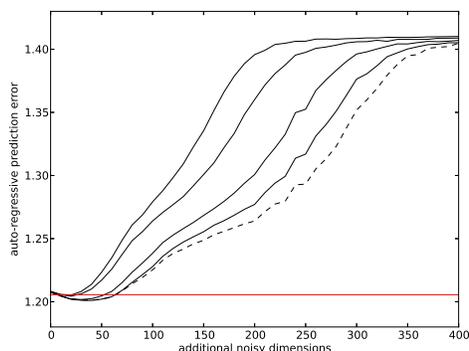


Figure 4: $k = 0, \dots, 4$ from left to right; $k = 4$ dashed; 1000 samples

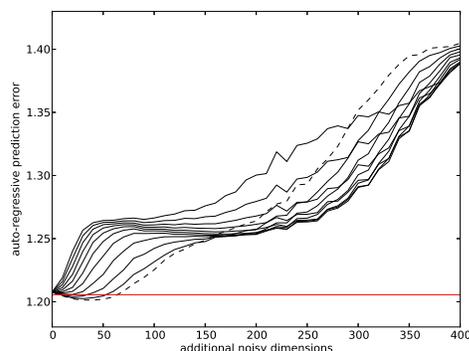


Figure 5: $k = 4, \dots, 14$ from right to left; $k = 4$ dashed; 1000 samples

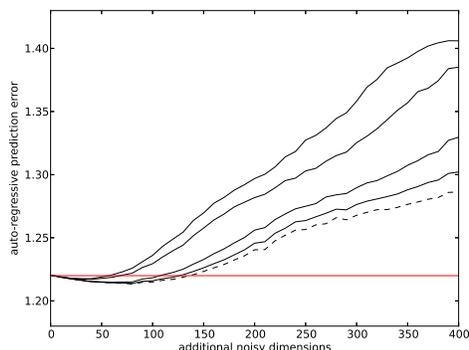


Figure 6: $k = 0, \dots, 4$ from left to right; $k = 4$ dashed; 2000 samples

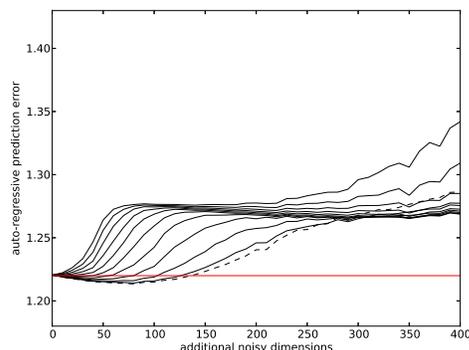


Figure 7: $k = 4, \dots, 14$ from right to left; $k = 4$ dashed; 2000 samples

Obviously $k = 4$ yields the best results in this example. The result for $k = 0$ with no additional noise can be considered to be the optimal solution. Solutions below the red line indicate overfitting conform with (19). This kind of overfitting can only be reduced by using larger samples. That it decreases for higher noise-dimension indicates that the best solution is not found any more. So we conclude that in our example, the algorithm is robust for about 50 dimensional noise if 1000 samples are used and for about 100-dimensional noise, if 2000 samples are used. Future work will include more detailed research about these relationships.

⁵Note that the vertical axis ranges from 1.15 onward to provide a better focus.

3 Criteria for suitable prediction models

In this section we discuss what properties of a prediction model are crucial to make the procedure described in 2.2 feasible.

Definition 1 (Orthogonal agnosticity criterion)

We say that a prediction-model \mathcal{G} is **orthogonal-agnostic** on Ω_t , if for every $\mathbf{A} \in O(n)$, $\mathbf{g} \in \mathcal{G}$:

$$\text{err}(\mathbf{z}) = \text{err}(\mathbf{A}^T \mathbf{z}) \quad (26)$$

(26) means that the model fits equally well to any orthogonal transformation of the data. In section 4 we will need a more restrictive variant of this criterion that additionally considers projections of the data to subspaces:

Definition 2 (Projective orthogonal agnosticity criterion)

We say that a prediction-model \mathcal{G} is **projective orthogonal-agnostic** on Ω_t , if for every $\mathbf{A} \in O(n)$, $r \leq n$ the following holds:

$$\langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{A}^T \mathbf{z}}^*(\mathbf{A}^T \text{hist}_{\mathbf{z},p})\|^2 \rangle = \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{A}_r^T \mathbf{g}_{\mathbf{z}}^*(\text{hist}_{\mathbf{z},p})\|^2 \rangle. \quad (27)$$

Note that for $r = n$, (27) simplifies to (26) and projective orthogonal agnosticity becomes equivalent to ordinary orthogonal agnosticity (since the Frobenius-norm is invariant under orthogonal transformations). An even stronger and very intuitive criterion is the following:

Definition 3 (Commuting with orthogonal transformations)

We say that a prediction-model \mathcal{G} **commutes with orthogonal transformations**, if for every $\mathbf{A} \in O(n)$ the following holds:

$$\mathbf{g}_{\mathbf{A}^T \mathbf{z}}^* = \mathbf{A}^T \mathbf{g}_{\mathbf{z}}^*. \quad (28)$$

It is rather obvious that this criterion implies projective and ordinary orthogonal agnosticity. To assure projective orthogonal agnosticity, it is a straightforward procedure to construct models such that they commute with orthogonal transformations.

Definition 4 (Information consistency criterion)

We say that a prediction-model \mathcal{G} is **information-consistent** on Ω_t , if for every $\mathbf{A} \in O(n)$, $r \leq n$ the following holds:

$$\langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{g}_{\mathbf{A}_r^T \mathbf{z}}^*(\mathbf{A}_r^T \text{hist}_{\mathbf{z},p})\|^2 \rangle \geq \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{A}_r^T \mathbf{g}_{\mathbf{z}}^*(\text{hist}_{\mathbf{z},p})\|^2 \rangle \quad (29)$$

An information-consistent model always benefits from more data rather than getting confused by it. Note that for $r = n$, (29) follows from orthogonal agnosticity.

Theorem 1

Model $\mathcal{G}_{(9)}$ is projective orthogonal agnostic and information consistent.

Proof. Projective orthogonal agnosticity follows, because the model commutes with orthogonal transformations, as $\mathbf{B}_{\mathbf{z}}(\mathbf{A}) = \mathbf{A}^T \mathbf{B}_{\mathbf{z}}(\mathbf{I}) \mathbf{A}$. To show that $\mathcal{G}_{(9)}$ is information consistent, we need the solution of the following optimization problem:

$$\underset{\mathbf{B} \in \mathbb{R}^{n \times np}}{\text{minimize}} \quad \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{B}}(\mathbf{A}^T \text{hist}_{\mathbf{z},p})\|^2 \rangle \quad (30)$$

Analytically we find a (not unique) solution to be $\mathbf{B}_z(\mathbf{A})$. Note that $\mathbf{B}_z(\mathbf{A}_r) \in \mathbb{R}^{r \times rp}$. We extend each $(r \times r)$ -block at the bottom and right with zeroes to get $(n \times n)$ -blocks and overall get an $(n \times np)$ -matrix $\mathbf{B}_z(\mathbf{A}_r)^{(n \times np)}$, which can be seen as a candidate to solve (30). Thus we have

$$\mathbf{B}_z(\mathbf{A}_r) = \mathbf{I}_r^T \mathbf{B}_z(\mathbf{A}_r)^{(n \times np)} \underline{\mathbf{I}}_r \quad (31)$$

which implies that

$$\langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{g}_{\mathbf{B}_z(\mathbf{A}_r)}(\mathbf{A}_r^T \text{hist}_{z,p})\|^2 \rangle = \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{B}_z(\mathbf{A}_r)^{(n \times np)}}(\mathbf{A}^T \text{hist}_{z,p})\|^2 \rangle \quad (32)$$

Since we know that $\mathbf{B}_z(\mathbf{A})$ is an optimal solution of (30), we have

$$\langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{B}_z(\mathbf{A}_r)^{(n \times np)}}(\mathbf{A}^T \text{hist}_{z,p})\|^2 \rangle \geq \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{B}_z(\mathbf{A})}(\mathbf{A}^T \text{hist}_{z,p})\|^2 \rangle \quad (33)$$

and thus

$$\langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{g}_{\mathbf{A}_r^T \mathbf{z}}^*(\mathbf{A}_r^T \text{hist}_{z,p})\|^2 \rangle \geq \langle \|\mathbf{A}_r^T \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{A}^T \mathbf{z}}^*(\mathbf{A}^T \text{hist}_{z,p})\|^2 \rangle \quad (34)$$

With the projective orthogonal agnosticity criterion (27), we can transform the right side of (34) into the right side of (29) and have formally shown information consistency for model (9). \square

3.1 General formulation of PFA

The criteria in section 3 assure that a problem analog to (19) can be relaxed to a tractable problem like (20) and that it can be solved like in the corresponding section. Additionally they assure that the theorem in section 4 holds and that the procedure from 2.3 is applicable. For any projective orthogonal agnostic and information consistent prediction model, (19) can be relaxed to

$$\underset{\mathbf{A} \in \mathcal{O}(n)}{\text{minimize}} \quad \langle \|\mathbf{A}_r^T (\mathbf{z} - \mathbf{g}_z^*(\text{hist}_{z,p}))\|^2 \rangle \quad (35)$$

which can be solved by diagonalizing $\langle (\mathbf{z} - \mathbf{g}_z^*(\text{hist}_{z,p})) (\mathbf{z} - \mathbf{g}_z^*(\text{hist}_{z,p}))^T \rangle$ and sorting the smallest eigenvalues to the upper left. To extend this generalization to (24), a version of (22) that only uses \mathbf{g}_z^* is needed. However this construction is straight forward, but can generally not be written as a matrix like in (22). One uses \mathbf{g}_z^* only for prediction of the first (i.e. the new) components of $\zeta(t+1)$, while the other components can be copied from $\zeta(t)$.

4 Relaxation Gap Theorem

Theorem 2

For any prediction model class \mathcal{G} that is projective orthogonal agnostic and information consistent, the following holds:

$$\text{If} \quad \exists r \leq n, \mathbf{A}^* \in \mathcal{O}(n): \text{err}(\mathbf{A}_r^{*T} \mathbf{z}) = 0 \quad (36)$$

$$\text{and} \quad \nexists \tilde{r} > r, \mathbf{A}^* \in \mathcal{O}(n): \text{err}(\mathbf{A}_{\tilde{r}}^{*T} \mathbf{z}) = 0 \quad (37)$$

$$\text{then} \quad \text{err}(\mathbf{A}_r^{(0)T} \mathbf{z}) = 0 \quad (38)$$

Line (37) has the purpose to ensure that the maximal r holding (36) is used in line (38).
 To prove the theorem, we need to setup some lemmas and the following definition:

Definition 5 (Space-partition preserving orthogonal transformations)

$$\mathbf{O}(r, s) := \text{diag}(\mathbf{O}(r), \mathbf{O}(s)) \quad (39)$$

This means that every $\tilde{\mathbf{A}} \in \mathbf{O}(r, s) \subset \mathbf{O}(r + s)$ has the form $\begin{pmatrix} \mathbf{A}_{rr} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{ss} \end{pmatrix}$ with $\mathbf{A}_{rr} \in \mathbf{O}(r)$ and $\mathbf{A}_{ss} \in \mathbf{O}(s)$. Now we can elegantly formulate the following lemma, which deals with the non-uniqueness of $\mathbf{A}^{(0)}$:

Lemma 1

Let $\mathbf{A}^{(0)}$ and $\mathbf{B}^{(0)}$ be two different global solutions of (35) and assume that the best r components are well defined (i.e. component $r + 1$ has worse error than component r).

$$\exists \tilde{\mathbf{A}} \in \mathbf{O}(r, n - r): \quad \mathbf{A}^{(0)} = \mathbf{B}^{(0)} \tilde{\mathbf{A}} \quad (40)$$

A problem would arise, if for instance the r th-worst component and the $(r + 1)$ th-worst component had equal error. In that case it would not be well defined, which signal space to extract and the lemma would not hold.

Proof of Lemma 1. As mentioned earlier, we can obtain one global solution by diagonalizing $\langle (\mathbf{z} - \mathbf{g}_z^*(\mathbf{Z})) (\mathbf{z} - \mathbf{g}_z^*(\mathbf{Z}))^T \rangle$ and sorting the r smallest eigenvalues to the upper left. Since lemma 1 requires to have a unique choice of r best components, every optimal solution must have the same set of eigenvalues in the upper left $(r \times r)$ -sub-matrix. Thus we can create every solution by orthogonally transforming the eigenspace of the r smallest eigenvalues in itself. In an analog way, the eigenspace of the $n - r$ largest eigenvalues may be transformed. The set of partition preserving orthogonal transformations $\mathbf{O}(r, n - r)$ is exactly defined to consist of the transformations performing this. \square

Lemma 2

$\forall r \leq n, \mathbf{A} \in \mathbf{O}(n), \tilde{\mathbf{A}} \in \mathbf{O}(r, n - r):$

$$\text{err}(\mathbf{I}_r^T \mathbf{A}^T \mathbf{z}) = \text{err}(\mathbf{I}_r^T (\mathbf{A} \tilde{\mathbf{A}})^T \mathbf{z}) \quad (41)$$

Lemma 2 implies that solutions of (19) stay solutions, if transformed by any $\tilde{\mathbf{A}} \in \mathbf{O}(r, n - r)$.

Proof of Lemma 2. First observe the following fact:

$$\exists \tilde{\mathbf{A}} \in \mathbf{O}(r, n - r): \quad \exists \mathbf{A}_{rr} \in \mathbf{O}(r): \quad \mathbf{I}_r^T \tilde{\mathbf{A}}^T = \mathbf{A}_{rr}^T \mathbf{I}_r^T \quad (42)$$

With (42) and the orthogonal agnosticity criterion it is straight forward to transform the right side of (41) into the left:

$$\text{err}(\mathbf{I}_r^T \tilde{\mathbf{A}}^T \mathbf{A}^T \mathbf{z}) \stackrel{(42)}{=} \text{err}(\mathbf{A}_{rr}^T \mathbf{I}_r^T \mathbf{A}^T \mathbf{z}) \stackrel{\substack{\text{orth.} \\ \text{agn.}}}{=} \text{err}(\mathbf{I}_r^T \mathbf{A}^T \mathbf{z}) \quad (43)$$

\square

Now we are ready to assemble the proof of the relaxation gap theorem:

Proof of the relaxation gap theorem. By condition (36), we have

$$\text{err}(\mathbf{A}_r^{*T} \mathbf{z}) = \left\langle \|\mathbf{A}_r^{*T} \mathbf{z} - \mathbf{g}_{\mathbf{A}_r^{*T} \mathbf{z}}^*(\mathbf{A}_r^{*T} \mathbf{Z})\|^2 \right\rangle = 0 \quad (44)$$

By information consistency, it follows that

$$\left\langle \|\mathbf{A}_r^{*T} \mathbf{z} - \mathbf{I}_r^T \mathbf{g}_{\mathbf{A}^{*T} \mathbf{z}}^*(\mathbf{A}^{*T} \mathbf{Z})\|^2 \right\rangle = 0 \quad (45)$$

So \mathbf{A}^* is a common global optimum of (19) and (35). Since (37) ensures maximality of r , we can apply lemma 1 and get:

$$\exists \tilde{\mathbf{A}} \in \mathcal{O}(r, n - r): \quad \mathbf{A}^{(0)} = \mathbf{A}^* \tilde{\mathbf{A}} \quad (46)$$

Finally by lemma 2 we conclude that $\mathbf{A}^{(0)}$ must also be an optimum of (19). \square

By continuity arguments, the implication of theorem 2 extends to signals with components of low error – the lower the error is, the more precise we can find an optimum of (19) by solving (35). Providing bounds for the steepness of this continuous relationship is still an open problem and may be subject of future work.

5 Future Work

An important aspect of our future work will be the application of PFA to real world problems. We plan to approach scenarios where SFA is known to produce good results, so we can compare PFA and SFA and get a clearer notion for the differences between the paradigms. On the other hand, we have new scenarios in mind, that specifically would benefit from predictable features. For instance we are working on applications related to robotic navigation and the lane-keeping problem of a simulated car.

As mentioned in some sections of this document, another fork of our future work will be to extend the analytic understanding of the heuristic aspects of the algorithm. This way we also aim to improve our methods to avoid overfitting.

Acknowledgments

This work is funded by a grant from the German Research Foundation (Deutsche Forschungsgemeinschaft, DFG) to L. Wiskott (SFB 874, TP B3) and supported by the German Federal Ministry of Education and Research within the National Network Computational Neuroscience - Bernstein Fokus: “Learning behavioral models: From human experiment to technical assistance”, grant FKZ 01GQ0951.

References

- [1] W. Bialek, I. Nemenman, and N. Tishby. Predictability, complexity, and learning. *Neural Comput*, 13:2409–2463, Nov 2001.
- [2] S. Dähne, N. Wilbert, and L. Wiskott. Self-organization of v1 complex cells based on slow feature analysis and retinal waves. In *Proc. Bernstein Conference on Computational Neuroscience, Sep 27–Oct 1, Berlin, Germany*, 2010.
- [3] Mathias Franzius, Niko Wilbert, and Laurenz Wiskott. Invariant object recognition and pose estimation with slow feature analysis. *Neural Computation*, 23(9):2289–2323, 2011.
- [4] Mathias Franzius, Henning Sprekeler, and Laurenz Wiskott. Slowness and sparseness lead to place, head-direction, and spatial-view cells. *PLoS Computational Biology*, 3(8):e166, 2007.
- [5] Laurenz Wiskott. Estimating driving forces of nonstationary time series with slow feature analysis. arXiv.org e-Print archive, 2003.
- [6] Tobias Blaschke, Tiziano Zito, and Laurenz Wiskott. Independent slow feature analysis and nonlinear blind source separation. *Neural Computation*, 19(4):994–1021, 2007.
- [7] Felix Creutzig, Amir Globerson, and Naftali Tishby. Past-future information bottleneck in dynamical systems. *Physical Review E*, 79:041925, 2009.
- [8] Felix Creutzig and Henning Sprekeler. Predictive coding and the slowness principle: An information-theoretic approach. *Neural Computation*, 20(4):1026–1041, 2008.
- [9] Alexander Gepperth. Simultaneous concept formation driven by predictability. In *ICDL-EPIROB’12*, pages 1–6, 2012.
- [10] Aapo Hyvärinen. Complexity pursuit: Separating interesting components from time series. *Neural Computation*, 13(4):883–898, 2001.
- [11] Georg Goerg. Forecastable component analysis. In Sanjoy Dasgupta and David Mcallester, editors, *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, volume 28, pages 64–72. JMLR Workshop and Conference Proceedings, May 2013.
- [12] Weghenkel Richthofer and Wiskott. Predictable feature analysis. In *Proc. Bernstein Conference on Computational Neuroscience, Sep 12–14, Munich, Germany*, 2012. Special issue of Frontiers in Computational Neuroscience 120.
- [13] L. Wiskott, P. Berkes, M. Franzius, H. Sprekeler, and N. Wilbert. Slow feature analysis. *Scholarpedia*, 6(4):5282, 2011.
- [14] G. E. P. Box and G. C. Tiao. A canonical analysis of multiple time series. *Biometrika*, 64(2):pp. 355–365, 1977.

A Appendix

A.1 Notation overview

This section gives an overview of the notation used in this paper.

$\mathbf{x}(t)$	denotes the raw input signal and might only be available for a discrete sequence of t 's.
Ω_t	$:= \{t_0, \dots, t_k\}$ denotes a discrete time sequence (considered as equidistant with step size normalized to 1). We usually refer to Ω_t as the <i>training phase</i> .
$\langle \mathbf{s}(t) \rangle_{t \in S}$	$:= \frac{1}{ S } \sum_{t \in S} \mathbf{s}(t)$ denotes the average of some signal \mathbf{s} over a finite set S . For $S = \Omega_t$ we just write $\langle \mathbf{s}(t) \rangle_t$ or even $\langle \mathbf{s} \rangle$, if it is obvious, what unbound variable is targeted.
$\mathbf{h}(\mathbf{x})$	denotes the expansion function and usually consists of a set of monomials of low degree.
$\mathbf{z}(t)$	denotes $\mathbf{h}(\mathbf{x}(t))$ after sphering it.
$\mathbf{m}(t)$	denotes the optimized output signal (\mathbf{m} for <i>model</i>).
n	denotes the number of components to be analyzed (after expansion).
r	denotes the number of extracted components ("features").
\mathbf{A}, \mathbf{a}	denotes the matrix (or vector if $r = 1$) holding the linear composition of the output-signal. We set $\mathbf{m}(t) = \mathbf{A}^T \mathbf{z}(t)$.
\mathbf{a}_i	denotes the i 'th column of \mathbf{A} , so we can write $m_i(t) = \mathbf{a}_i^T \mathbf{z}(t)$.
$O(n)$	$\subset \mathbb{R}^{n \times n}$ denotes the orthogonal group of dimension n , i.e. $\forall \mathbf{A} \in O(n): \mathbf{A}\mathbf{A}^T = \mathbf{A}^T\mathbf{A} = \mathbf{I}$
p	denotes the number of recent signal-values involved in the prediction. We also call it the <i>prediction-order</i> .
$\mathbf{I}_{s,r}$	denotes the $s \times r$ identity matrix (s counting rows, r counting columns). For $s = r$ this is a usual square identity, while in the non-square case it consists of a square identity block in the top or left area, filled up with zeroes to fit the given shape.
\mathbf{I}_r	$:= \mathbf{I}_{n,r}$
\mathbf{A}_r	$:= \mathbf{A}\mathbf{I}_r$

We frequently use the p -step time-history of a signal \mathbf{z} , which we formalize by the following function:

$$\text{hist}_{\mathbf{z},p,\Delta}(t) := \sum_{i=1}^p \mathbf{z}(t - i\Delta) \mathbf{e}_i^T \quad \text{with } \mathbf{e}_i \in \mathbb{R}^p \quad (47)$$

$$\text{hist}_{\mathbf{z},p}(t) := \text{hist}_{\mathbf{z},p,1}(t) \quad (48)$$

Here \mathbf{e}_i denotes the i -th p -dimensional euclidean unit vector, which is 1 at position i and 0 everywhere else.

Further more we sometimes use the Kronecker product \otimes and the vec-operator defined as follows:

For matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{k \times l}$ and with a_{ij} denoting the entries, \mathbf{a}_i the columns of \mathbf{A} :

$$\mathbf{A} \otimes \mathbf{B} := \begin{pmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{pmatrix} \in \mathbb{R}^{mk \times nl} \quad (49)$$

$$\text{vec}(\mathbf{A}) := \begin{pmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{pmatrix} \in \mathbb{R}^{mn} \quad (50)$$

Additionally, we sometimes make use of the following shortcut:

$$\underline{\mathbf{A}} := \mathbf{I}_{p,p} \otimes \mathbf{A} = \underbrace{\begin{pmatrix} \mathbf{A} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{A} \end{pmatrix}}_{p \text{ times } \mathbf{A}} \quad (51)$$

A.2 Extracting predictable single components

In section 2.2 we initially stated a prediction model that always scopes on single components. This idea was not suitable for PFA because it contradicts the orthogonal agnosticity criterion. In this section we propose a strategy to extract well predictable single components even though. We begin by recalling our initial notion of linear auto regressive predictability:

$$\mathbf{a}^T \mathbf{z}(t) \stackrel{!}{\approx} b_1 \mathbf{a}^T \mathbf{z}(t-1) + \dots + b_p \mathbf{a}^T \mathbf{z}(t-p) \quad (52)$$

$$= \mathbf{a}^T \text{hist}_{\mathbf{z},p}(t) \mathbf{b} \quad (53)$$

It is possible to write this for multiple dimensions by constraining the coefficient-matrices to be diagonal:

$$\mathbf{m}(t) \stackrel{!}{\approx} \mathbf{B}_1 \mathbf{m}(t-1) + \dots + \mathbf{B}_p \mathbf{m}(t-p) \quad \text{with } \mathbf{B}_i \in \mathbb{R}^{n \times n}, \text{ diagonal} \quad (54)$$

This model is not orthogonal agnostic, so a different approach than in section 2.2 is needed. To minimize the least-squares-error of (52), the following optimization problem needs to be solved:

$$\begin{aligned} & \underset{\mathbf{a} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}^p}{\text{minimize}} && \left\langle (\mathbf{a}^T (\mathbf{z} - \text{hist}_{\mathbf{z},p} \mathbf{b}))^2 \right\rangle \\ & \text{subject to} && \mathbf{a}^T \langle \mathbf{z} \rangle = 0 \quad (\text{zero mean}) \\ & && \mathbf{a}^T \langle \mathbf{z} \mathbf{z}^T \rangle \mathbf{a} = 1 \quad (\text{unit variance}) \end{aligned} \quad (55)$$

Via analytic optimization it is straight forward to find the optimal \mathbf{a} , if \mathbf{b} is fixed and vice versa:

If \mathbf{b} is fixed, choose \mathbf{a} as the eigenvector corresponding to the smallest eigenvalue in

$$\langle \mathbf{z} \mathbf{z}^T \rangle - \langle \mathbf{z} \mathbf{b}^T \text{hist}_{\mathbf{z},p} \rangle - \langle \text{hist}_{\mathbf{z},p}^T \mathbf{b} \mathbf{z}^T \rangle + \langle \text{hist}_{\mathbf{z},p} \mathbf{b} \mathbf{b}^T \text{hist}_{\mathbf{z},p}^T \rangle \quad (56)$$

If \mathbf{a} is fixed, choose \mathbf{b} as

$$\mathbf{b}^T := \langle \mathbf{z}^T \mathbf{a} \mathbf{a}^T \text{hist}_{\mathbf{z},p} \rangle \langle \text{hist}_{\mathbf{z},p}^T \mathbf{a} \mathbf{a}^T \text{hist}_{\mathbf{z},p} \rangle^{-1} \quad (57)$$

By inserting (57) into (55) one could obtain a problem written in \mathbf{a} only:

$$\begin{aligned} & \underset{\mathbf{a} \in \mathbb{R}^n}{\text{minimize}} && \left\langle \left(\mathbf{a}^T (\mathbf{z} - \text{hist}_{\mathbf{z},p} \langle \mathbf{z}^T \mathbf{a} \mathbf{a}^T \text{hist}_{\mathbf{z},p} \rangle \langle \text{hist}_{\mathbf{z},p}^T \mathbf{a} \mathbf{a}^T \text{hist}_{\mathbf{z},p} \rangle^{-1}) \right)^2 \right\rangle \\ & \text{subject to} && \mathbf{a}^T \langle \mathbf{z} \rangle = 0 \quad (\text{zero mean}) \\ & && \mathbf{a}^T \langle \mathbf{z} \mathbf{z}^T \rangle \mathbf{a} = 1 \quad (\text{unit variance}) \end{aligned} \quad (58)$$

Problem (58) is not efficiently globally solvable by any method known to us, which is mainly due to the occurrence of \mathbf{a} in a matrix-term under an inversion-symbol. However a possible strategy is to approximate the solution by choosing an initial value for \mathbf{a} or \mathbf{b} and applying (56) and (57) in turns until a stable state is reached.

As a reasonable initial value for this procedure we choose \mathbf{b} such that it is the best predictor of \mathbf{z} on average, in absence of any \mathbf{a} :

$$\mathbf{z}(t) \stackrel{!}{\approx} b_1 \mathbf{z}(t-1) + \dots + b_p \mathbf{z}(t-p) = \text{hist}_{\mathbf{z},p}(t) \mathbf{b} \quad (59)$$

To minimize the error of (59) on average over all components of \mathbf{z} , we propose the following least-squares optimization:

$$\underset{\mathbf{b} \in \mathbb{R}^p}{\text{minimize}} \quad \langle (\mathbf{z} - \text{hist}_{\mathbf{z},p} \mathbf{b})^T (\mathbf{z} - \text{hist}_{\mathbf{z},p} \mathbf{b}) \rangle \quad (60)$$

The solution of this problem is

$$\mathbf{b} := \langle \mathbf{z}^T \text{hist}_{\mathbf{z},p} \rangle \langle \text{hist}_{\mathbf{z},p}^T \text{hist}_{\mathbf{z},p} \rangle^{-1} \quad (61)$$

Solution (61) does not change, if we replace \mathbf{z} by $\mathbf{A}^T \mathbf{z}$ with any orthogonal, full ranked \mathbf{A} . However, one quickly finds examples, where the procedure stabilizes in sub-optimal states. Though one can partly overcome this issue by estimating better starting points, the method still has unknown success-probability.

Probably a better possibility is to solve (55) with PFA as described in section 2 for $r = 1$.

After extracting one component either way, one can project \mathbf{z} to the signal space uncorrelated (i.e. orthogonal) to the extracted component. The extraction- and projection-procedure can be repeated until any desired number of components is extracted.