

Non-linear, Sparse Dimensionality Reduction via Path Lasso Penalized Autoencoders

Oskar Allerbo

ALLERBO@CHALMERS.SE

Mathematical Sciences

University of Gothenburg and Chalmers University of Technology

SE-412 96 Gothenburg, Sweden

Rebecka Jörnsten

JORNSTEN@CHALMERS.SE

Mathematical Sciences

University of Gothenburg and Chalmers University of Technology

SE-412 96 Gothenburg, Sweden

Editor: Pradeep Ravikumar

Abstract

High-dimensional data sets are often analyzed and explored via the construction of a latent low-dimensional space which enables convenient visualization and efficient predictive modeling or clustering. For complex data structures, linear dimensionality reduction techniques like PCA may not be sufficiently flexible to enable low-dimensional representation. Non-linear dimension reduction techniques, like kernel PCA and autoencoders, suffer from loss of interpretability since each latent variable is dependent of all input dimensions. To address this limitation, we here present path lasso penalized autoencoders. This structured regularization enhances interpretability by penalizing each *path* through the encoder from an input to a latent variable, thus restricting how many input variables are represented in each latent dimension. Our algorithm uses a group lasso penalty and non-negative matrix factorization to construct a sparse, non-linear latent representation. We compare the path lasso regularized autoencoder to PCA, sparse PCA, autoencoders and sparse autoencoders on real and simulated data sets. We show that the algorithm exhibits much lower reconstruction errors than sparse PCA and parameter-wise lasso regularized autoencoders for low-dimensional representations. Moreover, path lasso representations provide a more accurate reconstruction match, i.e. preserved relative distance between objects in the original and reconstructed spaces.

Keywords: Sparse Dimensionality Reduction, Non-linear Dimensionality Reduction, Regularized Neural Networks, Group Lasso, Autoencoders

1. Introduction

Dimensionality reduction is a key component in data compression, data visualization and feature extraction. One of the most widely used techniques is principal component analysis (PCA), that uses the eigendecomposition of the sample covariance matrix to construct latent dimensions as linear combinations of the original dimensions. The interpretability of the latent representation is increased if each latent dimension consists only of a subset of the original dimensions, as in sparse PCA (Zou et al., 2006). Since the introduction of sparse PCA, several variants have been presented, such as non-negative sparse PCA (Zass

and Shashua, 2007), where the loadings are all non-negative; multilinear sparse PCA (Lai et al., 2014), that operates on tensors instead of vectors; and robust PCA (Meng et al., 2012; Croux et al., 2013), that is less affected by outliers.

While more interpretable than standard PCA, sparse PCA and its variants are still linear and therefore cannot capture more complex relations in the data. Furthermore, they have limitations on how efficiently they make use of the latent dimensions, something that is especially important when the number of latent dimensions is small. Several non-linear generalizations of PCA exist, the most important ones being kernel PCA (Schölkopf et al., 1998) and autoencoders (Kramer, 1991). In kernel PCA, a kernel function is used to implicitly and non-linearly map the data to a space with higher dimensionality than d_x (the original dimensionality) in which linear PCA is performed. Autoencoders use an hourglass shaped neural network with d_x input and output nodes and d_z (the latent dimensionality) nodes in the middle layer. The same data is used both as input and output, so the goal of the autoencoder is to reconstruct the original data from a lower dimensional representation, that is found in the middle layer.

Although there do exist several algorithms for sparse kernel PCA, see work by Tipping (2001), Smola et al. (2002), Wang and Tanaka (2016) and Guo et al. (2019), as well as for sparse autoencoders (e.g. Ng et al., 2011), the terminology might be a little confusing, since the algorithms are sparse in a different sense than sparse PCA. Instead of being sparse in the sense original to latent dimensions, they are sparse in the sense observations to latent dimensions, which means that each observation is active only in a subset of the latent dimensions (and vice versa, each latent dimension depends only on a subset of the data). To further illustrate this, we look at the linear case, with d_x and d_z as before and with n being the number of observations. For $\mathbf{Z} \in \mathbb{R}^{n \times d_z}$, $\mathbf{X} \in \mathbb{R}^{n \times d_x}$ and $\mathbf{W} \in \mathbb{R}^{d_x \times d_z}$, the latent representation \mathbf{Z} is obtained from the original representation \mathbf{X} as

$$\mathbf{Z} = \mathbf{X} \cdot \mathbf{W}. \tag{1}$$

Then sparsity in the sense original to latent dimensions means that \mathbf{W} is sparse, while sparsity in the sense observations to latent dimensions means that \mathbf{Z} is sparse. A third notion of sparsity is feature selection, where only a subset of the original dimension are included in the model. This corresponds to entire rows of \mathbf{W} being zero, allowing the remaining rows to be dense.

The sparse autoencoder generalizes Equation (1) to $\mathbf{Z} = \text{Enc}(\mathbf{X})$, where Enc is the non-linear encoder and, to produce a sparse \mathbf{Z} , it includes a constraint on how frequently a latent unit is allowed to be non-zero, where frequency is measured among observations. This means that the architecture of the encoder might still be dense, in contrast to sparsity corresponding to a sparse \mathbf{W} , which requires a sparse architecture.

Neural networks with sparse architectures are usually obtained using different versions of lasso, or l_1 -, regularization (Tibshirani, 1996), and there are numerous examples of these. Scardapane et al. (2017) used group lasso (Yuan and Lin, 2006) to remove all the links to or from a node, while Yoon and Hwang (2017) combined group lasso with exclusive lasso (Zhou et al., 2010) on filters in convolutional neural networks to, on one hand, totally eliminate some filters (using group lasso) and, on the other hand, make filters as different as possible (using exclusive lasso). Lasso regularized autoencoders include work by Wu et al. (2019), with a linear encoder and a lasso regularized decoder; Dabin et al. (2020), with a lasso

regularized encoder and a linear decoder; and Ainsworth et al. (2018), which uses group lasso and variational autoencoders to split the original dimensions into pre-defined groups and then uses one decoder per group, with a shared latent space.

In this paper we propose path lasso, that uses group lasso regularization to eliminate all connections between two nodes in two non-adjacent layers of a fully connected feedforward neural network. We apply path lasso, in combination with exclusive lasso, to an autoencoder to introduce sparsity between original and latent variables, obtaining non-linear, sparse dimensionality reduction in the same sense as in sparse PCA. Path lasso forces each latent dimension to be a function only of a subset of the original dimensions, while exclusive lasso encourages these subsets to differ. To the best of our knowledge this is the first non-linear dimensionality reduction algorithm that is sparse in this sense.

The rest of this paper is organized as follows: In Section 2 we introduce the path lasso penalty and the path lasso penalized autoencoder. In Section 3 we run experiments on real and simulated data sets, comparing path lasso to PCA, autoencoders, sparse autoencoders, sparse PCA and an autoencoder with parameter-wise l_1 -regularization. We show that, for a given sparsity, path lasso results in a lower reconstruction error and is better at reconstruction match, i.e. retaining relative positioning of objects in the reconstructed space as in the original space. We conclude with a discussion in Section 4.

2. Method

This section is structured in the following way: Sections 2.1 and 2.2 present short reviews of different flavours of the lasso algorithm and of proximal gradient descent, while Sections 2.3 to 2.6 describe different aspects of path lasso. Section 2.3 and 2.4 describe how paths between nodes in two non-adjacent layers are defined and penalized, and how the path penalties are transformed to individual link penalties. Section 2.5 discusses when and how to apply the path penalty and Section 2.6 describes how we adapt path lasso when using it in an autoencoder. In Appendix A we discuss different methods to accelerate training.

2.1 Review of Lasso Penalties

The lasso algorithm sets some model parameters exactly equal to zero, thus eliminating them from the model. There are different versions of the lasso, four of which are used in this paper. Here follows a short summary to these.

(Standard) Lasso applies an l_1 -penalty to all the parameters in the parameter vector $\boldsymbol{\theta} \in \mathbb{R}^d$, and is defined as

$$\lambda \|\boldsymbol{\theta}\|_1 = \lambda \sum_{i=1}^d |\theta_i|,$$

where $\lambda > 0$ is the regularization strength. Due to the non-differentiability of the absolute value at zero, some of the θ_i 's are set to exactly zero and thus eliminated from the model.

Adaptive Lasso applies an individual l_1 -penalty to all the parameters in the parameter vector $\boldsymbol{\theta} \in \mathbb{R}^d$, and is defined as

$$\sum_{i=1}^d \lambda_i |\theta_i|,$$

where $\lambda_i := \frac{\lambda}{|\hat{\theta}_i^R|^\gamma}$ for some $\gamma > 0$ (common practice is $\gamma = 2$), and $\hat{\theta}_i^R$ is the ridge regression estimate of θ_i . The idea is that important parameters will have larger values of $\hat{\theta}_i^R$ and thus be penalized less than unimportant parameters.

Group Lasso penalizes pre-defined groups of parameters together, which means that either all, or none, of the parameters in the group are set to zero. The group lasso penalty for a group $g \in \mathcal{G}$ is defined as

$$\lambda \|\boldsymbol{\theta}_g\|_2 = \lambda \sqrt{\sum_{i \in g} \theta_i^2},$$

where $\mathcal{G} = \{g_1, \dots, g_G\}$ is a disjoint partition of the index set $\{1, \dots, d\}$, i.e. each g is a set of indices defining a group, and, for a given $\boldsymbol{\theta} \in \mathbb{R}^d$, $\boldsymbol{\theta}_g$ is a d -dimensional vector with components equal to $\boldsymbol{\theta}$ for indices within g and zero otherwise. The total group lasso penalty is then taken as the sum of the penalties over the different groups. As seen, for a given group g ,

$$\sqrt{\sum_{i \in g} \theta_i^2} = 0 \iff \theta_i = 0 \forall i \in g.$$

The *Exclusive Lasso* can be seen as the opposite of the group lasso. Again, the parameters are split into pre-defined groups, but now the goal is instead to impose a similar number of non-zero parameters in every group. With g as before, its exclusive lasso penalty defined as

$$\lambda \|\boldsymbol{\theta}_g\|_1^2 = \lambda \left(\sum_{i \in g} |\theta_i| \right)^2,$$

and, again, the total penalty is defined as the sum over the groups. Since the number of mixed terms in the squared sum grows with the number of elements in the sum, the total number of mixed terms over all the groups is minimized when the non-zero elements are evenly distributed among the groups.

2.2 Review of Proximal Gradient Descent

The reason that some parameters in lasso penalized models are set exactly to zero, is that the derivative of the absolute value is not unique at zero: $\frac{\partial|\theta|}{\partial\theta} \Big|_{\theta=0} \in [-1, 1]$. Gradient descent methods are unable to use this non-uniqueness and as a consequence no parameters are set exactly to zero. Proximal gradient descent (Rockafellar, 1976) on the other hand takes the non-uniqueness into account, resulting in exact zeros.

If the objective function, $f(\boldsymbol{\theta})$, can be decomposed into $f(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) + h(\boldsymbol{\theta})$, where g is differentiable (typically a reconstruction error) and h is not (typically a lasso regularization term), then a standard gradient descent step, followed by a proximal gradient descent step, is defined as

$$\boldsymbol{\theta}^{t+1} = \text{prox}_{\alpha h}(\boldsymbol{\theta}^t - \alpha \nabla g(\boldsymbol{\theta}^t)),$$

where $\alpha > 0$ is the learning rate and prox is the proximal operator that depends on h .

For lasso, with $h(\boldsymbol{\theta}) = \lambda \|\boldsymbol{\theta}\|_1 = \sum_i \lambda |\theta_i|$, the proximal operator decomposes component-wise and is, with $(x)^+ := \max(x, 0)$,

$$\text{prox}_{\alpha h}(\theta_i) = \text{sign}(\theta_i) \cdot (|\theta_i| - \alpha \lambda)^+. \tag{2}$$

I.e., each θ_i is additively shrunk towards zero, and once it changes sign it becomes exactly zero.

For group lasso, with $h(\boldsymbol{\theta}) = \lambda \sum_{g \in \mathcal{G}} \|\boldsymbol{\theta}_g\|_2 = \lambda \sum_{g \in \mathcal{G}} \sqrt{\sum_{i \in g} \theta_i^2}$, where $\boldsymbol{\theta}_g$ are the θ_i 's that belong to group g and \mathcal{G} is the set of all groups, the proximal operator for θ_i is

$$\text{prox}_{\alpha h}(\theta_i) = \theta_i \cdot \left(1 - \frac{\alpha \lambda}{\sqrt{\sum_{j \in g_i} \theta_j^2}} \right)^+, \quad (3)$$

where g_i is the group that θ_i belongs to. Thus all members of the group are penalized equally and set to zero at the exact same time.

2.3 Path Penalties

Let $\{\mathbf{o}_l\}_{l=0}^L$, $\mathbf{o}_l \in \mathbb{R}^{d_l}$, denote the outputs of $L + 1$ consecutive layers in a fully connected feedforward neural network, where the first and last layers are also denoted by \mathbf{x} and \mathbf{y} respectively, i.e. $\mathbf{x} := \mathbf{o}_0$ and $\mathbf{y} := \mathbf{o}_L$. Then \mathbf{y} depends on \mathbf{x} as

$$\mathbf{y} = \Phi_L(\mathbf{W}_L \Phi_{L-1}(\mathbf{W}_{L-1} \Phi_{L-2}(\dots \Phi_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \dots) + \mathbf{b}_{L-1}) + \mathbf{b}_L), \quad (4)$$

where $\{\mathbf{W}_l\}_{l=1}^L$, $\mathbf{W}_l \in \mathbb{R}^{d_l \times d_{l-1}}$, are the weight matrices, $\{\mathbf{b}_l\}_{l=1}^L$, $\mathbf{b}_l \in \mathbb{R}^{d_l}$, the bias vectors, and $\{\Phi_l\}_{l=1}^L$ the (not necessarily identical) element-wise activation functions. Thinking of Equation (4) as a graph, each weight matrix element, $w_{i_1 i_2}^l := (\mathbf{W}_l)_{i_1 i_2}$, corresponds to a link between two nodes in two adjacent layers

By combining links from multiple weight matrices, paths between nodes in non-adjacent layers can be constructed. Between a given node in layer \mathbf{x} , x_{i_0} , and a node in layer \mathbf{y} , y_{i_L} , there are in total $\prod_{l=1}^{L-1} d_l$ paths, each consisting of L links, see Figure 1 for an illustration. A path is broken if at least one of its links has value zero and to disconnect the two nodes, all paths between them need to be broken. Just as Neyshabur et al. (2015), we define the value of a path as the product of its absolute valued links; see Definition 1.

Definition 1 (Path Value) For $k = 1, 2, \dots, \prod_{l=1}^{L-1} d_k$, where each k corresponds to a unique combination of the indices i_1 to i_{L-1}

$$p_{i_L i_0}^k := |w_{i_L i_{L-1}}^L| \cdot |w_{i_{L-1} i_{L-2}}^{L-1}| \cdot \dots \cdot |w_{i_1 i_0}^1|.$$

With this definition a broken path, where at least one link is zero, has the value zero. We further define a group of paths so that all paths connecting x_{i_0} to y_{i_L} form one group; then if all paths in the group are zero, the nodes are disconnected. Applying the group lasso proximal operator to path $p_{i_L i_0}^k$, belonging to group $g_{i_L i_0}$ then amounts to

$$p_{i_L i_0}^k \cdot \left(1 - \frac{\alpha \lambda}{\sqrt{\sum_{l \in g_{i_L i_0}} (p_{i_L i_0}^l)^2}} \right)^+, \quad (5)$$

which according to Proposition 3 can be written as

$$p_{i_L i_0}^k \cdot \left(1 - \frac{\alpha \lambda}{(\mathbf{W}_{\text{PL}})_{i_L i_0}} \right)^+$$

with \mathbf{W}_{PL} according to Definition 2.

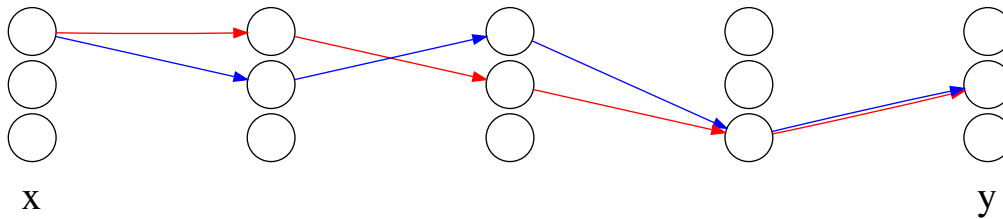


Figure 1: Illustration of two of the 3^3 (where the exponent is the number of inner layers and the base is their width) possible paths between node one in layer one and node two in layer five. Each arrow denotes a link, which corresponds to an element in a weight matrix.

Definition 2 (Path Lasso Matrix) *With the square and the square root taken element-wise,*

$$\mathbf{W}_{PL} := \sqrt{\prod_{l=L,\dots,1} (\mathbf{W}_l)^2} = \sqrt{(\mathbf{W}_L)^2 \cdot (\mathbf{W}_{L-1})^2 \cdot \dots \cdot (\mathbf{W}_1)^2}. \quad (6)$$

Proposition 3 *The element (i_L, i_0) in \mathbf{W}_{PL} is the group lasso penalty on the group consisting of all paths between nodes x_{i_0} and y_{i_L} , i.e.*

$$(\mathbf{W}_{PL})_{i_L i_0} = \sqrt{\sum_{k \in g_{i_L i_0}} (p_{i_L i_0}^k)^2}$$

For a proof, see Appendix B.

We call each element in \mathbf{W}_{PL} a connection, i.e. $(\mathbf{W}_{PL})_{i_L i_0}$ is the strength of the connection between nodes x_{i_0} and y_{i_L} and if it is zero, meaning that all paths between the nodes are broken, then the nodes are disconnected. Proposition 4 gives a theoretical justification for this definition of connections, stating that if there is no connection between x_{i_0} and y_{i_L} , then the derivative of y_{i_L} with respect to x_{i_0} is zero, regardless of the value of \mathbf{x} .

Proposition 4 *Let the vector \mathbf{y} depend on the vector \mathbf{x} as stated in Equation (4) and let \mathbf{W}_{PL} be the path lasso matrix defined in Equation (6). Then, if all weights and activations are bounded,*

$$(\mathbf{W}_{PL})_{i_L i_0} = 0 \implies \frac{\partial y_{i_L}}{\partial x_{i_0}} = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{i_L i_0} = 0 \quad \forall \mathbf{x} \in \mathbb{R}^{d_0}.$$

For a proof, see Appendix B.

2.4 From Paths to Links

Equation (5) describes how all paths between two nodes in the network are penalized towards zero, but it does not tell how this penalization affects the individual links. However, since the neural network is expressed in terms of the individual links, rather than in paths, the penalized paths must be translated into penalized links, i.e. the penalty must be expressed on link level, rather than on path level. Proposition 5 describes how this translation can be done, i.e. how applying the proximal group lasso operator to the paths can be translated into applying a standard lasso proximal operator to each individual link.

Proposition 5 *The group lasso proximal step in path space can be transformed to standard lasso proximal steps in link space by first solving the matrix equation*

$$\prod_{l=L,\dots,1} |(\mathbf{W}_l)^t| \odot \left(1 - \frac{\alpha\lambda}{(\mathbf{W}_{PL})^t}\right)^+ = \prod_{l=L,\dots,1} \tilde{\mathbf{W}}_l \quad (7)$$

for $\{\tilde{\mathbf{W}}_l\}_{l=1}^L$, and then set $(\mathbf{W}_l)^{t+1} := \text{sign}((\mathbf{W}_l)^t) \odot \tilde{\mathbf{W}}_l$, where $\tilde{w}_{ij}^l =: (|w_{ij}^l|^t - \alpha\lambda_{ij}^{l,t})^+$ for some $\lambda_{ij}^{l,t} > 0$.

The absolute value, sign and division are taken element-wise and \odot denotes element-wise multiplication.

The proof, which is presented in Appendix B, contains a step where the paths in each group are summed over. This step transforms the system of non-linear equations from one equation per path to one equation per connection, i.e. from $\prod_{l=0}^L d_l$ to $d_0 \cdot d_L$ equations, and to a form can be solved efficiently using non-negative matrix factorization, NMF, with the extra requirement that $\tilde{w}_{ii_{l-1}}^l \leq |(w_{ii_{l-1}}^l)^t|$, or equivalently $\lambda_{ii_{l-1}}^{l,t} \geq 0$. We describe this algorithm in Appendix C.

The reduced number of equations leads to an undetermined system (unless the hidden layers are very narrow). Therefore, all equations can hold, although there might be more than one solution. Since we are interested in a solution that lies close to the unpenalized weight matrices, we use $|\mathbf{W}_l|$ as the seed for each $\tilde{\mathbf{W}}_l$ in the matrix factorization.

An optimization step using proximal gradient descent on paths is summarized in Algorithm 1. The bottleneck of this algorithm is the matrix factorization step from Equation (7). In Appendix A, different methods to alleviate this bottleneck are discussed.

2.5 Applying the Path Lasso Penalty

Since training a neural network is a non-convex optimization problem, with multiple local optima, how and when the regularization is added might affect which optimum is found. If a too high penalty is added too early during training, the risk of getting stuck in a bad local optimum is larger than if the regularization is added later. To mitigate this, training was split into three stages:

- Following the adaptive lasso approach of Allerbo and Jörnsten (2020), we first trained the network without path regularization to obtain individual penalties for each connection during the path lasso stage.

Algorithm 1 Proximal Path Lasso Optimization Step

Input: Parameters at time t , $\{(\mathbf{W}_l)^t, (\mathbf{b}_l)^t\}_{l=1}^L$; data, \mathbf{x} ; learning rate, α ; regularization strength, λ .

Output: Parameters at time $t + 1$: $\{(\mathbf{W}_l)^{t+1}, (\mathbf{b}_l)^{t+1}\}_{l=1}^L$.

- 1: Update all weights and biases using one step of standard (stochastic) gradient descent:

$$\{(\mathbf{W}_l)^{t+\frac{1}{2}}, (\mathbf{b}_l)^{t+\frac{1}{2}}\}_{l=1}^L \leftarrow \{(\mathbf{W}_l)^t, (\mathbf{b}_l)^t\}_{l=1}^L - \alpha \cdot \nabla (\text{NN}(\{(\mathbf{W}_l)^t, (\mathbf{b}_l)^t\}_{l=1}^L; \mathbf{x})),$$

where NN denotes the neural network.

- 2: Penalize paths, by applying the group lasso proximal operator, and update the path values accordingly:

- a: Construct the path lasso penalty matrix for the weight outputs from step 1, according to Equation (6):

$$(\mathbf{W}_{\text{PL}})^{t+\frac{1}{2}} \leftarrow \sqrt{\prod_{l=L, \dots, 1} \left((\mathbf{W}_l)^{t+\frac{1}{2}} \right)^2}.$$

- b: Penalize paths according to Equation (7):

$$\mathbf{P}^{t+1} \leftarrow \prod_{l=L, \dots, 1} |(\mathbf{W}_l)^{t+\frac{1}{2}}| \odot \left(1 - \frac{\alpha \lambda}{(\mathbf{W}_{\text{PL}})^{t+\frac{1}{2}}} \right)^+.$$

- 3: Translate penalized paths to penalized links:

- a: Translate penalized paths to absolute valued penalized links, using modified non-negative matrix factorization:

$$\{(\tilde{\mathbf{W}}_l)^{t+1}\}_{l=1}^L \leftarrow \text{NMF}(\mathbf{P}^{t+1}).$$

- b: Restore signs:

$$(\mathbf{W}_l)^{t+1} \leftarrow \text{sign}((\mathbf{W}_l)^t) \odot (\tilde{\mathbf{W}}_l)^{t+1}.$$

- In a second stage, we added path regularization with an individual penalty to each connection, depending on the magnitude of the connection after the first stage:

$$\lambda_{i_L, i_0} := \frac{\lambda}{((\hat{\mathbf{W}}_{\text{PL}})_{i_L i_0})^\gamma},$$

where $\hat{\mathbf{W}}_{\text{PL}}$ is the value of \mathbf{W}_{PL} after the first optimization stage and $\gamma > 0$. Throughout this paper, $\gamma = 2$ was used.

- To reduce bias and thus improve performance, we finally added a stage of unregularized training after the path lasso stage, with the links set to zero in the previous stage kept to zero.

All stages were trained until convergence and stages two and three were warm started with the solutions from the previous stage.

2.6 Path Lasso for Dimensionality Reduction

Path lasso as described so far is applicable to any two non-adjacent layers in any feedforward neural network. To use it for sparse non-linear dimensionality reduction, it was applied to an autoencoder, with the following adaptations:

- Path penalties were applied between the input ($\mathbf{x} \in \mathbb{R}^{d_x}$) and latent ($\mathbf{z} \in \mathbb{R}^{d_z}$) variables, and between the latent and output ($\hat{\mathbf{x}} \in \mathbb{R}^{d_x}$) variables.
- To enforce the encoder and the decoder to be symmetric, the group lasso groups were defined as all paths connecting x_i and z_j , together with all paths connecting z_j and \hat{x}_i , i.e.

$$\mathbf{W}_{\text{PL}} := \sqrt{\prod_{l=L,\dots,1} (\mathbf{W}_l^E)^2 + \left(\prod_{l=L,\dots,1} (\mathbf{W}_l^D)^2 \right)^\top},$$

where $\{\mathbf{W}_l^E\}_{l=1}^L$ and $\{\mathbf{W}_l^D\}_{l=1}^L$ are the weight matrices of the encoder and the decoder, respectively. This means that if an input is disconnected to a latent variable, so is the corresponding output, by construction.

- To encourage the algorithm to make equal use of the latent dimensions, we added an exclusive lasso penalty to the elements in \mathbf{W}_{PL} , with as many groups as there are latent dimensions, each group being defined as the connections to a given latent dimension.

3. Experiments

In order to evaluate path lasso for dimensionality reduction we applied it to three different data sets, one with synthetic data, consisting of Gaussian clusters on a hypercube, one with text documents from newsgroup posts, and one with images of faces. In each experiment, 20 % of the data was set aside for testing and the remaining 80 % was split 90-10 into training and validation data; all visualizations were made using the testing data. All autoencoders used one hidden layer with tanh activations in the encoder and decoder respectively, and were trained with l_2 -loss. For optimization stages not using proximal gradient descent, the Adam optimizer (Kingma and Ba, 2014) was used.

In addition to path lasso, we used a standard autoencoder, a sparse autoencoder, an autoencoder with parameter-wise l_1 -regularization and thresholding (hereafter referred to as standard lasso), PCA and sparse PCA. The reason for adding thresholding in standard lasso is because unless proximal methods are used, no parameters are set exactly to zero, but instead to very small values, as discussed in Section 2.2. It should also be noted that the sparse autoencoder is sparse in terms of observations to latent dimensions, and not in terms of original to latent dimensions, as we are interested in; see the text associated to Equation (1) for details. The standard autoencoder and PCA are of course not sparse in any sense. Since "observation sparse" algorithms are not as relevant to us as the "truly sparse" algorithms, we omit them in some of the comparisons.

The following measures were calculated on the testing data:

- Reconstruction error as explained variance (R^2).

- Fraction of correctly identified reconstructions. A reconstructed observation is considered correctly identified if it is closer to its own original observation than any of the other ones, measured in l_2 -distance, i.e. $\|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2 < \|\hat{\mathbf{x}}_i - \mathbf{x}_j\|_2, i \neq j$. This is hereafter referred to as observation reconstruction match.
- Fraction of correctly reconstructed labels. The reconstructed label is defined as the label of the original observation that is closest to the reconstructed observation, where distance is measured as above, i.e. the label of \mathbf{x}_j , where $\|\hat{\mathbf{x}}_i - \mathbf{x}_j\|_2 < \|\hat{\mathbf{x}}_i - \mathbf{x}_k\|_2, j \neq k$. This is hereafter referred to as label reconstruction match.

3.1 Synthetic Data Set

Sixteen clusters were generated in \mathbb{R}^4 , centered at each of the sixteen vertices in the hypercube $\{0, 1\}^4$. For cluster i , 100 data points were sampled according to $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, 0.01 \cdot \mathbf{I}_4)$, where \mathbf{I}_4 is the identity matrix and $\boldsymbol{\mu}_i$ is one of the sixteen vertices in $\{0, 1\}^4$. The four dimensional data set was reduced down to two dimensions using the six different algorithms. For the four autoencoder based algorithms, the number of nodes in the five layers of the autoencoder were 4, 50, 2, 50 and 4, respectively. The three sparse algorithms (in the sense original to latent dimensions) were penalized so that four of the original eight connections remained. The experiment was performed twice, with and without added noise, distributed according to $\mathcal{N}(0, 0.3^2)$.

Ten different splits of the data into training and validation sets were done. For each split, three different optimization seeds were used and the seed resulting in the best R^2 value on the validation data was chosen. The resulting mean and standard deviations are presented in Table 1. The p-values come from the one-sided paired rank test, testing whether path lasso performs better than the competing algorithm. P-values smaller than 1 % are marked in bold. With noise added, path lasso performs significantly better than the other algorithms both in terms of R^2 and reconstruction match. Without noise path lasso still performs better than the dense algorithms in terms of observation reconstruction match, while all the four non-linear methods perform very well in terms of R^2 and label reconstruction.

The results with the best R^2 values are plotted in Figure 2, where clusters that are diagonal to each other in \mathbb{R}^4 (e.g. $(0, 1, 0, 0)$ and $(1, 0, 1, 1)$) are plotted using the same color, but with different markers - circles or crosses. For the three sparse algorithms, path lasso, standard lasso and sparse PCA, each of the two latent dimensions becomes a combination of two of the original four dimensions, which can be seen in the axis aligned data in the plots. Even with no added noise the linear algorithms, PCA and sparse PCA, are not able to fully separate the sixteen clusters, while all four non-linear algorithms, based on autoencoders, are.

3.2 Text - 20 Newsgroup Data Set

To test the algorithm on text data, the 20 newsgroups data set¹ was used. Out of the original 20 categories, the following 4 were selected: `soc.religion.christian`, `sci.space`, `comp.windows.x` and `rec.sport.hockey`, which resulted in 31225 documents. Then, for

1. Available at <http://qwone.com/~jason/20Newsgroups/>.

Algorithm	Noise	Connections	R^2		Reconstruction Match			
			Mean (std)	p-value	Observation		Label	
					Mean (std)	p-value	Mean (std)	p-value
Path Lasso	No	4	0.98 (0.0013)	-	0.37 (0.015)	-	1.0 (0.0014)	-
Standard Lasso	No	4	0.98 (0.0020)	0.28	0.36 (0.022)	0.070	1.0 (0.0014)	0.68
Sparse PCA	No	4	0.48 (0.0018)	0.00098	0.067 (0.019)	0.00098	0.30 (0.014)	0.0029
Autoencoder	No	8	0.98 (0.0026)	0.019	0.35 (0.013)	0.0046	1.0 (0.0021)	0.16
Sparse AE	No	8	0.98 (0.0027)	0.50	0.34 (0.0076)	0.0039	1.0 (0.0013)	0.018
PCA	No	8	0.48 (0.0014)	0.00098	0.067 (0.015)	0.00098	0.43 (0.044)	0.0029
Path Lasso	Yes	4	0.89 (0.017)	-	0.38 (0.038)	-	0.88 (0.019)	-
Standard Lasso	Yes	4	0.58 (0.097)	0.00098	0.11 (0.030)	0.00098	0.39 (0.11)	0.0029
Sparse PCA	Yes	4	0.50 (0.0016)	0.00098	0.11 (0.0061)	0.0029	0.37 (0.010)	0.0029
Autoencoder	Yes	8	0.86 (0.011)	0.0029	0.33 (0.014)	0.0095	0.86 (0.016)	0.0088
Sparse AE	Yes	8	0.85 (0.0038)	0.0020	0.30 (0.015)	0.00098	0.85 (0.0035)	0.0020
PCA	Yes	8	0.50 (0.0043)	0.00098	0.12 (0.0088)	0.00098	0.40 (0.030)	0.00098

Table 1: Number of remaining connections, explained variance, reconstruction match and p-value for six the algorithms when reducing 16 clusters from four to two dimensions with and without added noise. P-values are for the one-sided paired rank test, that tests whether path lasso performs better than the competing algorithm, with p-values smaller than 1 % marked in bold. For the three sparse algorithms, the number of connections is always four, by construction; for three the dense algorithms it is always eight. Especially for noisy data, path lasso outperforms the competing algorithms.

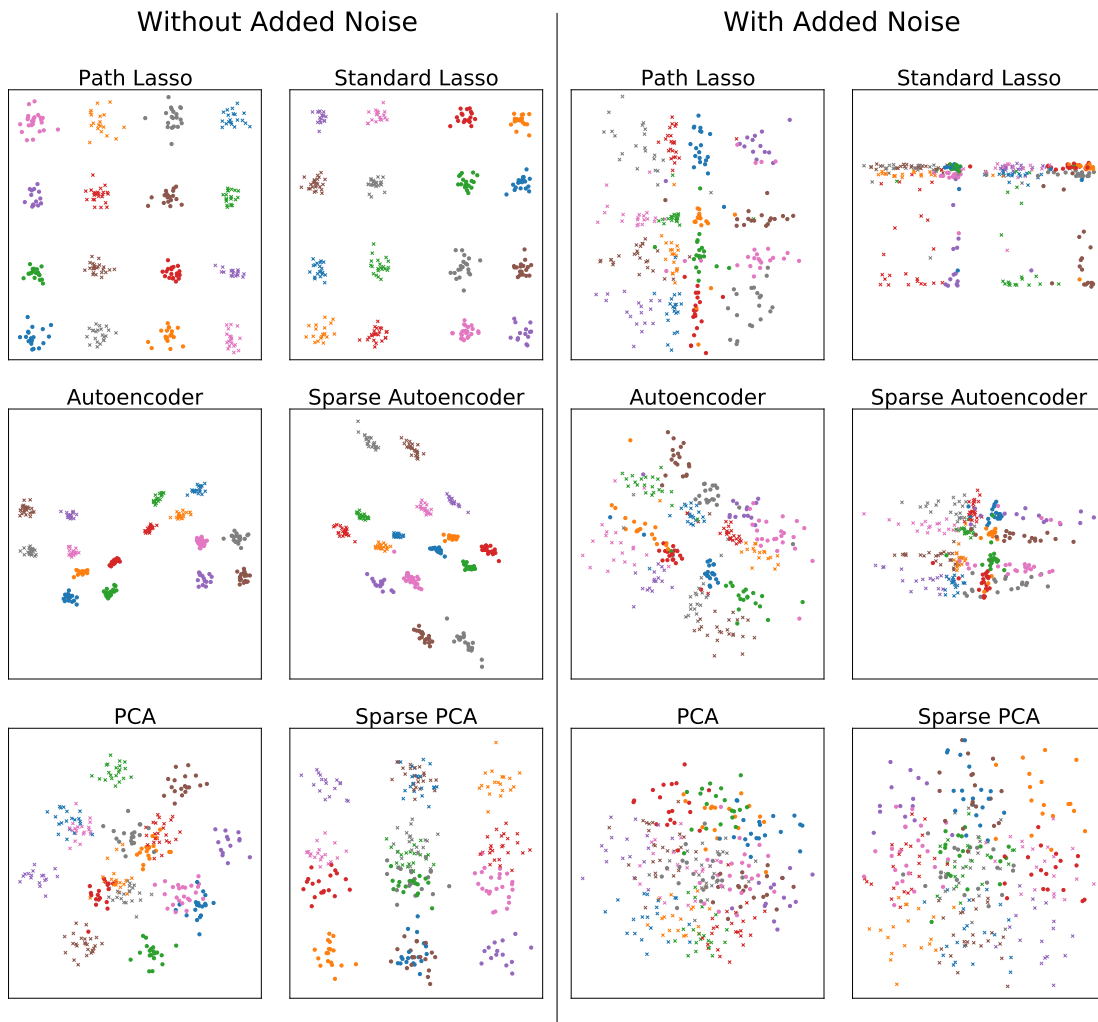


Figure 2: Reduction of 16 clusters from four to two dimension with and without added noise, using six different algorithms. Clusters that are diagonal to each other in \mathbb{R}^4 have the same color, but different markers.

each word the tf-idf score was calculated, and the 100 words with the highest score were kept, resulting in a 31225×100 data matrix. Using a standard autoencoder, path lasso, standard lasso and sparse PCA, the 100 dimensional data set was mapped down to two, four and 25 dimensions, where the 25D case was added for comparing the test measures at moderate dimensionality reductions. For the autoencoder based algorithms, the layer widths were 100, 50, 2 (4, 25), 50 and 100 nodes.

Two different sparsity levels were used; One sparse, to be able to compare interpretability, and one almost dense, to be able to compare to a standard autoencoder. In each case the penalties were set to obtain (approximately) the same sparsity, measured as number of connections, for all algorithms.

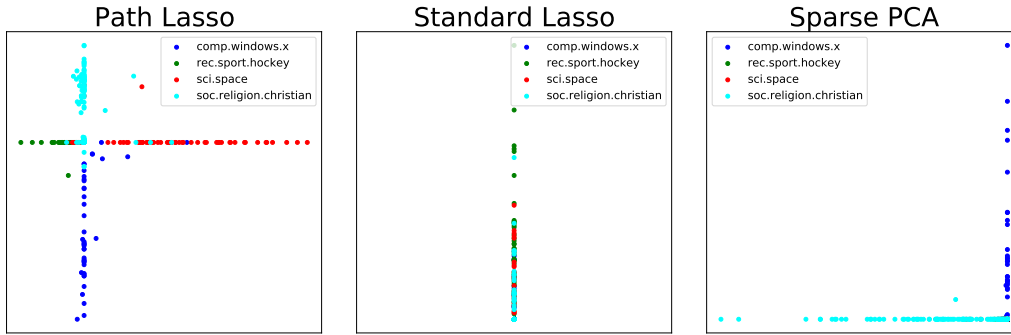
The latent spaces in the sparse case are shown in Figure 3. In the 4D case, standard lasso only uses one latent dimension and is incapable of distinguishing between the categories, sparse PCA maps one category to each latent dimension, while path lasso has a tendency to map two categories to each latent dimension, one to the positive and one to the negative axis. This is further accentuated when compressing to only two dimensions. While path lasso is able to identify all four categories, sparse PCA only identifies two of them. Standard lasso still only uses one dimension and is not able to identify any categories at all. The use of only one latent dimension by standard lasso is likely attributed to the flexibility of the autoencoder, and without the structure in the penalty imposed by the path lasso algorithm, there is less incentive to use all parts of the latent space. The same tendency is visible in Figure 4c.

To see which original dimensions (words) contribute to which latent dimensions, the non-zero elements of \mathbf{W}_{PL} can be used, but since all elements in this matrix are non-negative, it gives no information about the sign. Instead the corresponding signed matrix was created according to $\mathbf{W}_2 \cdot \mathbf{W}_1 + (\mathbf{W}_4 \cdot \mathbf{W}_3)^T$, where $(\mathbf{W}_1, \mathbf{W}_2)$ and $(\mathbf{W}_3, \mathbf{W}_4)$ are the weight matrices of the encoder and the decoder, respectively. The signed words in the latent dimensions are presented in Tables 2 and 3. The assignment of words to the 4 (8) latent half-axes done by path lasso and sparse PCA is consistent with the results in Figure 3. Path lasso also seems to identify a subcategory of `comp.windows.x` related to e-mails.

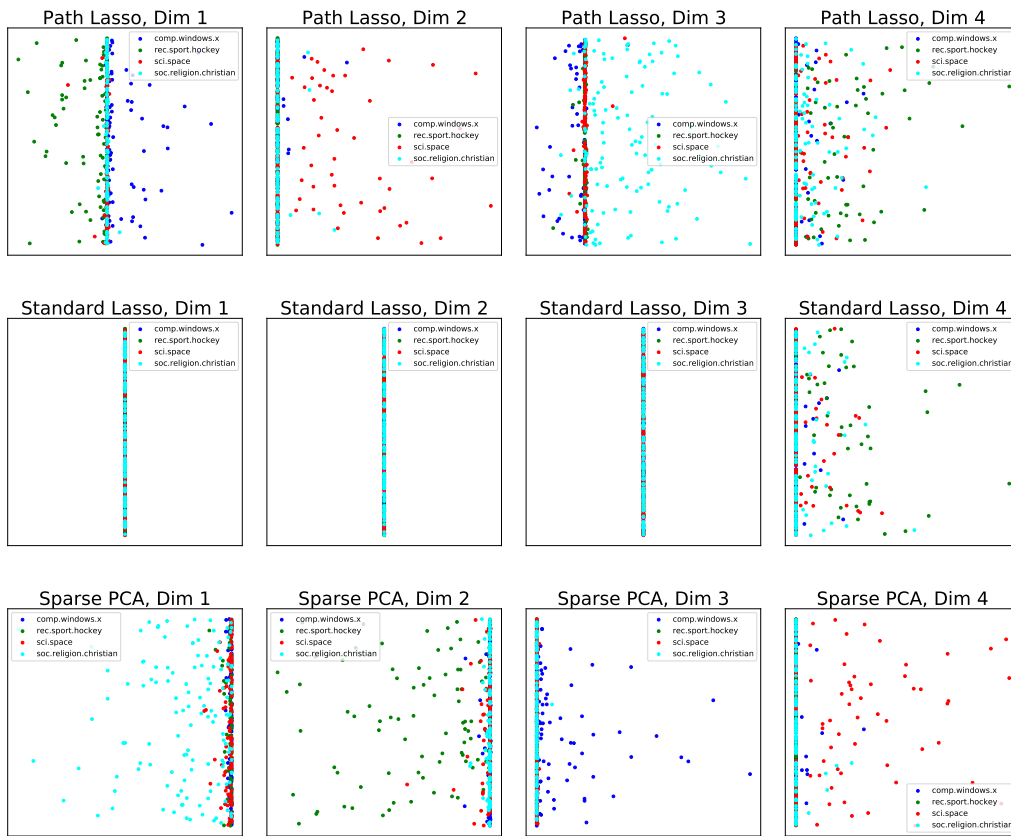
In Tables 4 and 5, remaining connections, explained variance and reconstruction match are presented for both sparsity levels. We conclude that path lasso performs best among the sparse algorithms both in terms of explained variance and reconstruction match. Compared to the standard autoencoder, path lasso performs slightly better in terms of reconstruction error and observation reconstruction match, which is in line with the results in Section 3.1. It is also noticeable that with increasing latent dimensionality the advantage of non-linear over linear models decreases.

3.3 Images - AT&T Face Database

We also tested the algorithm on the AT&T face database (Samaria and Harter, 1994), which contains 400 grayscale images of faces. The images were compressed to a size of 60×50 pixels, after which the 3000 dimensional images were reduced to 5 dimensions. Both the encoder and the decoder had 1000 units wide hidden layers, and to assure that a pixel that was disconnected from all the latent dimensions got a value of zero, no bias parameters were used.



(a) Latent 2D space for the newsgroup data



(b) Latent 4D space for the newsgroup data

Figure 3: Latent spaces of the news data, when compressed down to two and to four latent dimensions. In 3b the y-axis value is just a random number, added to increase readability. Path lasso uses the latent space more efficiently than the other algorithms, being capable of mapping one category to each half-axis.

Algorithm	Dimension	Axis	Words
Path Lasso	1 of 2	Negative	team, hockey, better, probably
		Positive	space
	2 of 2	Negative	window, hi
		Positive	god, church
Standard Lasso	1 of 2	Negative	-
		Positive	-
	2 of 2	Negative	-
		Positive	people, good, team, new, did, hockey, make, going, better, probably, lot
Sparse PCA	1 of 2	Negative	god, people, jesus, believe, bible, christ, faith, life
		Positive	-
	2 of 2	Negative	-
		Positive	window, application

Table 2: Positive and negative words in the two latent dimension, with sign calculated as $\text{sign}(W_2 \cdot W_1 + (W_4 \cdot W_3)^\top)$.

Again, each algorithm was penalized to obtain (approximately) the same number of connections, at two different sparsity levels. One low, with almost all of the $5 \cdot 3000 = 15000$ connections kept to be able to compare to a dense, l_2 -regularized autoencoder; and one high, with only a fourth of the connections kept, to compare the algorithms at more extreme sparsity levels.

The results are summarized in Table 6. Since the images are unlabeled the label reconstruction match was replaced by a 10 nearest neighbor reconstruction match, where an observation is considered correctly reconstructed if its own original observation is among the ten closest original observations.

Path lasso and standard lasso do much better than sparse PCA in terms of reconstruction, and at low sparsity they are on par with the dense autoencoder. For observation reconstruction match, path lasso outperforms the other two algorithms, and at low sparsity also the standard autoencoder. For the 10 nearest neighbors reconstruction match, all non-linear algorithms do very well at low sparsity, while at high sparsity path lasso does better than the other algorithms.

The five eigenfaces, calculated as $\text{Dec}(e_i)$ where Dec is the decoder and e_i is the i -th standard basis vector in \mathbb{R}^5 , are shown in the left column of Figure 4. For PCA, the decoder function corresponds to multiplication with the loadings matrix, whose i -th column is the i -th eigenface. A pixel with value zero is colored red. Sparse PCA distributes its non-zero pixels more evenly among the latent dimensions, while again standard lasso does not use all 5 latent dimensions. The right column shows the reconstruction of some of the images in the test set, using the four different algorithms. The reconstructed images using path lasso look more diversified than for the other sparse algorithms, something that is in line with the superior reconstruction match of path lasso.

Algorithm	Dimension	Axis	Words	
Path Lasso	1 of 4	Negative	game, year, hockey, play, games, players, season, nhl	
		Positive	window, use, need, server, program, motif, using, windows, application, widget	
	2 of 4	Negative	-	
		Positive	space	
	3 of 4	Negative	know, thanks, edu, mail, hi, list	
		Positive	god, don, think, people, jesus, say, believe, church, christians, christian, bible, christ, faith, life	
	4 of 4	Negative	-	
		Positive	team, new, hockey, going, better, probably, lot	
	Standard Lasso	1 of 4	Negative	-
			Positive	-
2 of 4		Negative	-	
		Positive	-	
3 of 4		Negative	-	
		Positive	-	
4 of 4		Negative	god, like, know, don, think, space, does, christ, time, window, thanks, use, jesus, way, say, sun, believe, need, want, problem, edu, server, hi, right, church, program, using, work, nasa, life, christians, true, bible, help, mail, used, actually	
		Positive	just, people, good, team, new, did, hockey, make, going, better, probably, lot	
Sparse PCA	1 of 4	Negative	god, don, think, people, does, jesus, say, believe, church, things, question, said, christians, true, christian, bible, come, world, point, christ, faith, life	
		Positive	-	
	2 of 4	Negative	game, good, team, year, hockey, play, games, players, season, better, best, nhl	
		Positive	use	
	3 of 4	Negative	-	
		Positive	window, server, program, motif, using, windows, application, widget	
	4 of 4	Negative	-	
		Positive	space, nasa, earth	

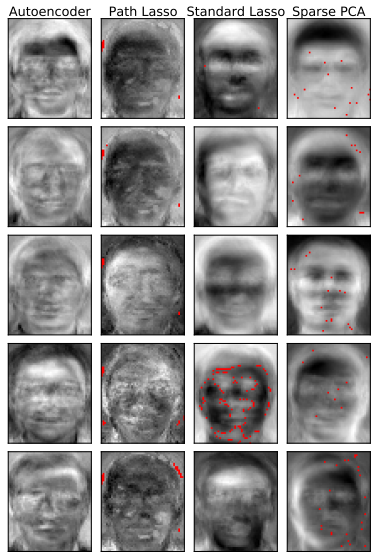
Table 3: Positive and negative words in the four latent dimension, with sign calculated as $\text{sign}(W_2 \cdot W_1 + (W_4 \cdot W_3)^\top)$.

Dimensions	Algorithm	Connections	R^2	Reconstruction Match	
				Observation	Label
2	Path lasso	9	0.13	0.021	0.47
2	Standard Lasso	11	0.015	0.0021	0.25
2	Sparse PCA	10	0.028	0.0084	0.29
4	Path lasso	46	0.20	0.027	0.55
4	Standard Lasso	47	0.015	0.0021	0.25
4	Sparse PCA	46	0.11	0.017	0.43
25	Path Lasso	102	0.55	0.34	0.66
25	Standard Lasso	104	0.069	0.0063	0.33
25	Sparse PCA	104	0.40	0.14	0.56

Table 4: Number of remaining connections, explained variance and reconstruction match for the three algorithms on the newsgroup data at high sparsity. Path lasso performs best both in terms of explained variance and reconstruction match.

Dimensions	Algorithm	Connections	R^2	Reconstruction Match	
				Observation	Label
2	Autoencoder	200	0.24	0.048	0.56
2	Path Lasso	194	0.27	0.055	0.54
2	Standard Lasso	196	0.073	0.0042	0.35
2	Sparse PCA	195	0.051	0.0063	0.32
4	Autoencoder	400	0.33	0.080	0.58
4	Path Lasso	396	0.34	0.10	0.56
4	Standard Lasso	398	0.19	0.019	0.45
4	Sparse PCA	397	0.14	0.015	0.43
25	Autoencoder	2500	0.60	0.47	0.71
25	Path Lasso	2475	0.61	0.49	0.72
25	Standard Lasso	2476	0.44	0.16	0.61
25	Sparse PCA	2486	0.44	0.15	0.57

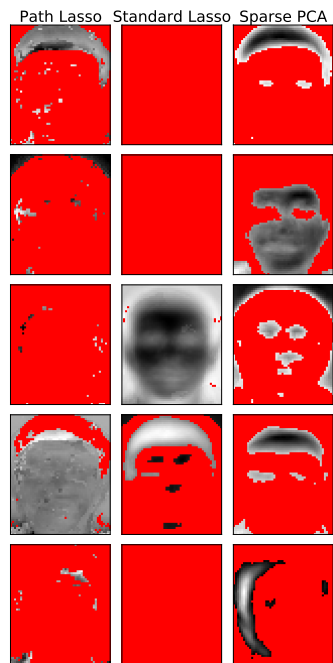
Table 5: Number of remaining connections, explained variance and reconstruction match for the four algorithms on the newsgroup data at low sparsity. Path lasso performs better than standard lasso and sparse PCA both in terms of explained variance and reconstruction match. Compared to the standard autoencoder path lasso performs slightly better in terms of explained variance and observation reconstruction match.



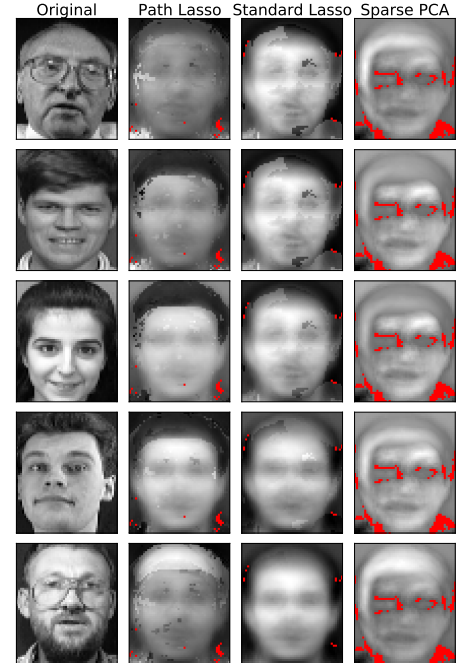
(a) Eigenfaces at low sparsity



(b) Reconstructions at low sparsity



(c) Eigenfaces at high sparsity



(d) Reconstructions at high sparsity

Figure 4: Eigenfaces and examples of reconstructions of test faces at small and large penalization. Zeros are indicated with red. The results of the dense autoencoder are presented together with the low sparsity results. The reconstructed images using path lasso look more diverse than those of the other sparse algorithms. At high sparsity standard lasso does not use all 5 latent dimensions.

Algorithm	Connections	R^2	Reconstruction Match	
			Observation	10 Nearest Neighbors
Autoencoder	15000	0.73	0.66	0.97
Path Lasso	14923	0.72	0.70	0.96
Standard Lasso	14936	0.70	0.55	0.96
Sparse PCA	14923	-0.24	0.16	0.60
Path Lasso	3770	0.57	0.21	0.75
Standard Lasso	3771	0.44	0.025	0.26
Sparse PCA	3770	-1.2	0.10	0.46

Table 6: Number of remaining connections, explained variance and reconstruction match for the four algorithms on the image data. Path lasso outperforms the competing algorithms in terms of observation reconstruction match, and at low sparsity even beats the fully connected autoencoder.

4. Conclusions

We proposed path lasso, a penalty that creates structured sparsity in feedforward neural networks, by using a group lasso penalty to remove all connections between two nodes in two non-adjacent layers. We applied path lasso to an autoencoder to obtain sparse, non-linear dimensionality reduction, and showed that this non-linearity makes path lasso much more flexible than sparse PCA, leading to a lower reconstruction error and a higher reconstruction match. Thanks to its higher flexibility path lasso is also able to use the latent space more efficiently, something that proved essential when the latent dimensions were few. Compared to an autoencoder with individually lasso penalized links, path lasso performed better in terms of reconstruction, reconstruction match and interpretation of the latent space. In addition, at low sparsity levels path lasso resulted in a better reconstruction match than the standard autoencoder.

Using path lasso for non-linear, sparse dimensionality reduction, flexibility and interpretability can be combined in a new way, enabling compression to fewer dimensions, with preserved interpretability. In this paper, we used the path lasso penalty in an autoencoder. However, path lasso can be used in many other types of feedforward neural networks with applications including non-linear, sparse multivariate regression, and non-linear, sparse network models. Such investigations are left for future work.

Code is available at https://github.com/allerbo/path_lasso.

Acknowledgments

We would like to thank the anonymous reviewers, whose suggestions helped improve and clarify the manuscript.

This research was supported by funding from the Swedish Research Council (VR), the Swedish Foundation for Strategic Research, the Wallenberg AI, Autonomous Systems and Software Program (WASP), and the Chalmers AI Research Center (CHAIR).

Appendix A. Accelerating the Non-Negative Matrix Factorization

To increase the speed of the matrix factorization step in Equation (7) the following three approaches were used: Substitution, parallelization and Boolean matrix factorization.

A.1 Substitution

By first applying an off-the-shelf optimization method to $f(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) + \lambda \mathbf{W}_{\text{PL}}$, where $g(\boldsymbol{\theta})$ is the reconstruction error, we obtain a solution where some of the paths have very small values, although non-zero. We then, as in adaptive lasso, use this solution to initialize a second optimization stage, using proximal gradient descent, with individual penalties for each connection, where the penalties depend on the magnitude of the connection after the first stage:

$$\lambda_{i_L, i_0} := \frac{\lambda}{((\hat{\mathbf{W}}_{\text{PL}})_{i_L i_0})^\gamma}$$

where $\hat{\mathbf{W}}_{\text{PL}}$ is the value of \mathbf{W}_{PL} after the optimization with the off-the-shelf optimizer and $\gamma > 0$. Just as in Section 2.5, $\gamma = 2$ was used. Choosing a relatively small value of λ , connections that are not close to zero after the first stage will be left more or less unpenalized, while connections close to zero will be penalized hard. Another advantage of this two stage procedure is that we are able to further increase the speed in the first stage by leveraging on momentum based optimization methods, which are not available for proximal gradient descent.

A.2 Parallelization

In general, the larger the weight matrices, the more time consuming is the factorization of the penalized path matrix into penalized weight matrices. By splitting the weight matrices into blocks, the factorization can instead be done in parallel for smaller sub-matrices. In the simplest example, three layers are split into two parts each. This corresponds to two weight matrices, whose rows and columns are both split into two blocks each, and is shown in Equation (8) below. Let $\mathbf{A} \in (\mathbb{R}_{\geq 0})^{d_3 \times d_2}$ and $\mathbf{B} \in (\mathbb{R}_{\geq 0})^{d_2 \times d_1}$ denote the absolute valued weight matrices and let $\mathbf{C} + \mathbf{D} \in (\mathbb{R}_{\geq 0})^{d_3 \times d_1}$ denote the penalized path matrix, where the sum of the paths in $(\mathbf{C} + \mathbf{D})_{ij}$ is split into \mathbf{C}_{ij} , containing the sum of some of the paths, and \mathbf{D}_{ij} , containing the sum of the remaining paths. Which paths belong to \mathbf{C} and which belong to \mathbf{D} depends on how \mathbf{A} and \mathbf{B} are split. Then, the factorization of $\mathbf{C} + \mathbf{D}$ into $\mathbf{A} \cdot \mathbf{B}$ can be divided into $2^3 = 8$ smaller factorizations which can be solved in parallel, i.e.

$$\left[\begin{array}{c|c} \mathbf{C}_{11} + \mathbf{D}_{11} & \mathbf{C}_{12} + \mathbf{D}_{12} \\ \hline \mathbf{C}_{21} + \mathbf{D}_{21} & \mathbf{C}_{22} + \mathbf{D}_{22} \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right] \cdot \left[\begin{array}{c|c} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \hline \mathbf{B}_{21} & \mathbf{B}_{22} \end{array} \right] \quad (8)$$

can be split into

$$\begin{aligned} \left[\begin{array}{c|c} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \hline \mathbf{C}_{21} & \mathbf{C}_{22} \end{array} \right] &= \left[\begin{array}{c|c} \mathbf{A}_{11} & \mathbf{0} \\ \hline \mathbf{A}_{21} & \mathbf{0} \end{array} \right] \cdot \left[\begin{array}{c|c} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_{11} \cdot \mathbf{B}_{11} & \mathbf{A}_{11} \cdot \mathbf{B}_{12} \\ \hline \mathbf{A}_{21} \cdot \mathbf{B}_{11} & \mathbf{A}_{21} \cdot \mathbf{B}_{12} \end{array} \right] \\ \left[\begin{array}{c|c} \mathbf{D}_{11} & \mathbf{D}_{12} \\ \hline \mathbf{D}_{21} & \mathbf{D}_{22} \end{array} \right] &= \left[\begin{array}{c|c} \mathbf{0} & \mathbf{A}_{12} \\ \hline \mathbf{0} & \mathbf{A}_{22} \end{array} \right] \cdot \left[\begin{array}{c|c} \mathbf{0} & \mathbf{0} \\ \hline \mathbf{B}_{21} & \mathbf{B}_{22} \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_{12} \cdot \mathbf{B}_{21} & \mathbf{A}_{12} \cdot \mathbf{B}_{22} \\ \hline \mathbf{A}_{22} \cdot \mathbf{B}_{21} & \mathbf{A}_{22} \cdot \mathbf{B}_{22} \end{array} \right] \end{aligned}$$

or equivalently

$$\begin{aligned} \mathbf{C}_{11} &= \mathbf{A}_{11} \cdot \mathbf{B}_{11}, \quad \mathbf{C}_{12} = \mathbf{A}_{11} \cdot \mathbf{B}_{12}, \quad \mathbf{C}_{21} = \mathbf{A}_{21} \cdot \mathbf{B}_{11}, \quad \mathbf{C}_{22} = \mathbf{A}_{21} \cdot \mathbf{B}_{12} \\ \mathbf{D}_{11} &= \mathbf{A}_{12} \cdot \mathbf{B}_{21}, \quad \mathbf{D}_{12} = \mathbf{A}_{12} \cdot \mathbf{B}_{22}, \quad \mathbf{D}_{21} = \mathbf{A}_{22} \cdot \mathbf{B}_{21}, \quad \mathbf{D}_{22} = \mathbf{A}_{22} \cdot \mathbf{B}_{22}. \end{aligned}$$

Here, $\mathbf{A}_{ij} \in (\mathbb{R}_{\geq 0})^{d_{3i} \times d_{2j}}$, $\mathbf{B}_{ij} \in (\mathbb{R}_{\geq 0})^{d_{2i} \times d_{1j}}$ and $\mathbf{C}_{ij}, \mathbf{D}_{ij} \in (\mathbb{R}_{\geq 0})^{d_{3i} \times d_{1j}}$, where $i, j \in \{1, 2\}$ and $d_{ki} + d_{kj} = d_k$ for $k = 1, 2, 3$, are all block matrices. This results in two, potentially different, solutions for each \mathbf{A}_{ij} and \mathbf{B}_{ij} , that need to be aggregated. To be conservative, and set a link to zero only when it is zero in both solutions, we use the maximum value element-wise for aggregation:

$$(\mathbf{A}_{ij})_{kl} = \max_{n \in \{1, 2\}} (\mathbf{A}_{ij}^n)_{kl}$$

where \mathbf{A}_{ij}^n is the n -th solution of \mathbf{A}_{ij} . Remember that all elements in all matrices are non-negative.

This generalizes trivially to more layers and splits and the number of equations is always the product of the splits over all the layers.

A.3 Boolean Matrix Factorization

Another way to speed up the computations is early stopping of the matrix factorization. This may however lead to some of the elements on the right hand side in Equation (7) being very close, but not equal, to zero, introducing the need to threshold. We determine the optimal threshold value using Boolean matrix multiplication, which differs from ordinary matrix multiplication in the following way:

- All elements are in $\{0, 1\}$.
- Instead of ordinary sum, Boolean sum (or Boolean **or**), \vee , is used: $0 \vee 0 = 0$, $0 \vee 1 = 1$, $1 \vee 0 = 1$, $1 \vee 1 = 1$.

We define a Boolean matrix for each matrix in Equation (7) as:

$$\begin{aligned} (\mathbf{B})_{ij} &:= \mathbb{I} \left[\left(\prod_{l=L, \dots, 1} |\mathbf{W}_l| \odot \left(1 - \frac{\alpha \lambda}{\mathbf{W}_{\text{PL}}} \right)^+ \right)_{ij} > 0 \right] \\ (\mathbf{B}_l)_{ij}(\tau) &:= \mathbb{I} \left[(\tilde{\mathbf{W}}_l)_{ij} > \tau \right], \quad \tau \geq 0 \end{aligned}$$

where $\mathbb{I}[\cdot]$ is the (element-wise) indicator function which equals 1 when its argument is true and 0 otherwise. Thus, in the matrix on the left hand side in Equation (7), zero elements become zeros and all other elements become ones, while for the matrices on the right hand side, we use a threshold to decide which elements should be set to zero. In the first case, an element of one means that there is a non-zero connection, while in the second case it means that there is a non-zero link. If $\mathbf{B} = \bigvee_{l=L, \dots, 1} \mathbf{B}_l$, where \bigvee denotes Boolean matrix multiplication, then the matrix factorization was successful in the sense that all connections in \mathbf{W}_{PL} are represented by appropriate links in the weight matrices. Each differing element means that either two nodes are disconnected in \mathbf{W}_{PL} but not in $\{\tilde{\mathbf{W}}_l\}_{l=1}^L$ or vice versa. To

minimize this quantity, we minimize $\sum |\mathbf{B} - \bigvee_{l=L,\dots,1} \mathbf{B}_l(\tau)|$, with the absolute value taken element-wise, with respect to τ , and then threshold $\{\tilde{\mathbf{W}}_l\}_{l=1}^L$ using the resulting τ :

$$\tau = \operatorname{argmin}_{\tau \geq 0} \sum_{\text{all elements}} \left| \mathbb{I} \left[\prod_{l=L,\dots,1} |\mathbf{w}_l| \odot \left(1 - \frac{\alpha\lambda}{\mathbf{w}_{\text{PL}}} \right)^+ > 0 \right] - \bigvee_{l=L,\dots,1} \mathbb{I}[\tilde{\mathbf{w}}_l > \tau] \right|$$

$$\tilde{\mathbf{W}}_l \leftarrow \tilde{\mathbf{W}}_l \odot \mathbb{I}[\tilde{\mathbf{W}}_l > \tau]$$

Since the objective function is piecewise constant and thus hard to optimize, we simply evaluate it for 20 logarithmically spaced values of $\tau \in [10^{-10}, 1]$ and pick the best one.

Appendix B. Proofs

Proof of Proposition 3 Using the definition of matrix multiplication, with k running over all paths connecting x_{i_0} and y_{i_L} , we obtain

$$\begin{aligned} \sqrt{\sum_{k \in g_{i_L i_0}} (p_{i_L i_0}^k)^2} &= \sqrt{\sum_{i_{L-1}=1}^{d_{L-1}} \sum_{i_{L-2}=1}^{d_{L-2}} \cdots \sum_{i_1=1}^{d_1} (|w_{i_L i_{L-1}}^L| \cdot |w_{i_{L-1} i_{L-2}}^{L-1}| \cdots |w_{i_1 i_0}^1|)^2} \\ &= \sqrt{\sum_{i_{L-1}=1}^{d_{L-1}} \sum_{i_{L-2}=1}^{d_{L-2}} \cdots \sum_{i_1=1}^{d_1} (w_{i_L i_{L-1}}^L)^2 \cdot (w_{i_{L-1} i_{L-2}}^{L-1})^2 \cdots (w_{i_1 i_0}^1)^2} \\ &= \left(\sqrt{(\mathbf{W}_L)^2 \cdot (\mathbf{W}_{L-1})^2 \cdots (\mathbf{W}_1)^2} \right)_{i_L i_0} = (\mathbf{W}_{\text{PL}})_{i_L i_0}. \end{aligned}$$

■

To prove Proposition 4 we begin with the following lemma:

Lemma 6 *Let $\{\mathbf{A}_k\}_{k=1}^n$ be a set of arbitrary matrices and $\{\mathbf{D}_k\}_{k=1}^n$ a set of diagonal matrices, where all elements in \mathbf{A}_k and \mathbf{D}_k are bounded. Let $f \geq 0$ an element-wise function where $f(x) = 0$ if and only if $x = 0$. Let the dimensions of the matrices be such that the matrix multiplications below make sense. Then*

$$(f(\mathbf{A}_1) \cdot f(\mathbf{A}_2) \cdots f(\mathbf{A}_n))_{ij} = 0 \implies (\mathbf{D}_1 \cdot \mathbf{A}_1 \cdot \mathbf{D}_2 \cdot \mathbf{A}_2 \cdot \mathbf{D}_2 \cdots \mathbf{D}_n \cdot \mathbf{A}_n)_{ij}.$$

Proof Denote $a_{ij}^k := (\mathbf{A}_k)_{ij}$ and $d_{ij}^k := (\mathbf{D}_k)_{ij}$. Then for the left hand side

$$(f(\mathbf{A}_1) \cdot f(\mathbf{A}_2) \cdots f(\mathbf{A}_n))_{ij} = \sum_{k_1} \sum_{k_2} \cdots \sum_{k_{n-1}} f(a_{i k_1}^1) \cdot f(a_{k_1 k_2}^2) \cdots f(a_{k_{n-1} j}^n)$$

and for the right hand side

$$(\mathbf{D}_1 \cdot \mathbf{A}_1 \cdot \mathbf{D}_2 \cdot \mathbf{A}_2 \cdots \mathbf{D}_n \cdot \mathbf{A}_n)_{ij} = \sum_{k_1} \sum_{k_2} \cdots \sum_{k_{n-1}} a_{i k_1}^1 \cdot a_{k_1 k_2}^2 \cdots a_{k_{n-1} j}^n \cdot d_{ii}^1 \cdot d_{k_1 k_1}^2 \cdots d_{k_{n-1} k_{n-1}}^n.$$

By the definition of f , $f(a_{ik_1}^1) \cdot f(a_{k_1k_2}^2) \cdot \dots \cdot f(a_{k_{n-1}j}^n) \geq 0$ with equality only if at least one of $a_{ik_1}^1, a_{k_1k_2}^2, \dots, a_{k_{n-1}j}^n$ is zero, which implies that $a_{ik_1}^1 \cdot a_{k_1k_2}^2 \cdot \dots \cdot a_{k_{n-1}j}^n = 0$.

The sum

$$\sum_{k_1} \sum_{k_2} \dots \sum_{k_{n-1}} f(a_{ik_1}^1) \cdot f(a_{k_1k_2}^2) \cdot \dots \cdot f(a_{k_{n-1}j}^n)$$

is zero if and only if all its terms are zero, which means that

$$a_{ik_1}^1 \cdot a_{k_1k_2}^2 \cdot \dots \cdot a_{k_{n-1}j}^n = 0$$

for all combinations of indices summed over. Multiplying with some bounded constant does not change that fact, so

$$f(a_{ik_1}^1) \cdot f(a_{k_1k_2}^2) \cdot \dots \cdot f(a_{k_{n-1}j}^n) = 0 \implies a_{ik_1}^1 \cdot a_{k_1k_2}^2 \cdot \dots \cdot a_{k_{n-1}j}^n \cdot d_{ii}^1 \cdot d_{k_1k_1}^2 \cdot \dots \cdot d_{k_{n-1}k_{n-1}}^n = 0.$$

Finally summing only zeros we get

$$\begin{aligned} 0 &= \sum_{k_1} \sum_{k_2} \dots \sum_{k_{n-1}} a_{ik_1}^1 \cdot a_{k_1k_2}^2 \cdot \dots \cdot a_{k_{n-1}j}^n \cdot d_{ii}^1 \cdot d_{k_1k_1}^2 \cdot \dots \cdot d_{k_{n-1}k_{n-1}}^n \\ &= (\mathbf{D}_1 \cdot \mathbf{A}_1 \cdot \mathbf{D}_2 \cdot \mathbf{A}_2 \cdot \dots \cdot \mathbf{D}_n \cdot \mathbf{A}_n)_{ij} \end{aligned}$$

■

Proof of Proposition 4 With $\{\mathbf{o}_l\}_{l=0}^L$ (where $\mathbf{o}_0 = \mathbf{x}$ and $\mathbf{o}_L = \mathbf{y}$) denoting the outputs and $\{\mathbf{i}_l\}_{l=1}^L$ the inputs of the activation functions $\{\Phi_l\}_{l=1}^L$, we can express Equation (4) recursively for $1 \leq l \leq L$ as

$$\begin{aligned} \mathbf{i}_l &= \mathbf{W}_l \mathbf{o}_{l-1} + \mathbf{b}_l \\ \mathbf{o}_l &= \Phi_k(\mathbf{i}_l) \end{aligned}$$

Applying the chain rule we obtain

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{o}_L}{\partial \mathbf{o}_0} = \frac{\partial \mathbf{o}_L}{\partial \mathbf{i}_L} \cdot \frac{\partial \mathbf{i}_L}{\partial \mathbf{o}_{L-1}} \cdot \frac{\partial \mathbf{o}_{L-1}}{\partial \mathbf{i}_{L-1}} \cdot \frac{\partial \mathbf{i}_{L-1}}{\partial \mathbf{o}_{L-2}} \cdot \dots \cdot \frac{\partial \mathbf{o}_1}{\partial \mathbf{i}_1} \cdot \frac{\partial \mathbf{i}_1}{\partial \mathbf{o}_0}, \quad (9)$$

where $\frac{\partial \mathbf{o}_l}{\partial \mathbf{i}_l}$ is a diagonal matrix, since Φ_l is an element-wise operator, and $\frac{\partial \mathbf{i}_l}{\partial \mathbf{o}_{l-1}} = \mathbf{W}_l$. We can now apply Lemma 6 with $f(x) = x^2$ to Equation (9) to obtain

$$\begin{aligned} (\mathbf{W}_{\text{PL}})_{i_L i_0} = 0 &\iff ((\mathbf{W}_L)^2 \cdot (\mathbf{W}_{L-1})^2 \cdot \dots \cdot (\mathbf{W}_1)^2)_{i_L i_0} = 0 \\ \implies 0 &= \left(\frac{\partial \mathbf{o}_L}{\partial \mathbf{i}_L} \cdot \mathbf{W}_L \cdot \frac{\partial \mathbf{o}_{L-1}}{\partial \mathbf{i}_{L-1}} \cdot \mathbf{W}_{L-1} \cdot \dots \cdot \frac{\partial \mathbf{o}_1}{\partial \mathbf{i}_1} \cdot \mathbf{W}_1 \right)_{i_L i_0} = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{i_L i_0} \end{aligned}$$

where the equivalence comes from the definition of \mathbf{W}_{PL} and the implication from Lemma 6. ■

Proof of Proposition 5 Assuming the penalized path can be written as a product of penalized links, we get the following equation for the k -th path between nodes i_0 and i_L , in total $\prod_{l=0}^L d_k$ equations:

$$\begin{aligned} p_{i_L i_0}^k \cdot \left(1 - \frac{\alpha\lambda}{(\mathbf{W}_{\text{PL}})_{i_L i_0}}\right)^+ &= |w_{i_L i_{L-1}^k}^L| \cdot |w_{i_{L-1}^k i_{L-2}^k}^{L-1}| \cdot \dots \cdot |w_{i_1^k i_0}^1| \cdot \left(1 - \frac{\alpha\lambda}{(\mathbf{W}_{\text{PL}})_{i_L i_0}}\right)^+ \\ &= \text{sign}(w_{i_L i_{L-1}^k}^L) \cdot (|w_{i_L i_{L-1}^k}^L| - \alpha\lambda_{i_L i_{L-1}^k}^L)^+ \cdot \dots \cdot \text{sign}(w_{i_1^k i_0}^1) \cdot (|w_{i_1^k i_0}^1| - \alpha\lambda_{i_1^k i_0}^1)^+ \\ &=: \text{sign}(w_{i_L i_{L-1}^k}^L) \cdot \tilde{w}_{i_L i_{L-1}^k}^L \cdot \dots \cdot \text{sign}(w_{i_1^k i_0}^1) \cdot \tilde{w}_{i_1^k i_0}^1. \end{aligned}$$

The second equality is our assumption, requiring that the proximal operator from Equation (3) can be written as a product of proximal operators from Equation (2). $\lambda_{i_l^k i_{l-1}^k}^l \in [0, |w_{i_l^k i_{l-1}^k}^l|/\alpha]$ is a, possibly unique, penalty for the weight $w_{i_l^k i_{l-1}^k}^l$ and $\tilde{w}_{i_l^k i_{l-1}^k}^l := (|w_{i_l^k i_{l-1}^k}^l| - \alpha\lambda_{i_l^k i_{l-1}^k}^l)^+$. The superscript k on the indices marks that different paths go through different nodes in the inner layers.

Taking absolute values of all equations and summing over the equations where the paths belong to the same connection, i.e. they share values for i_0 and i_L and thus have the same path penalty, we obtain, for a given combination of i_0 and i_L

$$\left(\sum_{i_{L-1}=1}^{d_{L-1}} \dots \sum_{i_1=1}^{d_1} |w_{i_L i_{L-1}}^L| \cdot \dots \cdot |w_{i_1 i_0}^1| \right) \cdot \left(1 - \frac{\alpha\lambda}{(\mathbf{W}_{\text{PL}})_{i_L i_0}}\right)^+ = \sum_{i_{L-1}=1}^{d_{L-1}} \dots \sum_{i_1=1}^{d_1} \tilde{w}_{i_L i_{L-1}}^L \cdot \dots \cdot \tilde{w}_{i_1 i_0}^1$$

which, using the definition of matrix multiplication can be written as

$$\left(\prod_{l=L, \dots, 1} |\mathbf{W}_l| \odot \left(1 - \frac{\alpha\lambda}{(\mathbf{W}_{\text{PL}})_{i_L i_0}}\right)^+ \right)_{i_L i_0} = \left(\prod_{l=L, \dots, 1} \tilde{\mathbf{W}}_l \right)_{i_L i_0}.$$

Thus, solving Equation (7) results in (absolute) links values being shifted towards zero in such a way that when multiplied into paths, the paths have the correct penalization. The only thing left to do is to restore the signs of the links. \blacksquare

Appendix C. Modified Non-Negative Matrix Factorization

We solve the non-negative matrix factorization problem in Equation (7) with coordinate descent, using a modified version of the solver in python's scikit-learn module (Pedregosa et al., 2011), which in turn is based on work by Cichocki and Phan (2009) and Hsieh and Dhillon (2011). The aim is to minimize $\|\mathbf{V} - \mathbf{W} \cdot \prod_{i=1}^I \mathbf{M}_i \cdot \mathbf{H}\|_F^2$ for $I \in \mathbb{N}_0$, keeping each entry in \mathbf{W} , \mathbf{M}_i and \mathbf{H} between zero and its seed.

Coordinate descent updates each matrix separately, keeping the others fixed. Since we can write $\mathbf{W} \cdot \prod_{i=1}^I \mathbf{M}_i \cdot \mathbf{H} = \tilde{\mathbf{W}} \cdot \tilde{\mathbf{M}} \cdot \tilde{\mathbf{H}}$, where the three matrices on the right hand side are products of the matrices of the left hand side, when updating a given matrix, we can always treat the problem as a product of only three matrices. Thus we need update rules

for $\tilde{\mathbf{W}}$ and $\tilde{\mathbf{M}}$, since we can use the same algorithm for $\tilde{\mathbf{W}}$ and $\tilde{\mathbf{H}}$ by taking transpose. Hsieh and Dhillon (2011) describe the case for two matrices, i.e. $I = 0$. Our contribution is the update rule for $\tilde{\mathbf{M}}$ (and trivially adding the upper constraint on the solution).

C.1 Update Rule for $\tilde{\mathbf{W}}$

For $\tilde{\mathbf{W}} \in (\mathbb{R}^+)^{d_i \times d_r}$ we want to solve, adding the possibility to add element-wise l_1 - and l_2 -penalties:

$$\tilde{\mathbf{W}}: \min_{0 \leq (\tilde{\mathbf{W}})_{ir} \leq (\tilde{\mathbf{W}}_0)_{ir}} \frac{1}{2} \|\mathbf{V} - \tilde{\mathbf{W}} \tilde{\mathbf{H}}\|_F^2 + \sum_{i,r} \left(\lambda_1 (\tilde{\mathbf{W}})_{ir} + \frac{\lambda_2}{2} (\tilde{\mathbf{W}})_{ir}^2 \right)$$

where $\tilde{\mathbf{W}}_0$ is the seed. We update each element $(\tilde{\mathbf{W}})_{ir}$ by adding $s \mathbf{E}_{ir}$, for an optimal s , where \mathbf{E}_{ir} is a matrix with all zeros, except element (i, r) , which is 1. Defining

$$g_{ir}^{\tilde{\mathbf{W}}}(s) := \frac{1}{2} \sum_j ((\mathbf{V})_{ij} - ((\tilde{\mathbf{W}} + s \mathbf{E}_{ir}) \tilde{\mathbf{H}})_{ij})^2 + \lambda_1 ((\tilde{\mathbf{W}})_{ir} + s) + \frac{\lambda_2}{2} ((\tilde{\mathbf{W}})_{ir} + s)^2$$

we get

$$\begin{aligned} g_{ir}^{\tilde{\mathbf{W}}}(s) &= \frac{1}{2} \sum_j ((\mathbf{V})_{ij} - ((\tilde{\mathbf{W}} + s \mathbf{E}_{ir}) \tilde{\mathbf{H}})_{ij})^2 + \lambda_1 ((\tilde{\mathbf{W}})_{ir} + s) + \frac{\lambda_2}{2} ((\tilde{\mathbf{W}})_{ir} + s)^2 \\ &= \frac{1}{2} \sum_j ((\mathbf{V})_{ij} - (\tilde{\mathbf{W}} \tilde{\mathbf{H}})_{ij} - s (\tilde{\mathbf{H}})_{rj})^2 + \lambda_1 ((\tilde{\mathbf{W}})_{ir} + s) + \frac{\lambda_2}{2} ((\tilde{\mathbf{W}})_{ir} + s)^2 \\ (g_{ir}^{\tilde{\mathbf{W}}})'(s) &= \sum_j (-(\mathbf{V})_{ij} (\tilde{\mathbf{H}})_{rj} + (\tilde{\mathbf{W}} \tilde{\mathbf{H}})_{ij} (\tilde{\mathbf{H}})_{rj} + 2s (\tilde{\mathbf{H}})_{rj}^2) + \lambda_1 + \lambda_2 ((\tilde{\mathbf{W}})_{ir} + s) \\ (g_{ir}^{\tilde{\mathbf{W}}})''(s) &= \sum_j (2 (\tilde{\mathbf{H}})_{rj}^2) + \lambda_2 \\ (g_{ir}^{\tilde{\mathbf{W}}})'(0) &= \sum_{i,j} (-(\mathbf{V})_{ij} (\tilde{\mathbf{H}}^\top)_{jr} + (\tilde{\mathbf{W}} \tilde{\mathbf{H}})_{ij} (\tilde{\mathbf{H}}^\top)_{jr}) + \lambda_1 + \lambda_2 (\tilde{\mathbf{W}})_{ir} \\ &= (-\mathbf{V} \tilde{\mathbf{H}}^\top + \tilde{\mathbf{W}} \tilde{\mathbf{H}} \tilde{\mathbf{H}}^\top)_{ir} + \lambda_1 + \lambda_2 (\tilde{\mathbf{W}})_{ir} \\ (g_{ir}^{\tilde{\mathbf{W}}})''(0) &= \sum_j (\tilde{\mathbf{H}})_{rj} (\tilde{\mathbf{H}}^\top)_{jr} + \lambda_2 = (\tilde{\mathbf{H}} \tilde{\mathbf{H}}^\top)_{rr} + \lambda_2. \end{aligned}$$

Since $g_{ir}^{\tilde{\mathbf{W}}}(s)$ is quadratic in s , its Taylor expansion only contains terms up to and including s^2 :

$$g_{ir}^{\tilde{\mathbf{W}}}(s) = g_{ir}^{\tilde{\mathbf{W}}}(0) + (g_{ir}^{\tilde{\mathbf{W}}})'(0) \cdot s + \frac{1}{2} (g_{ir}^{\tilde{\mathbf{W}}})''(0) \cdot s^2$$

which is minimized at

$$s^* = -\frac{(g_{ir}^{\tilde{\mathbf{W}}})'(0)}{(g_{ir}^{\tilde{\mathbf{W}}})''(0)} = \frac{(\mathbf{V} \tilde{\mathbf{H}}^\top - \tilde{\mathbf{W}} \tilde{\mathbf{H}} \tilde{\mathbf{H}}^\top)_{ir} - \lambda_1 - \lambda_2 (\tilde{\mathbf{W}})_{ir}}{(\tilde{\mathbf{H}} \tilde{\mathbf{H}}^\top)_{rr} + \lambda_2}$$

and, to make sure $(\tilde{\mathbf{W}})_{ir} \in [0, (\tilde{\mathbf{W}}_0)_{ir}]$:

$$(\tilde{\mathbf{W}})_{ir}^{t+1} = \max((\tilde{\mathbf{W}}_0)_{ir}, \min(0, (\tilde{\mathbf{W}})_{ir}^t + s^*)).$$

C.2 Update Rule for \tilde{M}

For $M \in (\mathbb{R}^+)^{d_p \times d_r}$ the corresponding equations become

$$\tilde{M}: \min_{0 \leq (\tilde{M})_{pr} \leq (\tilde{M}_0)_{pr}} \frac{1}{2} \|V - \tilde{W} \tilde{M} \tilde{H}\|_F^2 + \sum_{p,r} \left(\lambda_1 (\tilde{M})_{pr} + \frac{\lambda_2}{2} (\tilde{M})_{pr}^2 \right)$$

$$g_{pr}^{\tilde{M}}(s) := \frac{1}{2} \sum_{i,j} ((V)_{ij} - (\tilde{W}(\tilde{M} + sE_{pr})\tilde{H})_{ij})^2 + \lambda_1 ((\tilde{M})_{pr} + s) + \frac{\lambda_2}{2} ((\tilde{M})_{pr} + s)^2$$

$$\begin{aligned} &= \frac{1}{2} \sum_{i,j} ((V)_{ij} - (\tilde{W}\tilde{M}\tilde{H})_{ij} - s(\tilde{W})_{ip}(\tilde{H})_{rj})^2 + \lambda_1 ((\tilde{M})_{pr} + s) \\ &\quad + \frac{\lambda_2}{2} ((\tilde{M})_{pr} + s)^2 \end{aligned}$$

$$\begin{aligned} (g_{pr}^{\tilde{M}})'(s) &= \sum_{i,j} (-(V)_{ij}(\tilde{W})_{ip}(\tilde{H})_{rj} + (\tilde{W}\tilde{M}\tilde{H})_{ij}(\tilde{W})_{ip}(\tilde{H})_{rj} + 2s(\tilde{W})_{ip}^2(\tilde{H})_{rj}^2) \\ &\quad + \lambda_1 + \lambda_2((\tilde{M})_{pr} + s) \end{aligned}$$

$$(g_{pr}^{\tilde{M}})''(s) = \sum_{i,j} (2(\tilde{W})_{ip}(\tilde{H})_{rj}^2) + \lambda_2$$

$$\begin{aligned} (g_{pr}^{\tilde{M}})'(0) &= \sum_{i,j} (-(V)_{ij}(\tilde{W})_{ip}(\tilde{H})_{rj} + (\tilde{W}\tilde{M}\tilde{H})_{ij}(\tilde{W})_{ip}(\tilde{H})_{rj}) + \lambda_1 + \lambda_2(\tilde{M})_{pr} \\ &= (-\tilde{W}^\top V \tilde{H}^\top + \tilde{W}^\top \tilde{W} \tilde{M} \tilde{H} \tilde{H}^\top)_{pr} + \lambda_1 + \lambda_2(\tilde{M})_{pr} \end{aligned}$$

$$(g_{pr}^{\tilde{M}})''(0) = \sum_i (\tilde{W}^\top)_{pi} (\tilde{W})_{ip} \sum_j (\tilde{H})_{rj} (\tilde{H}^\top)_{jr} + \lambda_2 = (\tilde{W}^\top \tilde{W})_{pp} (\tilde{H} \tilde{H}^\top)_{rr} + \lambda_2$$

$$s^* = \frac{(\tilde{W}^\top V \tilde{H}^\top - \tilde{W}^\top \tilde{W} \tilde{M} \tilde{H} \tilde{H}^\top)_{pr} - \lambda_1 - \lambda_2(\tilde{M})_{pr}}{(\tilde{W}^\top \tilde{W})_{pp} (\tilde{H} \tilde{H}^\top)_{rr} + \lambda_2}$$

$$(\tilde{M})_{pr}^{t+1} = \max((\tilde{M}_0)_{pr}, \min(0, ((\tilde{M})_{pr})^t + s^*)).$$

References

- Samuel K Ainsworth, Nicholas J Foti, Adrian KC Lee, and Emily B Fox. oi-vae: Output interpretable vaes for nonlinear group factor analysis. In *International Conference on Machine Learning*, pages 119–128, 2018.
- Oskar Allerbo and Rebecka Jörnsten. Flexible, non-parametric modeling using regularized neural networks. *arXiv preprint arXiv:2012.11369*, 2020.
- Andrzej Cichocki and Anh-Huy Phan. Fast local algorithms for large scale nonnegative matrix and tensor factorizations. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, 92(3):708–721, 2009.
- Christophe Croux, Peter Filzmoser, and Heinrich Fritz. Robust sparse principal component analysis. *Technometrics*, 55(2):202–214, 2013.
- Jason A Dabin, Alexander M Haimovich, Justin Mauger, and Annan Dong. Blind source separation with l1 regularized sparse autoencoder. In *2020 29th Wireless and Optical Communications Conference (WOCC)*, pages 1–5. IEEE, 2020.
- Lingling Guo, Ping Wu, Jinfeng Gao, and Siwei Lou. Sparse kernel principal component analysis via sequential approach for nonlinear process monitoring. *IEEE Access*, 7:47550–47563, 2019.
- Cho-Jui Hsieh and Inderjit S Dhillon. Fast coordinate descent methods with variable selection for non-negative matrix factorization. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1064–1072, 2011.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Mark A Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991.
- Zhihui Lai, Yong Xu, Qingcai Chen, Jian Yang, and David Zhang. Multilinear sparse principal component analysis. *IEEE Transactions on Neural Networks and Learning Systems*, 25(10):1942–1950, 2014.
- Deyu Meng, Qian Zhao, and Zongben Xu. Improve robustness of sparse pca by l1-norm maximization. *Pattern Recognition*, 45(1):487–497, 2012.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In *Conference on Learning Theory*, pages 1376–1401, 2015.
- Andrew Ng et al. Sparse autoencoder. *CS294A Lecture Notes*, 72(2011):1–19, 2011.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

- R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, 14(5):877–898, 1976.
- Ferdinando S Samaria and Andy C Harter. Parameterisation of a stochastic model for human face identification. In *Proceedings of 1994 IEEE Workshop on Applications of Computer Vision*, pages 138–142. IEEE, 1994.
- Simone Scardapane, Danilo Comminiello, Amir Hussain, and Aurelio Uncini. Group sparse regularization for deep neural networks. *Neurocomputing*, 241:81–89, 2017.
- Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998.
- Alex J Smola, Olvi L Mangasarian, and Bernhard Schölkopf. Sparse kernel feature analysis. In *Classification, Automation, and New Media*, pages 167–178. Springer, 2002.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Michael E Tipping. Sparse kernel principal component analysis. In *Advances in Neural Information Processing Systems*, pages 633–639, 2001.
- Duo Wang and Toshihisa Tanaka. Sparse kernel principal component analysis based on elastic net regularization. In *2016 International Joint Conference on Neural Networks (IJCNN)*, pages 3703–3708. IEEE, 2016.
- Shanshan Wu, Alex Dimakis, Sujay Sanghavi, Felix Yu, Daniel Holtmann-Rice, Dmitry Storcheus, Afshin Rostamizadeh, and Sanjiv Kumar. Learning a compressed sensing measurement matrix via gradient unrolling. In *International Conference on Machine Learning*, pages 6828–6839. PMLR, 2019.
- Jaehong Yoon and Sung Ju Hwang. Combined group and exclusive sparsity for deep neural networks. In *International Conference on Machine Learning*, pages 3958–3966, 2017.
- Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- Ron Zass and Amnon Shashua. Nonnegative sparse pca. In *Advances in Neural Information Processing Systems*, pages 1561–1568, 2007.
- Yang Zhou, Rong Jin, and Steven Chu-Hong Hoi. Exclusive lasso for multi-task feature selection. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 988–995, 2010.
- Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15(2):265–286, 2006.