

NEURAL NETWORKS FOR PRINCIPAL COMPONENT ANALYSIS: A NEW LOSS FUNCTION PROVABLY YIELDS ORDERED EXACT EIGENVECTORS

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we propose a new loss function for performing principal component analysis (PCA) using linear autoencoders (LAEs). Optimizing the standard L_2 loss results in a decoder matrix that spans the principal subspace of the sample covariance of the data, but fails to identify the exact eigenvectors. This downside originates from an invariance that cancels out in the global map. Here, we prove that our loss function eliminates this issue, i.e. the decoder converges to the exact ordered unnormalized eigenvectors of the sample covariance matrix. For this new loss, we establish that all local minima are global optima and also show that computing the new loss (and also its gradients) has the same order of complexity as the classical loss. We report numerical results on both synthetic simulations, and a real-data PCA experiment on MNIST (i.e., a $60,000 \times 784$ matrix), demonstrating our approach to be practically applicable and rectify previous LAEs' downsides.

1 INTRODUCTION

Ranking among the most widely-used and valuable statistical tools, Principal Component Analysis (PCA) represents a given set of data within a new orthogonal coordinate system in which the data are uncorrelated and the variance of the data along each orthogonal axis is successively ordered from the highest to lowest. The projection of data along each axis gives what are called principal components. Theoretically, eigendecomposition of the covariance matrix provides exactly such a transformation. For large data sets, however, classical decomposition techniques are infeasible and other numerical methods, such as least squares approximation schemes, are practically employed. An especially notable instance is the problem of dimensionality reduction, where only the largest principal components—as the best representative of the data—are desired. Linear autoencoders (LAEs) are one such scheme for dimensionality reduction that is applicable to large data sets.

An LAE with a single fully-connected and linear hidden layer, and Mean Squared Error (MSE) loss function can discover the linear subspace spanned by the principal components. This subspace is the same as the one spanned by the weights of the decoder. However, its failure to identify the exact principal directions. This is due to the fact that, when the encoder is transformed by some matrix, transforming the decoder by the inverse of that matrix will yield no change in the loss. In other words, the loss possesses a symmetry under the action of a group of invertible matrices, so that directions (and orderings/permutations thereto) will not be discriminated.

The early work of Bourlard & Kamp (1988) and Baldi & Hornik (1989) connected LAEs and PCA and demonstrated the lack of identifiability of principal components. Several methods for neural networks compute the exact eigenvectors (Rubner & Tavan, 1989; Xu, 1993; Kung & Diamantaras, 1990; Oja et al., 1992), but they depend on either particular network structures or special optimization methods. It was recently observed (Plaut, 2018; Kunin et al., 2019) that regularization causes the left singular vectors of the decoder to become the exact eigenvectors, but recovering them still requires an extra decomposition step. As Plaut (2018) point out, no existent method recovers the eigenvectors from an LAE in an optimization-independent way on a standard network — this work fills that void.

Moreover, analyzing the loss surface for various architectures of linear/non-linear neural networks is a highly active and prominent area of research (e.g. Baldi & Hornik (1989); Kulin et al. (2019); Pretorius et al. (2018); Frye et al. (2019)). Most of these works extend the results of Baldi & Hornik (1989) for shallow LAEs to more complex networks. However, most retain the original MSE loss, and they prove the same critical point characterization for their specific architecture of interest. Most notably Zhou & Liang (2018) extends the results of Baldi & Hornik (1989) to deep linear networks and shallow RELU networks. In contrast in this work we are going after a loss with better loss surface properties.

We propose a new loss function for performing PCA using LAEs. We show that with the proposed loss function, the decoder converges to the exact ordered unnormalized eigenvectors of the sample covariance matrix. The idea is simple: for identifying p principal directions we build up a total loss function as a sum of p squared error losses, where the i^{th} loss function identifies only the first i principal directions. This approach breaks the symmetry since minimizing the first loss results in the first principal direction, which forces the second loss to find the first and the second. This constraint is propagated through the rest of the losses, resulting in all p principal components being identified. For the new loss we prove that all local minima are global minima.

Consequently, the proposed loss function has both theoretical and practical implications. Theoretically, it provides better understanding of the loss surface. Specifically, we show that any critical point of our loss L is a critical point of the original MSE loss but not vice versa, and conclude that L eliminates those undesirable global minima of the original loss (i.e., exactly those which suffer from the invariance). Given that the set of critical points of L is a subset of critical points of MSE loss, many of the previous work on loss surfaces of more complex networks likely extend. In light of the removal of undesirable global minima through L , examining more complex networks is certainly a very promising direction.

As for practical consequences, we show that the loss and its gradients can be compactly vectorized so that their computational complexity is no different from the MSE loss. Therefore, the loss L can be used to perform PCA/SVD on large datasets using any method of optimization such as Stochastic Gradient Descent (SGD). Chief among the compellingly reasons to perform PCA/SVD using this method is that, in recent years, there has been unprecedented gains in the performance of very large SGD optimizations, with autoencoders in particular successfully handling larger numbers of high-dimensional training data (e.g., images). The loss function we offer is attractive in terms of parallelizability and distributability, and does not prescribe any single specific algorithm or implementation, so stands to continue to benefit from the arms race between SGD and its competitors.

More importantly, this single loss function (without an additional post hoc processing step) fits seamlessly into optimization pipelines (where SGD is but one instance). The result is that the loss allows for PCA/SVD computation as single optimization layer, akin to an instance of a fully differentiable building block in a NN pipeline Amos & Kolter (2017), potentially as part of a much larger network.

2 THE PROPOSED LOSS FUNCTION AND REVIEW OF FINAL RESULTS

Let $\mathbf{X} \in \mathbb{R}^{n \times m}$ and $\mathbf{Y} \in \mathbb{R}^{n \times m}$ be the input and output matrices, where m centered sample points, each n -dimensional, are stacked column-wise. Let $\mathbf{x}_j \in \mathbb{R}^n$ and $\mathbf{y}_j \in \mathbb{R}^n$ be the j^{th} sample input and output (i.e. the j^{th} column of \mathbf{X} and \mathbf{Y} , respectively). Define the loss function $L(\mathbf{A}, \mathbf{B})$ as

$$L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \sum_{j=1}^m \|\mathbf{y}_j - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}\mathbf{x}_j\|_2^2 = \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}\mathbf{X}\|_F^2 \quad (1)$$

where $\langle \cdot, \cdot \rangle_F$ and $\|\cdot\|_F$ are the Frobenius inner product and norm, $\mathbf{I}_{i;p}$ is a $p \times p$ matrix with all elements zero except the first i diagonal elements being one. (Or, equivalently, the matrix obtained by setting the last $p - i$ diagonal elements of a $p \times p$ identity matrix to zero, e.g. $\mathbf{I}_{2;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$.) In

what follows, we shall denote the transpose of matrix \mathbf{M} by \mathbf{M}' . Moreover, the matrices $\mathbf{A} \in \mathbb{R}^{n \times p}$, and $\mathbf{B} \in \mathbb{R}^{p \times n}$ can be viewed as the weights of the decoder and encoder parts of an LAE.

The results are based on the following standard assumptions that hold generically:

Assumption 1. For an input X and an output Y , let $\Sigma_{xx} := XX'$, $\Sigma_{xy} := XY'$, $\Sigma_{yx} := \Sigma'_{xy}$ and $\Sigma_{yy} = YY'$ be their sample covariance matrices. We assume

- The input and output data are centered (zero mean).
- Σ_{xx} , Σ_{xy} , Σ_{yx} and Σ_{yy} are positive definite (of full rank and invertible).
- The covariance matrix $\Sigma := \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$ is of full rank with n distinct eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_n$.
- The decoder matrix A has no zero columns.

Claim. The main result of this work proved in Theorem 2 is as follows:

If the above assumptions hold then all the local minima of $L(A, B)$ are achieved iff A and B are of the form

$$\begin{aligned} A &= U_{1:p} D_p \\ B &= D_p^{-1} U'_{1:p} \Sigma_{yx} \Sigma_{xx}^{-1}, \end{aligned}$$

where the i^{th} column of $U_{1:p}$ is the unit eigenvector of $\Sigma := \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$ corresponding to the i^{th} largest eigenvalue and D_p is a diagonal matrix with nonzero diagonal elements. In other words, A contains ordered unnormalized eigenvectors of Σ corresponding to the p largest eigenvalues. Moreover, all the local minima are global minima with the value of the loss function at those global minima being

$$L(A, B) = p \operatorname{Tr}(\Sigma_{yy}) - \sum_{i=1}^p (p - i + 1) \lambda_i,$$

where λ_i is the i^{th} largest eigenvalue of $\Sigma := \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$. In the case of autoencoder ($Y = X$): $\Sigma = \Sigma_{xx}$. Finally, while $L(A, B)$ in the given form contains $O(p)$ matrix products, we will show that it can be evaluated with constant (less than 7) matrix products independent of the value p .

3 NOTATION

In this paper, the underlying field is always \mathbb{R} , and positive semidefinite matrices are symmetric by definition. The following constant matrices are used extensively throughout. The matrices $T_p \in \mathbb{R}^{p \times p}$ and $S_p \in \mathbb{R}^{p \times p}$ are defined as

$$(T_p)_{ij} = (p - i + 1) \delta_{ij}, \text{ i.e. } T_p = \operatorname{diag}(p, p - 1, \dots, 1), \quad (2)$$

$$(S_p)_{ij} = p - \max(i, j) + 1, \text{ i.e. } S_p = \begin{bmatrix} p & p-1 & \dots & 2 & 1 \\ p-1 & p-1 & \dots & 2 & 1 \\ \vdots & \vdots & \ddots & 2 & 1 \\ 2 & 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } S_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}. \quad (3)$$

Another matrix that will appear in the formulation is $\hat{S}_p := T_p^{-1} S_p T_p^{-1}$. Clearly, the diagonal matrix T_p is positive definite. As shown in Lemma 2, S_p and \hat{S}_p are positive definite as well.

4 MAIN THEOREMS

The general strategy to prove the above claim is as follows. First the analytical gradients of the loss is derived in a matrix form in Propositions 1 and 2. We compare the gradients with that of the original Minimum Square Error (MSE) loss. Next, we analyze the loss surface by solving the gradient equations which yields the general structure of critical points based on the rank of the decoder matrix A . Next, we delineate several interesting properties of the critical points, notably, any critical point of the loss is also a critical point for the MSE loss but not the other way around. Finally, by performing second order analysis on the loss in Theorem 2 the exact equations for local minima are derived which is shown to be as claimed.

Let $\tilde{L}(\mathbf{A}, \mathbf{B})$ and $L(\mathbf{A}, \mathbf{B})$ be the original loss, and the proposed loss function, respectively, i.e.,

$$\begin{aligned} \tilde{L}(\mathbf{A}, \mathbf{B}) &:= \sum_{j=1}^m \|y_j - \mathbf{A}\mathbf{B}x_j\|_2^2 \\ &= \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2 \end{aligned} \quad \left| \quad \begin{aligned} L(\mathbf{A}, \mathbf{B}) &:= \sum_{i=1}^p \sum_{j=1}^m \|y_j - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}x_j\|_2^2 \\ &= \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}\mathbf{X}\|_F^2 \end{aligned}$$

The first step is to calculate the gradients with respect to \mathbf{A} and \mathbf{B} and set them to zero to derive the implicit expressions for the critical points. In order to do so, first, in Lemma 5, for a fixed \mathbf{A} , we derive the directional (Gateaux) derivative of the loss with respect to \mathbf{B} along an arbitrary direction $\mathbf{W} \in \mathbb{R}^{p \times n}$, denoted as $d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W}$, i.e.

$$d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} = \lim_{\|\mathbf{W}\|_F \rightarrow 0} \frac{L(\mathbf{A}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B})}{\|\mathbf{W}\|_F}.$$

As shown in the proof of the lemma, $d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W}$ is derived by writing the norm in the loss as an inner product, opening it up using linearity of inner product, dismiss second order terms in \mathbf{W} (i.e. $O(\|\mathbf{W}\|^2)$) and rearrange the result as the inner product between the gradient with respect to \mathbf{B} , and the direction \mathbf{W} , which yields

$$\begin{aligned} d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2 \text{Tr}(\mathbf{W}'(\mathbf{T}_p\mathbf{A}'\Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))\mathbf{B}\Sigma_{xx})) \\ &= -2\langle \mathbf{T}_p\mathbf{A}'\Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))\mathbf{B}\Sigma_{xx}, \mathbf{W} \rangle_F, \end{aligned} \quad (4)$$

where, \circ is the Hadamard product and the constant matrices \mathbf{T}_p and \mathbf{S}_p , were defined in the beginning. Second, the same process is done in Lemma 6, to derive $d_{\mathbf{A}}L(\mathbf{A}, \mathbf{B})\mathbf{V}$; the derivative of L with respect to \mathbf{A} in an arbitrary direction $\mathbf{V} \in \mathbb{R}^{n \times p}$, for a fixed \mathbf{B} , which is then set to zero to derive the implicit expressions for the critical points. The results are formally stated in the two following propositions.

Proposition 1. *For any fixed matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ the function $L(\mathbf{A}, \mathbf{B})$ is convex in the coefficients of \mathbf{B} and attains its minimum for any \mathbf{B} satisfying the equation*

$$(\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))\mathbf{B}\Sigma_{xx} = \mathbf{T}_p\mathbf{A}'\Sigma_{yx}, \quad (5)$$

where \circ is the Hadamard (element-wise) product operator, and \mathbf{S}_p and \mathbf{T}_p are constant matrices defined in the previous section. Further, if \mathbf{A} has no zero column, then $L(\mathbf{A}, \mathbf{B})$ is strictly convex in \mathbf{B} and has a unique minimum when the critical \mathbf{B} is

$$\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A}) = (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1}\mathbf{T}_p\mathbf{A}'\Sigma_{yx}\Sigma_{xx}^{-1}, \quad (6)$$

and in the autoencoder case it becomes

$$\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A}) = (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1}\mathbf{T}_p\mathbf{A}'. \quad (6')$$

The proof is given in appendix A.2.

Remark 1. Note that as long as \mathbf{A} has no zero column, $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is nonsingular (we will explain the reason soon). In practice, \mathbf{A} with zero columns can always be avoided by nudging the zero columns of \mathbf{A} during the gradient decent process.

Proposition 2. *For any fixed matrix $\mathbf{B} \in \mathbb{R}^{p \times n}$ the function $L(\mathbf{A}, \mathbf{B})$ is a convex function in \mathbf{A} . Moreover, for a fixed \mathbf{B} , the matrix \mathbf{A} that satisfies*

$$\mathbf{A}(\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx}\mathbf{B}')) = \Sigma_{yx}\mathbf{B}'\mathbf{T}_p \quad (7)$$

is a critical point of $L(\mathbf{A}, \mathbf{B})$.

The proof is given in appendix A.3.

The pair (\mathbf{A}, \mathbf{B}) is a critical point of L if they make $d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W}$ and $d_{\mathbf{A}}L(\mathbf{A}, \mathbf{B})\mathbf{V}$ zero for any pair of directions (\mathbf{V}, \mathbf{W}) . Therefore, the implicit equations for critical points are given below, next to the ones derived by Baldi & Hornik (1989) for $\tilde{L}(\mathbf{A}, \mathbf{B})$.

$$\begin{array}{l|l}
\text{For } \tilde{L}(\mathbf{A}, \mathbf{B}): & \text{For } L(\mathbf{A}, \mathbf{B}): \\
\mathbf{A}'\mathbf{A}\mathbf{B}\Sigma_{xx} = \mathbf{A}'\Sigma_{yx}, & (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))\mathbf{B}\Sigma_{xx} = \mathbf{T}_p\mathbf{A}'\Sigma_{yx}, \\
\mathbf{A}\mathbf{B}\Sigma_{xx}\mathbf{B}' = \Sigma_{yx}\mathbf{B}', & \mathbf{A}(\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx}\mathbf{B}')) = \Sigma_{yx}\mathbf{B}'\mathbf{T}_p.
\end{array}$$

Remark 2. Notice the similarity, and the difference only being the presence of Hadamard product by \mathbf{S}_p in the left and by diagonal \mathbf{T}_p in the right. Therefore, practically, the added computational cost of evaluating the gradients is negligible compare to that of MSE loss.

The next step is to determine the structure of (\mathbf{A}, \mathbf{B}) that satisfies the above equations, and find the subset of those solutions that account for local minima. For the original loss, the first expression ($\mathbf{A}'\mathbf{A}\mathbf{B}\Sigma_{xx} = \mathbf{A}'\Sigma_{yx}$) is used to solve for \mathbf{B} and put it in the second to derive an expression solely based on \mathbf{A} . Obviously, in order to solve the first expression for \mathbf{B} , two cases are considered separately: the case where \mathbf{A} is of full rank p , so $\mathbf{A}'\mathbf{A}$ is invertible, and the case of \mathbf{A} being of rank $r < p$. Here we do the same but there is a twist; for us there is only one case. The reason is as long as (not necessarily full rank) \mathbf{A} has no zero column, $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is positive definite and hence, invertible. This is discussed in detail in Lemma 2 and we briefly explain it here. As shown in the lemma, \mathbf{S}_p is positive definite and by Shur product theorem for any \mathbf{A} (of any rank), $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is positive semidefinite. However, as a result of Oppenheim inequality (Horn & Johnson (2012), Thm 7.8.16), that in our case translates to $\det(\mathbf{S}_p) \prod_i (\mathbf{A}'\mathbf{A})_{ii} \leq \det(\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))$, as long as \mathbf{A} has no zero column, $\prod_i (\mathbf{A}'\mathbf{A})_{ii} > 0$ and therefore, $\det(\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})) > 0$. Here, we assume \mathbf{A} of any rank $r \leq p$ has no zero column (since this can be easily avoided in practice) and consider $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ to be always invertible. Therefore, (\mathbf{A}, \mathbf{B}) define a critical point of losses \tilde{L} and L if

$$\begin{array}{l|l}
\text{For } \tilde{L}(\mathbf{A}, \mathbf{B}) \text{ and full rank } \mathbf{A}: & \text{For } L(\mathbf{A}, \mathbf{B}) \text{ and no zero column } \mathbf{A}: \\
\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A}) = (\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'\Sigma_{yx}\Sigma_{xx}^{-1}, & \mathbf{B} = \hat{\mathbf{B}}(\mathbf{A}) = (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1}\mathbf{T}_p\mathbf{A}'\Sigma_{yx}\Sigma_{xx}^{-1}, \\
\mathbf{A}\mathbf{B}\Sigma_{xx}\mathbf{B}' = \Sigma_{yx}\mathbf{B}', & \mathbf{A}(\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx}\mathbf{B}')) = \Sigma_{yx}\mathbf{B}'\mathbf{T}_p.
\end{array}$$

Before, we state the main theorem we need the following definitions. First, a rectangular permutation matrix $\mathbf{\Pi}_r \in \mathbb{R}^{r \times p}$ is a matrix that each column consists of at most one nonzero element with the value 1. If the rank of $\mathbf{\Pi}_r$ is r with $r < p$ then clearly, $\mathbf{\Pi}_r$ has $p - r$ zero columns. Also, by taking away those zero columns the resultant $r \times r$ submatrix of $\mathbf{\Pi}_r$ is a standard square permutation matrix.

Second, under the conditions provided in Assumption 1, the matrix $\Sigma := \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$ has an eigenvalue decomposition $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$, where the i^{th} column of \mathbf{U} , denoted as \mathbf{u}_i , is an eigenvector of Σ corresponding to the i^{th} largest eigenvalue of Σ , denoted as λ_i . Also, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the diagonal vector of ordered eigenvalues of Σ , with $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$. We use the following notation to organize a subset of eigenvectors of Σ into a rectangular matrix. Let for any $r \leq p$, $\mathbb{I}_r = \{i_1, \dots, i_r\} (1 \leq i_1 < \dots < i_r < n)$ be any ordered r -index set. Define $\mathbf{U}_{\mathbb{I}_r} \in \mathbb{R}^{n \times p}$ as $\mathbf{U}_{\mathbb{I}_r} = [\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_r}]$. That is the columns of $\mathbf{U}_{\mathbb{I}_r}$ are the ordered orthonormal eigenvectors of Σ associated with eigenvalues $\lambda_{i_1} < \dots < \lambda_{i_r}$. Clearly, when $r = p$, we have $\mathbf{U}_{\mathbb{I}_r} = [\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_p}]$ corresponding to an p -index set $\mathbb{I}_p = \{i_1, \dots, i_p\} (1 \leq i_1 < \dots < i_p < n)$. Similarly, we define $\mathbf{\Lambda}_{\mathbb{I}_r} \in \mathbb{R}^{p \times p}$ as $\mathbf{\Lambda}_{\mathbb{I}_r} = \text{diag}(\lambda_{i_1}, \dots, \lambda_{i_r})$.

Theorem 1. Let $\mathbf{A} \in \mathbb{R}^{n \times p}$ and $\mathbf{B} \in \mathbb{R}^{p \times n}$ such that \mathbf{A} is of rank $r \leq p$. Under the conditions provided in Assumption 1 and the above notation, The matrices \mathbf{A} and \mathbf{B} define a critical point of $L(\mathbf{A}, \mathbf{B})$ if and only if for any given r -index set \mathbb{I}_r , and a nonsingular diagonal matrix $\mathbf{D} \in \mathbb{R}^{r \times r}$, \mathbf{A} and \mathbf{B} are of the form

$$\mathbf{A} = \mathbf{U}_{\mathbb{I}_r}\mathbf{C}\mathbf{D}, \quad (8)$$

$$\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A}) = \mathbf{D}^{-1}\mathbf{\Pi}_C\mathbf{U}'_{\mathbb{I}_r}\Sigma_{yx}\Sigma_{xx}^{-1}, \quad (9)$$

where, $\mathbf{C} \in \mathbb{R}^{r \times p}$ is of full rank r with nonzero and normalized columns such that $\mathbf{\Pi}_C := (\mathbf{S}_p \circ (\mathbf{C}'\mathbf{C}))^{-1}\mathbf{T}_p\mathbf{C}'$ is a rectangular permutation matrix of rank r and $\mathbf{C}\mathbf{\Pi}_C = \mathbf{I}_r$. For all $1 \leq r \leq p$, such \mathbf{C} always exists. In particular, if matrix \mathbf{A} is of full rank p , i.e. $r = p$, the two given conditions on $\mathbf{\Pi}_C$ are satisfied iff the invertible matrix \mathbf{C} is any squared $p \times p$ permutation matrix $\mathbf{\Pi}$. In this case (\mathbf{A}, \mathbf{B}) define a critical point of $L(\mathbf{A}, \mathbf{B})$ iff they are of the form

$$\mathbf{A} = \mathbf{U}_{\mathbb{I}_p}\mathbf{\Pi}\mathbf{D}, \quad (10)$$

$$\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A}) = \mathbf{D}^{-1} \mathbf{\Pi}' \mathbf{U}'_{\mathbb{I}_p} \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1}. \quad (11)$$

The proof is given in appendix A.4.

Remark 3. The above theorem provides explicit equations for the critical points of the loss surface in terms of the rank of the decoder matrix \mathbf{A} and the eigenvectors of $\mathbf{\Sigma}$. This explicit structure allows us to further analyze the loss surface and its local/global minima.

Here, we provide a proof sketch for the above theorem to make the claims more clear. Again as a reminder, the EVD of $\mathbf{\Sigma} := \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1} \mathbf{\Sigma}_{xy}$ is $\mathbf{\Sigma} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}'$. For both \tilde{L} and L , the corresponding $\hat{\mathbf{B}}(\mathbf{A})$ is replaced by \mathbf{B} on the RHS of critical point equations. For the loss $L(\mathbf{A}, \mathbf{B})$, as shown in the proof of the theorem, results in the following identity

$$\mathbf{U}' \mathbf{A} \left(\mathbf{S}_p \circ \left(\hat{\mathbf{B}} \mathbf{\Sigma}_{xx} \hat{\mathbf{B}}' \right) \right) \mathbf{A}' \mathbf{U} = \mathbf{\Lambda} \mathbf{\Delta}, \quad (12)$$

where $\mathbf{\Delta} := \mathbf{U}' \mathbf{A} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p \mathbf{A}' \mathbf{U}$ is symmetric and positive semidefinite. The LHS of eq. (12) is symmetric so the RHS is symmetric too, so $\mathbf{\Lambda} \mathbf{\Delta} = (\mathbf{\Lambda} \mathbf{\Delta})' = \mathbf{\Delta}' \mathbf{\Lambda}' = \mathbf{\Delta} \mathbf{\Lambda}$. Therefore, $\mathbf{\Delta}$ commutes with the diagonal matrix of eigenvalues $\mathbf{\Lambda}$. Since eigenvalues are assumed to be distinct, $\mathbf{\Delta}$ has to be diagonal as well. By Lemma 2 $\mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p$ is positive definite and \mathbf{U} is an orthogonal matrix. Therefore, $r = \text{rank}(\mathbf{A}) = \text{rank}(\mathbf{\Delta}) = \text{rank}(\mathbf{U}' \mathbf{\Delta} \mathbf{U})$, which implies that the diagonal matrix $\mathbf{\Delta}$, has r nonzero and *positive* diagonal entries. There exists an r -index set \mathbb{I}_r , corresponding to the nonzero diagonal elements of $\mathbf{\Delta}$. Forming a diagonal matrix $\mathbf{\Delta}_{\mathbb{I}_r} \in \mathbb{R}^{r \times r}$ by filling its diagonal entries (in order) by the nonzero diagonal elements of $\mathbf{\Delta}$, we have

$$\begin{aligned} \mathbf{U} \mathbf{\Delta} \mathbf{U}' &= \mathbf{U}_{\mathbb{I}_r} \mathbf{\Delta}_{\mathbb{I}_r} \mathbf{U}'_{\mathbb{I}_r} \stackrel{\text{Def of } \mathbf{\Delta}}{\Longrightarrow} \\ \mathbf{A} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p \mathbf{A}' &= \mathbf{U}_{\mathbb{I}_r} \mathbf{\Delta}_{\mathbb{I}_r} \mathbf{U}'_{\mathbb{I}_r}, \end{aligned} \quad (13)$$

which indicates that the matrix \mathbf{A} has the same column space as $\mathbf{U}_{\mathbb{I}_r}$. Therefore, there exists a full rank matrix $\tilde{\mathbf{C}} \in \mathbb{R}^{r \times p}$ such that $\mathbf{A} = \mathbf{U}_{\mathbb{I}_r} \tilde{\mathbf{C}}$. Since \mathbf{A} has no zero column, $\tilde{\mathbf{C}}$ has no zero column. Further, by normalizing the columns of $\tilde{\mathbf{C}}$ we can write $\mathbf{A} = \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D}$, where $\mathbf{D} \in \mathbb{R}^{p \times p}$ is diagonal that contains the norms of columns of $\tilde{\mathbf{C}}$.

Baldi & Hornik (1989) did something similar for full rank \mathbf{A} for the loss \tilde{L} to derive $(\mathbf{A}_{\tilde{L}} = \mathbf{U}_{\mathbb{I}_p} \tilde{\mathbf{C}})$. But their $\tilde{\mathbf{C}}$ can be any invertible $p \times p$ matrix. However, in our case, the matrix $\mathbf{C} \in \mathbb{R}^{r \times p}$ corresponding to rank $r \leq p$ matrix \mathbf{A} , has to satisfy eq. (13) by replacing \mathbf{A} by $\mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D}$ and eq. (12) by replacing $\hat{\mathbf{B}}(\mathbf{A})$ by $\hat{\mathbf{B}}(\mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D})$. In the case of Baldi & Hornik (1989), for the original loss \tilde{L} , equations similar to eq. (13) and eq. (12) appear but they are satisfied trivially by any invertible matrix $\tilde{\mathbf{C}}$. Simplifying those equations by using $\mathbf{A} = \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D}$ after some algebraic manipulation results in the following two conditions for \mathbf{C} :

$$\mathbf{C} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{C}' \mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}' = \mathbf{\Delta}_{\mathbb{I}_r}, \text{ and} \quad (14)$$

$$\mathbf{C} (\mathbf{S}_p \circ ((\mathbf{S}_p \circ (\mathbf{C}' \mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}' \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{C} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{C}' \mathbf{C}))^{-1})) \mathbf{C}' = \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Delta}_{\mathbb{I}_r}. \quad (15)$$

As detailed in proof of Theorem 1, solving for \mathbf{C} leads to its specific structure as laid out in the theorem.

Remark 4. Note that when \mathbf{A} is of rank $r < p$ with no zero columns then the invariant matrix \mathbf{C} is not necessarily a rectangular permutation matrix but $\mathbf{\Pi}_{\mathbf{C}} := (\mathbf{S}_p \circ (\mathbf{C}' \mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}'$ is a rectangular permutation matrix with $\mathbf{C} \mathbf{\Pi}_{\mathbf{C}} = \mathbf{I}_r$. It is only when $r = p$ that the invariant matrix \mathbf{C} becomes a permutation matrix. Nevertheless, as we show in the following corollary, the global map is always $\forall r \leq p : \mathbf{G} = \mathbf{A} \mathbf{B} = \mathbf{U}_{\mathbb{I}_r} \mathbf{U}'_{\mathbb{I}_r} \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1}$. It is possible to find further structure (in terms of block matrices) for the invariant matrix \mathbf{C} when $r < p$. However, this is not necessary as we soon show that all rank deficient matrix \mathbf{A} s are saddle points for the loss and ideally should be passed by during the gradient decent process. Based on some numerical results our conjecture is that when $r < p$ the matrix \mathbf{C} can only start with a $r \times k$ rectangular permutation matrix of rank r with $r \leq k \leq p$ and the rest of $p - k$ columns of \mathbf{C} is arbitrary as long as none of the columns are identically zero.

Corollary 1. *Let (\mathbf{A}, \mathbf{B}) be a critical point of $L(\mathbf{A}, \mathbf{B})$ under the conditions provided in Assumption 1 and $\text{rank} \mathbf{A} = r \leq p$. Then the following are true*

1. The matrix $\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}'$ is a $p \times p$ diagonal matrix of rank r .

2. For all $1 \leq r \leq p$, for any critical pair (\mathbf{A}, \mathbf{B}) , the global map $\mathbf{G} := \mathbf{AB}$ becomes

$$\mathbf{G} = \mathbf{U}_{\mathbb{I}_r} \mathbf{U}'_{\mathbb{I}_r} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}. \quad (16)$$

For the autoencoder case ($\mathbf{Y} = \mathbf{X}$) the global map is simply $\mathbf{G} = \mathbf{U}_{\mathbb{I}_r} \mathbf{U}'_{\mathbb{I}_r}$.

3. (\mathbf{A}, \mathbf{B}) is also the critical point of the classical loss $\tilde{L}(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^p \|\mathbf{Y} - \mathbf{ABX}\|_F^2$.

The proof is given in appendix A.5.

Remark 5. The above corollary implies that $L(\mathbf{A}, \mathbf{B})$ not only does not add any extra critical point compare to the original loss $\tilde{L}(\mathbf{A}, \mathbf{B})$, it provides the same global map $\mathbf{G} := \mathbf{AB}$. It only limits the structure of the invariance matrix \mathbf{C} as described in Theorem 1 so that the decoder matrix \mathbf{A} can recover the exact eigenvectors of $\boldsymbol{\Sigma}$.

Lemma 1. The loss function $L(\mathbf{A}, \mathbf{B})$ can be written as

$$L(\mathbf{A}, \mathbf{B}) = p \operatorname{Tr}(\boldsymbol{\Sigma}_{yy}) - 2 \operatorname{Tr}(\mathbf{A} \mathbf{T}_p \mathbf{B} \boldsymbol{\Sigma}_{xy}) + \operatorname{Tr}(\mathbf{B}' (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \boldsymbol{\Sigma}_{xx}). \quad (17)$$

The above identity shows that the number of matrix operations required for computing the loss $L(\mathbf{A}, \mathbf{B})$ is constant and thereby independent of the value of p .

The proof is given in appendix A.6.

Theorem 2. Let $\mathbf{A}^* \in \mathbb{R}^{n \times p}$ and $\mathbf{B}^* \in \mathbb{R}^{p \times n}$ such that \mathbf{A}^* is of rank $r \leq p$. Under the conditions provided in Assumption 1, $(\mathbf{A}^*, \mathbf{B}^*)$ define a local minima of the proposed loss function iff they are of the form

$$\mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{D}_p \quad (18)$$

$$\mathbf{B}^* = \mathbf{D}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \quad (19)$$

where the i^{th} column of $\mathbf{U}_{1:p}$ is a unit eigenvector of $\boldsymbol{\Sigma} := \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$ corresponding the i^{th} largest eigenvalue and \mathbf{D}_p is a diagonal matrix with nonzero diagonal elements. In other words, \mathbf{A}^* contains ordered unnormalized eigenvectors of $\boldsymbol{\Sigma}$ corresponding to the p largest eigenvalues. Moreover, all the local minima are global minima with the value of the loss function at those global minima being

$$L(\mathbf{A}^*, \mathbf{B}^*) = p \operatorname{Tr}(\boldsymbol{\Sigma}_{yy}) - \sum_{i=1}^p (p - i + 1) \lambda_i, \quad (20)$$

where λ_i is the i^{th} largest eigenvalue of $\boldsymbol{\Sigma}$.

The proof is given in appendix A.7.

Remark 6. Finally, the second and third assumptions we made in the beginning in Assumption 1 can be relaxed by requiring only $\boldsymbol{\Sigma}_{xx}$ to be full rank. The output data can have a different dimension than the input. That is $\mathbf{Y} \in \mathbb{R}^{n \times m}$ and $\mathbf{X} \in \mathbb{R}^{n' \times m}$, where $n \neq n'$. The reason is that the given loss function structurally is very similar to MSE loss and can be represented as a Frobenius norm on the space of $n \times m$ matrices. In this case the covariance matrix $\boldsymbol{\Sigma} := \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$ is still $n \times n$. Clearly, for under-constrained systems with $n < n'$ the full rank assumption of $\boldsymbol{\Sigma}$ holds. For the overdetermined case, where $n' > n$ the second and third assumptions in Assumption 1 can be relaxed: we only require $\boldsymbol{\Sigma}_{xx}$ to be full rank since this is the only matrix that is inverted in the theorems. Note that if $p > \min(n', n)$ then $\boldsymbol{\Lambda}_{\mathbb{I}_p}$: the $p \times p$ diagonal matrix of eigenvalues of $\boldsymbol{\Sigma}$ for a p -index-set \mathbb{I}_p bounds to have some zeros and will be say rank $r < p$, which in turn, results in the encoder \mathbf{A} with rank r . However, the Theorem 1 is proved for encoder of any rank $r \leq p$. Finally, following theorem 2 then the first r columns of the encoder \mathbf{A} converges to ordered eigenvectors of $\boldsymbol{\Sigma}$ while the $p - r$ remaining columns span the kernel space of $\boldsymbol{\Sigma}$. Moreover, $\boldsymbol{\Sigma}$ need not to have distinct eigenvectors. In this case $\boldsymbol{\Delta}_{\mathbb{I}_r}$ becomes a block diagonal matrix, where the blocks correspond to identical eigenvalues $\boldsymbol{\Sigma}_{\mathbb{I}_r}$. In this case, the corresponding eigenvectors in \mathbf{A}^* are not unique but they span the respective eigenspace.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

LAEs with Two Loss functions We will verify the loss function $L(\mathbf{A}, \mathbf{B})$ defined in eq. (1) by setting the input matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$ equal to the output matrix $\mathbf{Y} \in \mathbb{R}^{n \times m}$ ($\mathbf{Y} = \mathbf{X}$), where the

linear autodecoder (LAE) becomes a solution to PCA. In order for comparison, we train another LAE using the MSE loss $\tilde{L}(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ defined as $\tilde{L}(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) = \left\| \mathbf{Y} - \tilde{\mathbf{A}}\tilde{\mathbf{B}}\mathbf{X} \right\|_F^2$, where $\mathbf{Y} = \mathbf{X}$ is also applied in our experiments.

The weights of networks are initialized to random numbers with a small enough standard deviation (10^{-7} in our case). We choose to use the Adam optimizer with a scheduled learning rate (starting from 10^{-3} and ending with 10^{-6} in our case), which empirically benefits the optimization process. The two training processes are stopped at the same iteration at which one of the models firstly finds all of the principal directions. As a side note, we feed all data samples to the network at one time with batch size equal to m , although mini-batch implementations are apparently amendable.

Evaluation Metrics We use the classical PCA approach to get the ground truth principal direction matrix $\mathbf{A}^* \in \mathbb{R}^{n \times p}$, by conducting Eigen Value Decomposition (EVD) to $\mathbf{X}\mathbf{X}' \in \mathbb{R}^{n \times n}$ or Singular Value Decomposition (SVD) to $\mathbf{X} \in \mathbb{R}^{n \times m}$. As a reminder, $\mathbf{A} \in \mathbb{R}^{n \times p}$ stands for the decoder weight matrix of an trained LAE given a loss function L . To measure the distance between \mathbf{A}^* and \mathbf{A} , we propose an absolute cosine similarity (ACS) matrix inspired by mutual coherence (Donoho et al., 2005), which is defined as:

$$\mathbf{ACS}_{ij} = \frac{|\langle \mathbf{A}_i^*, \mathbf{A}_j \rangle|}{\|\mathbf{A}_i^*\| \cdot \|\mathbf{A}_j\|}, \quad (21)$$

where $\mathbf{A}_i^* \in \mathbb{R}^{n \times 1}$ denotes the i^{th} ground truth principal direction, and $\mathbf{A}_j \in \mathbb{R}^{n \times 1}$ denotes the j^{th} column of the decoder \mathbf{A} , $i, j = 1, 2, \dots, p$. The elements of $\mathbf{ACS} \in \mathbb{R}^{p \times p}$ in eq. (21) take values between $[0, 1]$, measuring pair-wise similarity across two sets of vectors. The absolute value absorbs the sign ambiguity of principal directions.

The performances of LAEs are evaluated by defining the following metrics:

$$\mathbf{Ratio}_{TP} = \sum_{i=1}^p I[\mathbf{ACS}_{ii} > 1 - \epsilon] / p \quad (22)$$

$$\mathbf{Ratio}_{FP} = \sum_{\substack{i,j=1 \\ i \neq j}}^p I[\mathbf{ACS}_{ij} > 1 - \epsilon] / p, \text{ and} \quad (23)$$

$$\mathbf{Ratio}_{Total} = \mathbf{Ratio}_{TP} + \mathbf{Ratio}_{FP}, \quad (24)$$

where I is the indicator function and ϵ is a manual tolerance threshold ($\epsilon = 0.01$ in our case). If two vectors have absolute cosine similarity over $1 - \epsilon$, they are deemed equal. Considering some columns of decoder may be correct principal directions but not in the right order, we introduce \mathbf{Ratio}_{TP} and \mathbf{Ratio}_{FP} in eqs. (22) and (23) to check the ratio of correct in-place and out-of-place principal directions respectively. Then \mathbf{Ratio}_{Total} in eq. (24) measures the total ratio of the correctly obtained principal directions by the LAE regardless of the order.

Datasets As a proof-of-concept, both synthetic data and real data are used. For the synthetic data, 2000 zero-centered data samples are generated from a 1000-dimension zero mean multivariate normal distribution with the covariance matrix being $\text{diag}(\mathbb{N}_p)$. For the real data, we choose to use MNIST dataset (LeCun et al., 1998), which includes 60,000 grayscale handwritten digits images, each of dimension $28 \times 28 = 784$.

5.2 EVALUATION AND ANALYSIS

Synthetic Data Experiments In our experiment, p , the number of desired principal components (PCs), is set to 100, i.e. the dimension is to be reduced from 1000 to 100. Figures 1 and 2 demonstrate a few conclusions. First, during the training process, the *loss ratio* of both losses continuously decreases to 1, i.e. they both converge to the optimal loss value. However, when both get close enough, L require more iterations since the optimizer is forced to find the right directions: it fully converges only after it has found all the principal directions in the right order.

Second, using the loss L results in finding more and more correct principal directions, with \mathbf{Ratio}_{TP} continuously rising; and ultimately affords all correct and ordered principal directions,

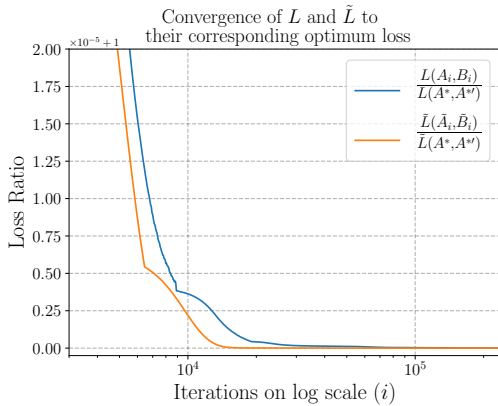


Figure 1: Convergence of losses to their corresponding optimal loss. Note that the correct shift and scaling of the y -axis tick values is printed at the top left corner of the figure.

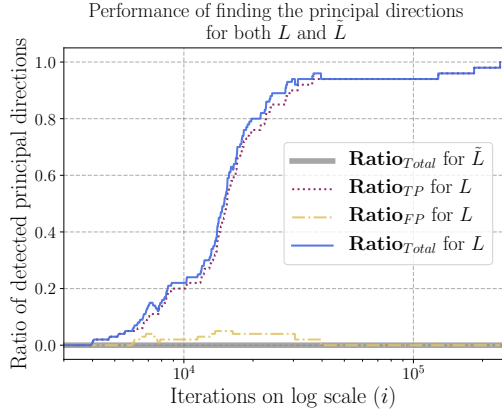


Figure 2: Performance of both losses L and \tilde{L} in finding the principal directions at the columns of their respective decoders.

with $Ratio_{TP}$ ending with 100%. Notice that occasionally and temporarily, some of the principal directions is found but not at their correct position, which is indicated by the rise of $Ratio_{FP}$ in the figure. However, as optimization continues they are shifted to the right column, which results in $Ratio_{FP}$ going back to zero, and $Ratio_{TP}$ reaching one. As for \tilde{L} , it fails to identify any principal directions; both $Ratio_{TP}$ and $Ratio_{FP}$ for \tilde{L} stay at 0, which indicates that none of the columns of the decoder \tilde{A} , aligns with any principal direction.

Third, as shown in the figure, while the optimizer finds almost all the principal directions rather quickly, it requires much more iterations to find some final ones. This is because some eigenvalues in the empirical covariance matrix of the finite 2000 samples become very close (the difference becomes less than 1). Therefore, the loss has to get very close to the optimal loss, making the gradient of the loss hard to distinguish between the two.

Real Data: MNIST Experiments We set the number of principal components (PCs) as 100, i.e., the dimension is to be reduced from 784 to 100. We also try to reconstruct with the top-10 columns found in this case. As in Fig. 3, the reconstruction performance of L is consistently better than \tilde{L} . That also reflects that \tilde{L} does not identify PCs, while L is directly applicable to performing PCA without bells and whistles.

6 CONCLUSION

In this paper, we have introduced a loss function for performing principal component analysis and linear regression using linear autoencoders. We have proved that the optimizing with the given loss results in the decoder matrix converges to the exact ordered unnormalized eigenvectors of the sample covariance matrix. We have also demonstrated the claims on a synthetic data set of random samples drawn from a multivariate normal distribution and on MNIST data set. There are several possible generalizations of this approach we are currently working on. One is improving performance when the corresponding eigenvalues of two principal directions are very close and another is generalization of the loss for tensor decomposition.

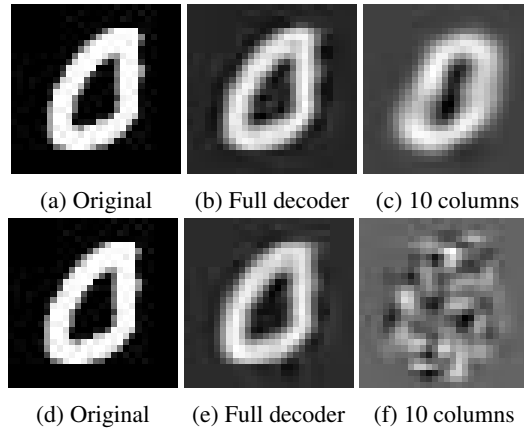


Figure 3: Real data experimental comparison in the reconstruction performance of MNIST images. First column: original image. Second column: reconstructed image using full columns of the decoder. Third column: reconstructed image using the first 10 columns of the decoder. Top row: using L . Bottom row: using \tilde{L} .

REFERENCES

- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *International Conference on Machine Learning*, pp. 136–145, 2017.
- Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- Hervé Bourlard and Yves Kamp. Auto-association by multilayer perceptrons and singular value decomposition. *Biological cybernetics*, 59(4-5):291–294, 1988.
- David L Donoho, Michael Elad, and Vladimir N Temlyakov. Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Transactions on information theory*, 52(1):6–18, 2005.
- Charles G Frye, Neha S Wadia, Michael R DeWeese, and Kristofer E Bouchard. Numerically recovering the critical points of a deep linear autoencoder. *arXiv preprint arXiv:1901.10603*, 2019.
- R.A. Horn and C.R. Johnson. *Matrix Analysis*. Cambridge University Press, 2012. ISBN 9781139788885.
- Sun-Yuan Kung and KI Diamantaras. A neural network learning algorithm for adaptive principal component extraction (apex). In *International Conference on Acoustics, Speech, and Signal Processing*, pp. 861–864. IEEE, 1990.
- Daniel Kunin, Jonathan M Bloom, Aleksandrina Goeva, and Cotton Seed. Loss landscapes of regularized linear autoencoders. *arXiv preprint arXiv:1901.08168*, 2019.
- Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Erkki Oja, Hidemitsu Ogawa, and Jaroonsakdi Wangviattana. Principal component analysis by homogeneous neural networks, part 1: The weighted subspace criterion. *IEICE Transactions on Information and Systems*, 75(3):366–375, 1992.
- Elad Plaut. From principal subspaces to principal components with linear autoencoders. *arXiv preprint arXiv:1804.10253*, 2018.
- Arnu Pretorius, Steve Kroon, and Herman Kamper. Learning dynamics of linear denoising autoencoders. In *International Conference on Machine Learning*, pp. 4138–4147, 2018.
- Jeanne Rubner and Paul Tavan. A self-organizing network for principal-component analysis. *EPL (Europhysics Letters)*, 10(7):693, 1989.
- Lei Xu. Least mean square error reconstruction principle for self-organizing neural-nets. *Neural networks*, 6(5):627–648, 1993.
- E. Zeidler. *Applied Functional Analysis: Main Principles and Their Applications*. Applied Mathematical Sciences. Springer New York, 1995. ISBN 9780387944227.
- Y Zhou and Y Liang. Critical points of linear neural networks: Analytical forms and landscape properties. In *Proc. Sixth International Conference on Learning Representations (ICLR)*, 2018.

APPENDIX

A PROOFS

A.1 PRELIMINARIES

Before we present the proof for the main theorems, the following two lemmas introduce some notations and basic relations that are required for the proofs.

Lemma 2. *The constant matrices $\mathbf{T}_p \in \mathbb{R}^{p \times p}$ and $\mathbf{S}_p \in \mathbb{R}^{p \times p}$ are defined as*

$$(\mathbf{T}_p)_{ij} = (p - i + 1) \delta_{ij}, \text{ i.e. } \mathbf{T}_p = \text{diag}(p, p - 1, \dots, 1),$$

$$(\mathbf{S}_p)_{ij} = p - \max(i, j) + 1, \text{ i.e. } \mathbf{S}_p = \begin{bmatrix} p & p-1 & \dots & 2 & 1 \\ p-1 & p-1 & \dots & 2 & 1 \\ \vdots & \vdots & \ddots & 2 & 1 \\ 2 & 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } \mathbf{S}_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

Clearly, the diagonal matrix \mathbf{T}_p is positive definite. Another matrix that will appear in the formulation is $\hat{\mathbf{S}}_p := \mathbf{T}_p^{-1} \mathbf{S}_p \mathbf{T}_p^{-1}$

$$(\hat{\mathbf{S}}_p)_{ij} = (\mathbf{T}_p^{-1} \mathbf{S}_p \mathbf{T}_p^{-1})_{ij} = \frac{1}{p - \min(i, j) + 1} \text{ i.e. } \mathbf{T}_p^{-1} \mathbf{S}_p \mathbf{T}_p^{-1} = \begin{bmatrix} \frac{1}{p} & \frac{1}{p} & \dots & \frac{1}{p} & \frac{1}{p} \\ \frac{1}{p} & \frac{1}{p-1} & \dots & \frac{1}{p-1} & \frac{1}{p-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{1}{p} & \frac{1}{p-1} & \dots & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{p} & \frac{1}{p-1} & \dots & \frac{1}{2} & 1 \end{bmatrix},$$

$$\text{e.g. } \hat{\mathbf{S}}_4 = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{3} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{3} & \frac{1}{2} & 1 \end{bmatrix}.$$

The following properties of Hadamard product and matrices \mathbf{T}_p and \mathbf{S}_p are used throughout:

1. For any arbitrary matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$,

$$\sum_{i=1}^p \mathbf{I}_{i:p} = \mathbf{T}_p, \text{ and} \quad (25)$$

$$\sum_{i=1}^p \mathbf{I}_{i:p} \mathbf{A}' \mathbf{A} \mathbf{I}_{i:p} = \mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}), \quad (26)$$

where, \circ is the Hadamard (element-wise) product.

2. For any matrices $\mathbf{M}_1, \mathbf{M}_2 \in \mathbb{R}^{p \times p}$ and diagonal matrices $\mathcal{D}, \mathcal{E} \in \mathbb{R}^{p \times p}$,

$$\mathcal{D} (\mathbf{M}_1 \circ \mathbf{M}_2) \mathcal{E} = (\mathcal{D} \mathbf{M}_1 \mathcal{E}) \circ \mathbf{M}_2 = \mathbf{M}_1 \circ (\mathcal{D} \mathbf{M}_2 \mathcal{E}).$$

Moreover, if $\mathbf{\Pi}_1, \mathbf{\Pi}_2 \in \mathbb{R}^{p \times p}$ are permutation matrices then

$$\mathbf{\Pi}_1 (\mathbf{M}_1 \circ \mathbf{M}_2) \mathbf{\Pi}_2 = (\mathbf{\Pi}_1 \mathbf{M}_1 \mathbf{\Pi}_2) \circ (\mathbf{\Pi}_1 \mathbf{M}_2 \mathbf{\Pi}_2).$$

3. \mathbf{S}_p is invertible and its inverse is a symmetric tridiagonal matrix

$$(\mathbf{S}_p^{-1})_{ij} = \begin{cases} 1 & i = j = 1 \\ 2 & i = j \neq 1 \\ -1 & |i - j| = 1 \\ 0 & \text{otherwise} \end{cases}, \text{ i.e. } \mathbf{S}_p^{-1} = \begin{bmatrix} 1 & -1 & \dots & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{bmatrix}.$$

4. \mathbf{S}_p is positive definite.
5. For any matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$, $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is positive semidefinite. If (not necessarily full rank) \mathbf{A} has no zero column then $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is positive definite.
6. For any diagonal matrix $\mathcal{D} \in \mathbb{R}^{p \times p}$

$$\mathbf{S}_p \circ \mathcal{D} = \mathbf{T}_p \mathcal{D}, \text{ and} \quad (27)$$

$$\hat{\mathbf{S}}_p \circ \mathcal{D} = \mathbf{T}_p^{-1} \mathcal{D}. \quad (28)$$

7. Let $\mathcal{D}, \mathcal{E} \in \mathbb{R}^{p \times p}$ be positive semidefinite matrices, where \mathcal{E} has no zero diagonal element, and \mathcal{D} is of rank $r \leq p$. Also, let for any $r \leq p$, $\mathbb{J}_r = \{i_1, \dots, i_r\} (1 \leq i_1 < \dots < i_r < n)$ be any ordered r -index set. Then \mathcal{D} and \mathcal{E} satisfy

$$\mathcal{E} \left(\hat{\mathbf{S}}_p \circ \mathcal{D} \right) = \left(\hat{\mathbf{S}}_p \circ \mathcal{E} \right) \mathcal{D},$$

if and only if, the following two conditions are satisfied:

- (a) The matrix \mathcal{D} is diagonal with $p - r$ zero diagonal elements and r positive diagonal elements indexed by the set \mathbb{J}_r . That is for any $i \in \mathbb{J}_r : (\mathcal{D})_{ii} > 0$ and the rest of elements of \mathcal{D} are zero.
- (b) For any $i, j \in \mathbb{J}_r$ and $i \neq j$ we have $(\mathcal{E})_{i,j} = 0$.

Clearly, if \mathcal{D} is positive definite then $\mathbb{J}_r = \mathbb{N}_p$ and hence, both \mathcal{D} and \mathcal{E} are diagonal.

Proof. . The proof of the properties are as follows.

1. eq. (25) is trivial. For eq. (26) note that $\mathbf{A}\mathbf{I}_{i,p}$ selects the first i columns of \mathbf{A} (zeros out the rest), and similarly, $\mathbf{I}_{i,p}\mathbf{A}'$ selects the first i rows of \mathbf{A} (zeros out the rest). Therefore, $\mathbf{I}_{i,p}\mathbf{A}'\mathbf{A}\mathbf{I}_{i,p}$ is a $p \times p$ matrix that its Leading Principal Submatrix of order i (LPS $_i$)¹ is the same as the LPS $_i$ of $\mathbf{A}'\mathbf{A}$ (and the rest of the elements are zero). Hence, $\sum_{i=1}^p \mathbf{I}_{i,p}\mathbf{A}'\mathbf{A}\mathbf{I}_{i,p}$ (counting backwards) adds LPS $_p$ of $\mathbf{A}'\mathbf{A}$ (i.e. $\mathbf{A}'\mathbf{A}$ itself) with LPS $_{p-1}$ that doubles LPS $_{p-1}$ part of the result and then adds LPS $_{p-2}$ that triples the LPS $_{p-2}$ part of result, the process continues until by the last addition LPS $_1$ is added to the result for the p^{th} times. This is exactly the same as evaluating $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$.
2. This is a standard result (Horn & Johnson, 2012), and no proof is needed.
3. Directly compute $\mathbf{S}_p \mathbf{S}_p^{-1}$:

$$\begin{aligned} (\mathbf{S}_p \mathbf{S}_p^{-1})_{ij} &= \sum_{k=1}^p (\mathbf{S}_p)_{ik} (\mathbf{S}_p^{-1})_{kj} \xrightarrow{\forall |k-j|>1: (\mathbf{S}_p^{-1})_{kj}=0} \\ &= \begin{cases} (\mathbf{S}_p)_{i,j-1} (\mathbf{S}_p^{-1})_{j-1,j} + (\mathbf{S}_p)_{i,j} (\mathbf{S}_p^{-1})_{j,j} + (\mathbf{S}_p)_{i,j+1} (\mathbf{S}_p^{-1})_{j+1,j} & 2 \leq j \text{ \&} \\ & j \leq p-1 \\ (\mathbf{S}_p)_{i,p-1} (\mathbf{S}_p^{-1})_{p-1,p} + (\mathbf{S}_p)_{i,p} (\mathbf{S}_p^{-1})_{p,p} & j = p \\ (\mathbf{S}_p)_{i,1} (\mathbf{S}_p^{-1})_{1,1} + (\mathbf{S}_p)_{i,2} (\mathbf{S}_p^{-1})_{2,1} & j = 1 \end{cases} \\ &= \begin{cases} -(\mathbf{S}_p)_{i,j-1} + 2(\mathbf{S}_p)_{i,j} - (\mathbf{S}_p)_{i,j+1} & 2 \leq j \leq p-1 \\ -(\mathbf{S}_p)_{i,p-1} + 2(\mathbf{S}_p)_{i,p} & j = p \\ (\mathbf{S}_p)_{i,1} - (\mathbf{S}_p)_{i,2} & j = 1 \end{cases} \\ &= \begin{cases} \max(i, j-1) - 2\max(i, j) + \max(i, j+1) & 2 \leq j \leq p-1 \\ -(p - \max(i, p-1) + 1) + 2(p - \max(i, p) + 1) & j = p \\ -\max(i, 1) + \max(i, 2) & j = 1 \end{cases} \end{aligned}$$

¹For a $p \times p$ matrix, the leading principal submatrix of order i is an $i \times i$ matrix derived by removing the last $p - i$ rows and columns of the original matrix (Horn & Johnson (2012), P17)

$$\begin{aligned}
&= \begin{cases} \max(i, j-1) - 2\max(i, j) + \max(i, j+1) & 2 \leq j \leq p-1 \\ 1-p + \max(i, p-1) & j = p \\ \max(i, 2) - \max(i, 1) & j = 1 \end{cases} \\
&= \begin{cases} \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} & 1 < j < p \\ \begin{cases} 1 & i = p \\ 0 & i \neq p \end{cases} & j = p \\ \begin{cases} 1 & i = 1 \\ 0 & i \geq 2 \end{cases} & j = 1 \end{cases} = (\mathbf{I}_p)_{ij}.
\end{aligned}$$

4. Firstly, note that \mathbf{S}_p^{-1} is symmetric and nonsingular so all the eigenvalues are real and nonzero. It is also a diagonally dominant matrix (Horn & Johnson (2012), Def 6.1.9) since

$$\forall i \in \{1, \dots, p\} : C_i := |(\mathbf{S}_p^{-1})_{ii}| \geq \sum_{j=1, j \neq i} |(\mathbf{S}_p^{-1})_{ij}| =: R_i,$$

where the inequality is strict for the first and the last row and it is equal for the rows in the middle. Moreover, by Gersgorin circle theorem (Horn & Johnson (2012), Thm 6.1.1) for every eigenvalue l_i of \mathbf{S}_p^{-1} there exists i such that $l_i \in [C_i - R_i, C_i + R_i]$. Since $\forall i : C_i \geq R_i$ we have all the eigenvalues are non-negative. They are also nonzero, hence, \mathbf{S}_p^{-1} is positive definite, which implies \mathbf{S}_p is also positive definite.

5. For any matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$, $\mathbf{A}'\mathbf{A}$ is positive semidefinite. Also, \mathbf{S}_p is positive definite so by Schur product theorem (Horn & Johnson (2012), Thm 7.5.3(a)), $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is positive semidefinite. Moreover, if all diagonal elements of $\mathbf{A}'\mathbf{A}$ are positive (i.e. \mathbf{A} has no zero column) by the extension of Schur product theorem (Horn & Johnson (2012), Thm 7.5.3(b)) it is positive definite. This can also be easily deduced using the Oppenheim inequality (Horn & Johnson (2012), Thm 7.8.16); that is for positive semidefinite matrices \mathbf{S}_p and $\mathbf{A}'\mathbf{A}$: $\det(\mathbf{S}_p) \prod_i (\mathbf{A}'\mathbf{A})_{ii} \leq \det(\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))$. Since, \mathbf{S}_p is positive definite, $\det(\mathbf{S}_p) > 0$ (in fact it is 1 for any p) and if $\mathbf{A}'\mathbf{A}$ has no zero diagonal then $\det(\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})) > 0$ and therefore, $\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})$ is positive definite.
6. Clearly, the matrix \mathbf{T}_p is achieved by setting the off-diagonal elements of \mathbf{S}_p to zero. Hence, for any diagonal matrix $\mathcal{D} \in \mathbb{R}^{p \times p}$: $\mathbf{S}_p \circ \mathcal{D} = \mathbf{T}_p \circ \mathcal{D}$. For the diagonal matrices Hadamard product and matrix product are interchangeable so the latter may also be written as $\mathbf{T}_p \mathcal{D}$. The same argument applies for the second identity.
7. This property can easily be proved by induction on p and careful bookkeeping of indices.

□

Lemma 3 (Simultaneous diagonalization by congruence). *Let $\mathbf{M}_1, \mathbf{M}_2 \in \mathbb{R}^{p \times p}$, where \mathbf{M}_1 is positive definite and \mathbf{M}_2 is positive semidefinite. Also, let $\mathcal{D}, \mathcal{E} \in \mathbb{R}^{r \times r}$ be positive definite diagonal matrices with $r \leq p$. Further, assume there is a $\mathbf{C} \in \mathbb{R}^{r \times p}$ of rank $r \leq p$ such that*

$$\begin{aligned}
\mathbf{C}\mathbf{M}_1\mathbf{C}' &= \mathcal{D} \text{ and} \\
\mathbf{C}\mathbf{M}_2\mathbf{C}' &= \mathcal{D}\mathcal{E}.
\end{aligned}$$

Then there exists a nonsingular $\bar{\mathbf{C}} \in \mathbb{R}^{p \times p}$ that its first r rows are the matrix \mathbf{C} and

$$\begin{aligned}
\bar{\mathbf{C}}\mathbf{M}_1\bar{\mathbf{C}}' &= \bar{\mathcal{D}} \text{ and} \\
\bar{\mathbf{C}}\mathbf{M}_2\bar{\mathbf{C}}' &= \bar{\mathcal{D}}\bar{\mathcal{E}},
\end{aligned}$$

where, $\bar{\mathcal{D}} = \mathcal{D} \oplus \mathbf{I}_{r-p}$ is a $p \times p$ diagonal matrix and $\bar{\mathcal{E}} = \mathcal{E} \oplus \underline{\mathcal{E}}$ is another $p \times p$ diagonal matrix, in which $\underline{\mathcal{E}} \in \mathbb{R}^{p-r \times p-r}$ is a nonnegative diagonal matrix. Clearly, the rank of \mathbf{M}_2 is r plus the number of nonzero diagonal elements of $\underline{\mathcal{E}}$.

Proof. The proof is rather straightforward since this lemma is the direct consequence of Theorem 7.6.4 in Horn & Johnson (2012). The theorem basically states that if $\mathbf{M}_1, \mathbf{M}_2 \in \mathbb{R}^{p \times p}$ is symmetric

and M_1 is positive definite then there exists an invertible $S \in \mathbb{R}^{p \times p}$ such that $SM_1S' = I_p$ and SM_2S' is a diagonal matrix with the same inertia as M_2 . Here, we have M_2 that is positive semidefinite and $C \in \mathbb{R}^{r \times p}$ of rank $r \leq p$ such that

$$\begin{aligned} \left(\mathcal{D}^{-\frac{1}{2}}C\right)M_1\left(\mathcal{D}^{-\frac{1}{2}}C\right)' &= I_r \text{ and} \\ \left(\mathcal{D}^{-\frac{1}{2}}C\right)M_2\left(\mathcal{D}^{-\frac{1}{2}}C\right)' &= \mathcal{E}. \end{aligned}$$

Therefore, since S is of full rank p and $\mathcal{D}^{-\frac{1}{2}}C$ is of rank $r \leq p$, there exists $p - r$ rows in S that are linearly independent of rows of $\mathcal{D}^{-\frac{1}{2}}C$. Establish $\bar{C} \in \mathbb{R}^{p \times p}$ by adding those $p - r$ rows to C . Then \bar{C} has p linearly independent rows so it is nonsingular, and fulfills the lemma's proposition that is

$$\begin{aligned} \bar{C}M_1\bar{C}' &= \bar{\mathcal{D}} \text{ and} \\ \bar{C}M_2\bar{C}' &= \bar{\mathcal{D}}\bar{\mathcal{E}}, \end{aligned}$$

where, $\bar{\mathcal{D}} = \mathcal{D} \oplus I_{r-p}$ is a $p \times p$ diagonal matrix and $\bar{\mathcal{E}} = \mathcal{E} \oplus \underline{\mathcal{E}}$ is another $p \times p$ diagonal matrix, in which $\underline{\mathcal{E}} \in \mathbb{R}^{p-r \times p-r}$ is a nonnegative diagonal matrix. \square

Lemma 4. Let A and B define a critical point of L . Further, let $V \in \mathbb{R}^{n \times p}$ and $W \in \mathbb{R}^{p \times n}$ are such that $\|V\|_F, \|W\|_F = O(\varepsilon)$ for some $\varepsilon > 0$. Then

$$\begin{aligned} L(A+V, B+W) - L(A, B) &= \langle VT_p B \Sigma_{xx} B', V \rangle_F \\ &\quad - 2 \langle \Sigma_{yx} W' T_p - A (S_p \circ (B \Sigma_{xx} W' + W \Sigma_{xx} B')), V \rangle_F \\ &\quad + \langle (S_p \circ (A' A)) W \Sigma_{xx}, W \rangle_F + O(\varepsilon^3). \end{aligned} \quad (29)$$

Further, for $W = \bar{W} := (S_p \circ (A' A))^{-1} T_p V' \Sigma_{yx} \Sigma_{xx}^{-1}$, the above equation becomes

$$\begin{aligned} L(A+V, B+\bar{W}) - L(A, B) &= \text{Tr}(V' V T_p B \Sigma_{xx} B') - \text{Tr}\left(V' \Sigma V T_p (S_p \circ (A' A))^{-1} T_p\right) \\ &\quad + 2 \text{Tr}\left(V' A \left(S_p \circ \left(B \Sigma_{xy} V T_p (S_p \circ (A' A))^{-1}\right.\right.\right. \\ &\quad \left.\left.\left.+ (S_p \circ (A' A))^{-1} T_p V' \Sigma_{yx} B'\right)\right)\right) + O(\varepsilon^3). \end{aligned} \quad (30)$$

Finally, in case the critical A is of full rank p and so, $(A, B) = (U_{\mathbb{I}_p} \Pi D, \hat{B}(U_{\mathbb{I}_p} \Pi D))$, for the encoder direction V with $\|V\|_F = O(\varepsilon)$ and $W = \bar{W}$ we have,

$$\begin{aligned} L(A+V, B+W) - L(A, B) &= \text{Tr}(V' V \Pi' \Lambda_{\mathbb{I}_p} \Pi T_p D^{-2}) - \text{Tr}(V' \Sigma V T_p D^{-2}) \\ &\quad + 2 \text{Tr}\left(V' U_{\mathbb{I}_p} \Pi D \left(S_p \circ \left(D^{-1} \Pi' U_{\mathbb{I}_p}' \Sigma V D^{-2}\right)\right)\right) \\ &\quad + 2 \text{Tr}\left(V' U_{\mathbb{I}_p} \Pi D \left(S_p \circ \left(D^{-2} V' \Sigma U_{\mathbb{I}_p} \Pi D^{-1}\right)\right)\right) \\ &\quad + O(\varepsilon^3). \end{aligned} \quad (31)$$

Proof. As described in appendix B.1, the second order Taylor expansion for the loss $L(A, B)$ is then given by eq. (63), i.e.

$$\begin{aligned} L(A+V, B+W) - L(A, B) &= d_A L(A, B) V + d_B L(A, B) W + \frac{1}{2} d_A^2 L(A, B) V^2 \\ &\quad + d_{AB} L(A, B) V W + \frac{1}{2} d_B^2 L(A, B) W^2 + R_{V, W}(A, B). \end{aligned}$$

If $\|V\|_F, \|W\|_F = O(\varepsilon)$ then $\|R(V, W)\| = O(\varepsilon^3)$. Moreover, when A and B define a critical point of L we have $d_A L(A, B) V = d_B L(A, B) W = 0$. By setting the derivatives $d_A^2 L(A, B) V^2$, $d_{AB} L(A, B) V W$, $d_B^2 L(A, B) W^2$ that are given by eq. (69), eq. (68), and eq. (66) respectively, the above equation simplifies to

$$\begin{aligned}
L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \langle \mathbf{V} (\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx}\mathbf{B}')), \mathbf{V} \rangle_F \\
&\quad - 2\langle \Sigma_{yx} \mathbf{W}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx} \mathbf{W}' + \mathbf{W}\Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \\
&\quad + \langle (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})) \mathbf{W}\Sigma_{xx}, \mathbf{W} \rangle_F + O(\varepsilon^3).
\end{aligned}$$

Now, based on the first item in Corollary 1, $\mathbf{B}\Sigma_{xx}\mathbf{B}'$ is a $p \times p$ diagonal matrix, so based on eq. (27): $\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx}\mathbf{B}') = \mathbf{T}_p \mathbf{B}\Sigma_{xx}\mathbf{B}'$. The substitution then yields eq. (29). Finally, in the above equation replace \mathbf{W} with $\bar{\mathbf{W}} = (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \Sigma_{yx} \Sigma_{xx}^{-1}$. We have

$$\begin{aligned}
L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \bar{\mathbf{W}}) - L(\mathbf{A}, \mathbf{B}) &= \\
&= \langle \mathbf{V} \mathbf{T}_p \mathbf{B}\Sigma_{xx}\mathbf{B}', \mathbf{V} \rangle_F - 2\langle \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p, \mathbf{V} \rangle_F \\
&\quad + 2\langle \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xx} \Sigma_{xx}^{-1} \Sigma_{xy} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} + (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \\
&\quad + \langle (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A})) (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xx}, (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \Sigma_{yx} \Sigma_{xx}^{-1} \rangle_F + O(\varepsilon^3) \\
&= \text{Tr}(\mathbf{V}' \mathbf{V} \mathbf{T}_p \mathbf{B}\Sigma_{xx}\mathbf{B}') - \text{Tr}(\mathbf{V}' \Sigma \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p) \\
&\quad + 2 \text{Tr}(\mathbf{V}' \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xy} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} + (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \Sigma_{yx} \mathbf{B}'))) + O(\varepsilon^3),
\end{aligned}$$

which is eq. (30). For the final equation, we have

$$\begin{aligned}
\mathbf{T}_p \mathbf{B}\Sigma_{xx}\mathbf{B}' &= \mathbf{T}_p \mathbf{D}^{-1} \Pi' \mathbf{U}'_{\mathbb{I}_p} \underbrace{\Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xx} \Sigma_{xx}^{-1} \Sigma_{xy}}_{\Sigma \mathbf{U}'_{\mathbb{I}_p} \Sigma} \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}^{-1} \\
&= \mathbf{T}_p \mathbf{D}^{-1} \Pi' \underbrace{\mathbf{U}'_{\mathbb{I}_p} \Sigma \mathbf{U}_{\mathbb{I}_p}}_{\Pi} \Pi \mathbf{D}^{-1} = \mathbf{T}_p \mathbf{D}^{-1} \underbrace{\Pi' \Lambda_{\mathbb{I}_p} \Pi}_{\Pi} \mathbf{D}^{-1} \\
&= \Pi' \Lambda_{\mathbb{I}_p} \Pi \mathbf{T}_p \mathbf{D}^{-2}, \text{ and} \tag{32}
\end{aligned}$$

$$\begin{aligned}
\mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p &= \mathbf{T}_p \left(\mathbf{S}_p \circ \left(\mathbf{D} \underbrace{\Pi' \mathbf{U}'_{\mathbb{I}_p} \mathbf{U}_{\mathbb{I}_p} \Pi}_{\mathbf{D}} \mathbf{D} \right) \right)^{-1} \mathbf{T}_p \\
&= \mathbf{T}_p (\mathbf{S}_p \circ \mathbf{D}^2)^{-1} \mathbf{T}_p = \mathbf{T}_p \mathbf{T}_p^{-1} \mathbf{D}^{-2} \mathbf{T}_p = \mathbf{T}_p \mathbf{D}^{-2}. \tag{33}
\end{aligned}$$

Replace the above in eq. (30) and simplify:

$$\begin{aligned}
L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\mathbf{V}' \mathbf{V} \mathbf{T}_p \mathbf{B}\Sigma_{xx}\mathbf{B}') - \text{Tr}(\mathbf{V}' \Sigma \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p) \\
&\quad + 2 \text{Tr}(\mathbf{V}' \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xy} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \\
&\quad + (\mathbf{S}_p \circ (\mathbf{A}'\mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \Sigma_{yx} \mathbf{B}'))) + O(\varepsilon^3) \stackrel{\text{eq. (32)}}{\underset{\text{eq. (33)}}{=}} \\
L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\mathbf{V}' \mathbf{V} \Pi' \Lambda_{\mathbb{I}_p} \Pi \mathbf{T}_p \mathbf{D}^{-2}) - \text{Tr}(\mathbf{V}' \Sigma \mathbf{V} \mathbf{T}_p \mathbf{D}^{-2}) \\
&\quad + 2 \text{Tr}(\mathbf{V}' \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B}\Sigma_{xy} \mathbf{V} \mathbf{D}^{-2} + \mathbf{D}^{-2} \mathbf{V}' \Sigma_{yx} \mathbf{B}'))) \\
&\quad + O(\varepsilon^3) \xrightarrow[\mathbf{B} = \hat{\mathbf{B}}(\mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D})]{\mathbf{A} = \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}} \\
L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\mathbf{V}' \mathbf{V} \Pi' \Lambda_{\mathbb{I}_p} \Pi \mathbf{T}_p \mathbf{D}^{-2}) - \text{Tr}(\mathbf{V}' \Sigma \mathbf{V} \mathbf{T}_p \mathbf{D}^{-2}) \\
&\quad + 2 \text{Tr}(\mathbf{V}' \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D} (\mathbf{S}_p \circ (\mathbf{D}^{-1} \Pi' \mathbf{U}'_{\mathbb{I}_p} \Sigma \mathbf{V} \mathbf{D}^{-2}))) \\
&\quad + 2 \text{Tr}(\mathbf{V}' \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D} (\mathbf{S}_p \circ (\mathbf{D}^{-2} \mathbf{V}' \Sigma \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}^{-1}))) \\
&\quad + O(\varepsilon^3),
\end{aligned}$$

which finalizes the proof. \square

A.2 PROOF OF PROPOSITION 1

For this proof we use the first and second order derivatives for $L(\mathbf{A}, \mathbf{B})$ wrt \mathbf{B} derived in Lemma 5. From eq. (66), we have that for a given \mathbf{A} the second derivative wrt to \mathbf{B} of the cost $L(\mathbf{A}, \mathbf{B})$ at \mathbf{B} ,

and in the direction \mathbf{W} is the quadratic form

$$d_{\mathbf{B}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{W}^2 = 2 \text{Tr}(\mathbf{W}' (\mathbf{S}_p \circ \mathbf{A}' \mathbf{A}) \mathbf{W} \Sigma_{xx}).$$

The matrix Σ_{xx} is positive-definite and by Lemma 2, $\mathbf{S}_p \circ \mathbf{A}' \mathbf{A}$ is positive-semidefinite. Hence, $d_{\mathbf{B}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{W}^2$ is clearly non-negative for all $\mathbf{W} \in \mathbb{R}^{p \times n}$. Therefore, $L(\mathbf{A}, \mathbf{B})$ is convex in coefficients of \mathbf{B} for a fixed matrix \mathbf{A} . Also the critical points of $L(\mathbf{A}, \mathbf{B})$ for a fixed \mathbf{A} is a matrix \mathbf{B} that satisfies $\forall \mathbf{W} \in \mathbb{R}^{p \times n} : d_{\mathbf{B}} L(\mathbf{A}, \mathbf{B}) \mathbf{W} = 0$ and hence, from eq. (64) we have

$$-2 \langle \mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx}, \mathbf{W} \rangle_F = 0.$$

Setting $\mathbf{W} = \mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx}$ we have

$$\mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx} = 0.$$

For a fixed \mathbf{A} , the cost $L(\mathbf{A}, \mathbf{B})$ is convex in \mathbf{B} , so any matrix \mathbf{B} that satisfies the above equation corresponds to a minimum of $L(\mathbf{A}, \mathbf{B})$. Further, if \mathbf{A} has no zero column then by Lemma 2, $\mathbf{S}_p \circ \mathbf{A}' \mathbf{A}$ is positive definite. Hence, $\forall \mathbf{W} \in \mathbb{R}^{p \times n} : d_{\mathbf{B}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{W}^2 = 2 \text{Tr}(\mathbf{W}' (\mathbf{S}_p \circ \mathbf{A}' \mathbf{A}) \mathbf{W} \Sigma_{xx})$ is positive. Therefore, the cost $L(\mathbf{A}, \mathbf{B})$ becomes strictly convex and the unique global minimum is achieved at $\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A})$ as defined in eq. (6).

A.3 PROOF OF PROPOSITION 2

For this proof we use the first and second order derivatives for $L(\mathbf{A}, \mathbf{B})$ wrt \mathbf{A} derived in Lemma 6. For a fixed \mathbf{B} , based on eq. (69) the second derivative wrt to \mathbf{A} of $L(\mathbf{A}, \mathbf{B})$ at \mathbf{A} , and in the direction \mathbf{V} is the quadratic form

$$d_{\mathbf{A}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{V}^2 = 2 \langle \mathbf{V} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F = 2 \text{Tr}(\mathbf{V} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')) \mathbf{V}').$$

The matrix Σ_{xx} is positive-definite and by Lemma 2, $\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')$ is positive-semidefinite. Hence, $d_{\mathbf{A}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{V}^2$ is non-negative for all $\mathbf{V} \in \mathbb{R}^{n \times p}$. Therefore, $L(\mathbf{A}, \mathbf{B})$ is convex in coefficients of \mathbf{A} for a fixed matrix \mathbf{B} . Based on eq. (67) the critical point of $L(\mathbf{A}, \mathbf{B})$ for a fixed \mathbf{B} is a matrix \mathbf{A} that satisfies for all $\mathbf{V} \in \mathbb{R}^{n \times p}$

$$\begin{aligned} d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} &= \langle -2 (\Sigma_{yx} \mathbf{B}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}'))), \mathbf{V} \rangle_F = 0 \implies \\ \Sigma_{yx} \mathbf{B}' \mathbf{T}_p &= \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \end{aligned}$$

which is eq. (7).

A.4 PROOF OF THEOREM 1

Before we start, a reminder on notation and some useful identities that are used throughout the proof. The matrix $\Sigma := \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$ has an eigenvalue decomposition $\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}'$, where the i^{th} column of \mathbf{U} , denoted as \mathbf{u}_i , is an eigenvector of Σ corresponding to the i^{th} largest eigenvalue of Σ , denoted as λ_i . Also, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the diagonal vector of ordered eigenvalues of Σ , with $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$. We use the following notation to organize a subset of eigenvectors of Σ into a rectangular matrix. Let for any $r \leq p$, $\mathbb{I}_r = \{i_1, \dots, i_r\} (1 \leq i_1 < \dots < i_r < n)$ be any *ordered* r -index set. Define $\mathbf{U}_{\mathbb{I}_r} \in \mathbb{R}^{n \times p}$ as $\mathbf{U}_{\mathbb{I}_r} = [\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_r}]$. That is the columns of $\mathbf{U}_{\mathbb{I}_r}$ are the ordered orthonormal eigenvectors of Σ associated with eigenvalues $\lambda_{i_1} < \dots < \lambda_{i_r}$. The following identities are then easy to verify:

$$\begin{aligned} \mathbf{U}_{\mathbb{I}_r}' \mathbf{U}_{\mathbb{I}_r} &= \mathbf{I}_r, \\ \Sigma \mathbf{U}_{\mathbb{I}_r} &= \mathbf{U}_{\mathbb{I}_r} \mathbf{\Lambda}_{\mathbb{I}_r}, \end{aligned} \tag{34}$$

$$\mathbf{U}_{\mathbb{I}_r}' \Sigma \mathbf{U}_{\mathbb{I}_r} = \mathbf{\Lambda}_{\mathbb{I}_r}. \tag{35}$$

The sufficient condition:

Let $\mathbf{A} \in \mathbb{R}^{n \times p}$ of rank $r \leq p$ and no zero column be given by eq. (8), $\mathbf{B} \in \mathbb{R}^{p \times n}$ given by eq. (9), and the accompanying conditions are met. Notice that $\mathbf{U}_{\mathbb{I}_r}' \mathbf{U}_{\mathbb{I}_r} = \mathbf{I}_r$ implies that $\mathbf{D} \mathbf{C}' \mathbf{C} \mathbf{D} = \mathbf{D} \mathbf{C}' \mathbf{U}_{\mathbb{I}_r}' \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D} = \mathbf{A}' \mathbf{A}$, so

$$\mathbf{B} = \mathbf{D}^{-1} \mathbf{\Pi}_C \mathbf{U}_{\mathbb{I}_r}' \Sigma_{yx} \Sigma_{xx}^{-1} \xrightarrow{\mathbf{\Pi}_C := (\mathbf{S}_p \circ (\mathbf{C}' \mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}'} \mathbf{D}^{-1} \mathbf{D} = \mathbf{I}_p$$

$$\begin{aligned}
B &= D^{-1} (S_p \circ (C' C))^{-1} D^{-1} D T_p C' U'_{\mathbb{I}_r} \Sigma_{yx} \Sigma_{xx}^{-1} \xrightarrow[\text{DT}_p = T_p D]{\text{Lemma 2-2}} \\
B &= \left(S_p \circ \underbrace{(DC' CD)} \right)^{-1} T_p \underbrace{DC' U'_{\mathbb{I}_r}} \Sigma_{yx} \Sigma_{xx}^{-1} \xrightarrow[\text{DC' CD} = A' A]{A' = D' C' U'_{\mathbb{I}_r}} \\
B &= (S_p \circ (A' A))^{-1} T_p A' \Sigma_{yx} \Sigma_{xx}^{-1} = \hat{B}(A),
\end{aligned}$$

which is eq. (6). Therefore, based on Proposition 1, for the given A , the matrix B defines a critical point of $L(A, B)$. For the gradient wrt to A , first note that with B given by eq. (9) we have

$$\begin{aligned}
B \Sigma_{xx} B' &= D^{-1} \Pi_C U'_{\mathbb{I}_r} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xx} \Sigma_{xx}^{-1} \Sigma_{xy} U_{\mathbb{I}_r} \Pi'_C D^{-1} \\
&= D^{-1} \Pi_C \underbrace{U'_{\mathbb{I}_r} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} U_{\mathbb{I}_r}} \Pi'_C D^{-1} \xrightarrow{\text{eq. (35)}} \\
B \Sigma_{xx} B' &= D^{-1} \Pi_C \Lambda_{\mathbb{I}_r} \Pi'_C D^{-1}. \tag{36}
\end{aligned}$$

The matrix Π_C is a rectangular permutation matrix so $\Pi_C \Lambda_{\mathbb{I}_r} \Pi'_C$ is diagonal so as $D^{-1} \Pi_C \Lambda_{\mathbb{I}_r} \Pi'_C D^{-1}$. Therefore, $B \Sigma_{xx} B'$ is diagonal and by eq. (27) in Lemma 2-6 we have

$$\begin{aligned}
S_p \circ (B \Sigma_{xx} B') &= T_p B \Sigma_{xx} B' = B \Sigma_{xx} B' T_p \\
&= D^{-1} \Pi_C \Lambda_{\mathbb{I}_r} \Pi'_C D^{-1} T_p \xrightarrow{A \times} \\
A (S_p \circ (B \Sigma_{xx} B')) &= A D^{-1} \Pi_C \Lambda_{\mathbb{I}_r} \Pi'_C D^{-1} T_p \xrightarrow[A = U_{\mathbb{I}_r} C D]{A = U_{\mathbb{I}_r} C D} \\
A (S_p \circ (B \Sigma_{xx} B')) &= U_{\mathbb{I}_r} C D D^{-1} \Pi_C \Lambda_{\mathbb{I}_r} \Pi'_C D^{-1} T_p \xrightarrow[A = U_{\mathbb{I}_r} C D]{A = U_{\mathbb{I}_r} C D} \\
&= U_{\mathbb{I}_r} \underbrace{C \Pi_C} \Lambda_{\mathbb{I}_r} \Pi'_C D^{-1} T_p \xrightarrow{C \Pi_C = I_r} \\
A (S_p \circ (B \Sigma_{xx} B')) &= \underbrace{U_{\mathbb{I}_r} \Lambda_{\mathbb{I}_r}} \Pi'_C D^{-1} T_p \xrightarrow{\text{eq. (34)}} \\
&= \Sigma U_{\mathbb{I}_r} \Pi'_C D^{-1} T_p \\
&= \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} U_{\mathbb{I}_r} \Pi'_C D^{-1} T_p \\
&= \Sigma_{yx} \underbrace{(D^{-1} \Pi_C U'_{\mathbb{I}_r} \Sigma_{yx} \Sigma_{xx}^{-1})}' T_p \\
&= \Sigma_{yx} B' T_p,
\end{aligned}$$

which is eq. (7). Therefore, based on Proposition Proposition 2, for the given B , the matrix A define a critical point of $L(A, B)$. Hence, A and B together define a critical point of $L(A, B)$.

The necessary condition:

Based on Proposition 1 and Proposition 2, for A (with no zero column) and B , to define a critical point of $L(A, B)$, B has to be $\hat{B}(A)$ given by eq. (6), and A has to satisfy eq. (7). That is

$$\begin{aligned}
A \left(S_p \circ \left(\hat{B} \Sigma_{xx} \hat{B}' \right) \right) &= \Sigma_{yx} \hat{B}' T_p \xrightarrow{\hat{B}(A) \text{ on RHS}} \\
A \left(S_p \circ \left(\hat{B} \Sigma_{xx} \hat{B}' \right) \right) &= \Sigma_{xy} \Sigma_{xx}^{-1} \Sigma_{yx} A T_p (S_p \circ (A' A))^{-1} T_p \xrightarrow[\Sigma = \Sigma_{xy} \Sigma_{xx}^{-1} \Sigma_{yx}]{\times A'} \\
A \left(S_p \circ \left(\hat{B} \Sigma_{xx} \hat{B}' \right) \right) A' &= \Sigma A T_p (S_p \circ (A' A))^{-1} T_p A' \xrightarrow[\times U, U' \times]{\Sigma = U \Lambda U''} \\
U' A \left(S_p \circ \left(\hat{B} \Sigma_{xx} \hat{B}' \right) \right) A' U &= U' U \Lambda U' A T_p (S_p \circ (A' A))^{-1} T_p A' U \xrightarrow{U' U = I_r} \\
U' A \left(S_p \circ \left(\hat{B} \Sigma_{xx} \hat{B}' \right) \right) A' U &= \Lambda \Delta, \tag{37}
\end{aligned}$$

where, $\Delta := U' A T_p (S_p \circ (A' A))^{-1} T_p A' U$ is symmetric and positive semidefinite. The LHS of the above equation is symmetric so the RHS is symmetric too, so $\Lambda \Delta = (\Lambda \Delta)' = \Delta' \Lambda' = \Delta \Lambda$. Therefore, Δ commutes with the diagonal matrix of eigenvalues Λ . Since, eigenvalues are assumed to be distinct, Δ has to be diagonal as well. By Lemma 2 $T_p (S_p \circ (A' A))^{-1} T_p$ is positive definite and U is an orthogonal matrix. Therefore, $r = \text{rank}(A) = \text{rank}(\Delta) = \text{rank}(U' \Delta U)$, which implies that the diagonal matrix Δ , has r nonzero and *positive* diagonal entries. There exists an

r -index set \mathbb{I}_r , corresponding to the nonzero diagonal elements of Δ . Forming a diagonal matrix $\Delta_{\mathbb{I}_r} \in \mathbb{R}^{r \times r}$ by filling its diagonal entries (in order) by the nonzero diagonal elements of Δ we have

$$\begin{aligned} U\Delta U' &= U_{\mathbb{I}_r} \Delta_{\mathbb{I}_r} U_{\mathbb{I}_r}' \xrightarrow{\text{Def of } \Delta} \\ U U' A T_p (S_p \circ (A' A))^{-1} T_p A' U U' &= U_{\mathbb{I}_r} \Delta_{\mathbb{I}_r} U_{\mathbb{I}_r}' \xrightarrow{U U' = I_r} \\ A T_p (S_p \circ (A' A))^{-1} T_p A' &= U_{\mathbb{I}_r} \Delta_{\mathbb{I}_r} U_{\mathbb{I}_r}', \end{aligned} \quad (38)$$

which indicates that the matrix A has the same column space as $U_{\mathbb{I}_r}$. Therefore, there exists a full rank matrix $\tilde{C} \in \mathbb{R}^{r \times p}$ such that $A = U_{\mathbb{I}_r} \tilde{C}$. Since A has no zero column, \tilde{C} has no zero column. Further, by normalizing the columns of \tilde{C} we can write $A = U_{\mathbb{I}_r} C D$, where $D \in \mathbb{R}^{p \times p}$ is diagonal that contains the norms of columns of \tilde{C} . Therefore, A is exactly in the form given by eq. (8). The matrix C has to satisfy eq. (38) that is

$$\begin{aligned} A T_p (S_p \circ (A' A))^{-1} T_p A' &= U_{\mathbb{I}_r} \Delta_{\mathbb{I}_r} U_{\mathbb{I}_r}' \xrightarrow{A = U_{\mathbb{I}_r} C} \\ U_{\mathbb{I}_r} C D T_p (S_p \circ (A' A))^{-1} T_p D C' U_{\mathbb{I}_r}' &= U_{\mathbb{I}_r} \Delta_{\mathbb{I}_r} U_{\mathbb{I}_r}' \xrightarrow{\begin{smallmatrix} \times U_{\mathbb{I}_r}, U_{\mathbb{I}_r} \times \\ A' A = D C' C D \end{smallmatrix}} \\ C D T_p (S_p \circ (D C' C D))^{-1} T_p C' D &= \Delta_{\mathbb{I}_r} \xrightarrow{\text{Lemma 2-2}} \\ C T_p D D^{-1} (S_p \circ (C' C))^{-1} D^{-1} D T_p C' &= \Delta_{\mathbb{I}_r} \implies \\ C T_p (S_p \circ (C' C))^{-1} T_p C' &= \Delta_{\mathbb{I}_r}. \end{aligned} \quad (39)$$

Now that the structure of A has been identified, evaluate $\hat{B}(A)$ of eq. (6) by setting $A = U_{\mathbb{I}_r} C D$, that is

$$\begin{aligned} B &= \hat{B}(A) = (S_p \circ (A' A))^{-1} T_p A' \Sigma_{yx} \Sigma_{xx}^{-1} \\ &= (S_p \circ (D C' C D))^{-1} T_p D C' U_{\mathbb{I}_r}' \Sigma_{yx} \Sigma_{xx}^{-1} \xrightarrow{\text{Lemma 2-2}} \\ B &= D^{-1} (S_p \circ (C' C))^{-1} T_p C' U_{\mathbb{I}_r}' \Sigma_{yx} \Sigma_{xx}^{-1}, \end{aligned}$$

which by defining $\Pi_C := (S_p \circ (C' C))^{-1} T_p C'$ gives eq. (34) for B as claimed. While C has to satisfy eq. (39), A and B in the given form have to satisfy eq. (37) that provides another condition for C as follows. First, note that

$$\begin{aligned} S_p \circ (\hat{B} \Sigma_{xx} \hat{B}') &= S_p \circ (D^{-1} (S_p \circ (C' C))^{-1} T_p C' U_{\mathbb{I}_r}' \Sigma U_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1} D^{-1}) \\ &= S_p \circ (D^{-1} (S_p \circ (C' C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1} D^{-1}) \xrightarrow{\text{Lemma 2-2}} \\ &= D^{-1} (S_p \circ ((S_p \circ (C' C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1})) D^{-1} \end{aligned}$$

Now, replace A and B in eq. (37) by their respective identities that we just derived. Performing the same process for eq. (37) we have

$$\begin{aligned} U' A \left(S_p \circ (\hat{B} \Sigma_{xx} \hat{B}') \right) A' U &= \Lambda \Delta \xrightarrow{\begin{smallmatrix} A = U_{\mathbb{I}_r} C D \\ \times U', U \times \end{smallmatrix}} \\ U_{\mathbb{I}_r}' C \left(S_p \circ ((S_p \circ (C' C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1}) \right) C' U_{\mathbb{I}_r}' &= U \Lambda \Delta U' \xrightarrow{\begin{smallmatrix} \times U_{\mathbb{I}_r} \\ U_{\mathbb{I}_r}' \times \end{smallmatrix}} \\ C \left(S_p \circ ((S_p \circ (C' C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1}) \right) C' &= U_{\mathbb{I}_r}' U \Lambda \Delta U' U_{\mathbb{I}_r} \implies \\ C \left(S_p \circ ((S_p \circ (C' C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1}) \right) C' &= \Lambda_{\mathbb{I}_r} \Delta_{\mathbb{I}_r}. \end{aligned} \quad (40)$$

Now we have to find C such that it satisfies eq. (39) and eq. (40). To make the process easier to follow, let's have them in one place. The matrix $C \in \mathbb{R}^{r \times p}$ have to satisfy

$$C T_p (S_p \circ (C' C))^{-1} T_p C' = \Delta_{\mathbb{I}_r} \text{ and} \quad (41)$$

$$C \left(S_p \circ ((S_p \circ (C' C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C' C))^{-1}) \right) C' = \Lambda_{\mathbb{I}_r} \Delta_{\mathbb{I}_r}. \quad (42)$$

Since C is a rectangular matrix, solving above equations for C in this form seems intractable. We use a trick to temporarily extend C into an invertible square matrix as follows.

- Temporarily, let $M_1 = T_p (S_p \circ (C'C))^{-1} T_p$, and $M_2 = S_p \circ ((S_p \circ (C'C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C'C))^{-1})$. Then M_1 is positive definite and M_2 is positive semidefinite, so they are simultaneously diagonalizable by congruence that is based on Lemma 3 and eq. (41) and eq. (42), there exists a nonsingular $\bar{C} \in \mathbb{R}^{p \times p}$ such that C consists of the first r rows of \bar{C} and

$$\bar{C} T_p (S_p \circ (C'C))^{-1} T_p \bar{C}' = \bar{\Delta}_{\mathbb{I}_r}, \quad (43)$$

$$\bar{C} \left(S_p \circ \left((S_p \circ (C'C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C'C))^{-1} \right) \right) \bar{C}' = \bar{\Lambda}_{\mathbb{I}_r} \bar{\Delta}_{\mathbb{I}_r}, \quad (44)$$

where, $\bar{\Delta}_{\mathbb{I}_r} = \Delta_{\mathbb{I}_r} \oplus I_{r-p}$ is a $p \times p$ diagonal matrix and $\bar{\Lambda}_{\mathbb{I}_r} = \Lambda_{\mathbb{I}_r} \oplus \underline{\Lambda}$ is another $p \times p$ diagonal matrix, in which $\underline{\Lambda} \in \mathbb{R}^{r-p \times r-p}$ is a nonnegative diagonal matrix.

- Substitute $\bar{\Delta}_{\mathbb{I}_r}$ from eq. (43) in eq. (44), then left multiply by \bar{C}'^{-1} , and right multiply by $\bar{C}' I_{r;p}$:

$$\begin{aligned} \bar{C}' \left(S_p \circ \left((S_p \circ (C'C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C'C))^{-1} \right) \right) \bar{C}' &= \\ \bar{\Lambda}_{\mathbb{I}_r} \bar{C} T_p (S_p \circ (C'C))^{-1} T_p \bar{C}' &\xrightarrow[\times \bar{C}'^{-1}]{\bar{C}' I_{r;p} \times} \\ \bar{C}' I_{r;p} \bar{C} \left(S_p \circ \left((S_p \circ (C'C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C'C))^{-1} \right) \right) &= \\ \bar{C}' I_{r;p} \bar{\Lambda}_{\mathbb{I}_r} \bar{C} T_p (S_p \circ (C'C))^{-1} T_p &. \end{aligned}$$

- Now we can revert back everything to C again. Since C consists of the first r rows of \bar{C} we have $\bar{C}' I_{r;p} \bar{C} = C'C$, and $\bar{C}' I_{r;p} \bar{\Lambda}_{\mathbb{I}_r} \bar{C} = C' \Lambda_{\mathbb{I}_r} C$, which turns the above equation into

$$\begin{aligned} C'C \left(S_p \circ \left(I_p (S_p \circ (C'C))^{-1} T_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C'C))^{-1} I_p \right) \right) &= \\ I_p C' \Lambda_{\mathbb{I}_r} C T_p (S_p \circ (C'C))^{-1} T_p &. \end{aligned}$$

- In the above equation, replace I_p by $T_p^{-1} T_p$ in LHS and by $T_p^{-1} (S_p \circ (C'C)) T_p^{-1} T_p (S_p \circ (C'C))^{-1} T_p$ in the RHS. Use $\Pi_C := (S_p \circ (C'C))^{-1} T_p C'$ to shrink it into :

$$C'C (S_p \circ (T_p^{-1} T_p \Pi_C \Lambda_{\mathbb{I}_r} \Pi_C' T_p T_p^{-1})) = T_p^{-1} (S_p \circ (C'C)) T_p^{-1} T_p \Pi_C \Lambda_{\mathbb{I}_r} \Pi_C' T_p.$$

- By the second property of Lemma 2 we can collect diagonal matrices T_p^{-1} 's around S_p to arrive at

$$(C'C) \left(\hat{S}_p \circ (T_p \Pi_C \Lambda_{\mathbb{I}_r} \Pi_C' T_p) \right) = \left(\hat{S}_p \circ (C'C) \right) (T_p \Pi_C \Lambda_{\mathbb{I}_r} \Pi_C' T_p),$$

where, $\hat{S}_p := T_p^{-1} S_p T_p^{-1}$.

- Define $p \times p$ matrices $\mathcal{E}_r := C'C$ and $\mathcal{D}_r := T_p \Pi_C \Lambda_{\mathbb{I}_r} \Pi_C' T_p$. Substitute in the above to arrive at:

$$\mathcal{E}_r \left(\hat{S}_p \circ \mathcal{D}_r \right) = \left(\hat{S}_p \circ \mathcal{E}_r \right) \mathcal{D}_r.$$

Both \mathcal{D}_r and \mathcal{E}_r in the above identity are positive semidefinite. Moreover, since by assumption C has no zero columns, \mathcal{E}_r has no zero diagonal element. Then the 7th property of Lemma 2 implies the following two conclusions:

- The matrix \mathcal{D}_r is diagonal. The rank of \mathcal{D}_r is r so it has exactly r positive diagonal elements and the rest is zero. This argument is true for $T_p^{-1} \mathcal{D}_r T_p^{-1} = \Pi_C \Lambda_{\mathbb{I}_r} \Pi_C'$. Since $\Lambda_{\mathbb{I}_r}$ is a diagonal positive definite matrix, the $p \times r$ matrix $\Pi_C := (S_p \circ (C'C))^{-1} T_p C'$ of rank r should have $p - r$ zero rows. Let \mathbb{J}_r be an r -index set corresponding to nonzero diagonal elements of $\Pi_C \Lambda_{\mathbb{I}_r} \Pi_C'$. Then the matrix $\Pi_C[\mathbb{J}_r, \mathbb{N}_r]$ ($r \times r$ submatrix of Π_C consist of its \mathbb{J}_r rows) is nonsingular.
- For every $i, j \in \mathbb{J}_r$ and $i \neq j$, $(\mathcal{E}_r)_{i,j} = 0$. Since $\mathcal{E}_r := C'C$ and so $(\mathcal{E}_r)_{i,j}$ is the inner product of i^{th} and j^{th} columns of C , we conclude that the columns of $C[\mathbb{N}_r, \mathbb{J}_r]$ ($r \times r$ submatrix of C consist of its \mathbb{J}_r columns) are orthogonal or in other words $C[\mathbb{N}_r, \mathbb{J}_r]' C[\mathbb{N}_r, \mathbb{J}_r]$ is diagonal. The columns of C are normalized. Therefore, $C[\mathbb{N}_r, \mathbb{J}_r]' C[\mathbb{N}_r, \mathbb{J}_r] = I_r$ and hence, $C[\mathbb{N}_r, \mathbb{J}_r]$ is an orthogonal matrix.

- We use the two conclusions to solve the original eq. (41) and eq. (42). First use $\mathbf{\Pi}_C := (\mathbf{S}_p \circ (\mathbf{C}'\mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}'$ to shrink them into :

$$\mathbf{C} \mathbf{T}_p \mathbf{\Pi}_C = \mathbf{\Delta}_{\mathbb{I}_r}, \quad (45)$$

$$\mathbf{C} (\mathbf{S}_p \circ (\mathbf{\Pi}_C \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C)) \mathbf{C}' = \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Delta}_{\mathbb{I}_r}. \quad (46)$$

Next, by the first conclusion, the matrix $\mathbf{T}_p^{-1} \mathcal{D}_r \mathbf{T}_p^{-1} = \mathbf{\Pi}_C \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C$ is diagonal and so eq. (46) becomes

$$\begin{aligned} \underbrace{\mathbf{C} \mathbf{T}_p \mathbf{\Pi}_C}_{\mathbf{\Delta}_{\mathbb{I}_r}} \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C \mathbf{C}' &= \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Delta}_{\mathbb{I}_r} \xrightarrow{\text{eq. (45)}} \\ \mathbf{\Delta}_{\mathbb{I}_r} \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C \mathbf{C}' &= \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Delta}_{\mathbb{I}_r} \implies \\ \mathbf{\Pi}'_C \mathbf{C}' &= \mathbf{C} \mathbf{\Pi}_C = \mathbf{I}_r, \end{aligned} \quad (47)$$

which is one of the two claimed conditions. What is left is to show that $\mathbf{\Pi}_C$ is a rectangular permutation matrix. From the first conclusion we also have $\mathbf{\Pi}_C$ has exactly r nonzero columns indexed by \mathbb{J}_r so

$$\mathbf{C}[\mathbb{N}_r, \mathbb{J}_r] \mathbf{\Pi}_C[\mathbb{J}_r, \mathbb{N}_r] = \mathbf{I}_r.$$

By the second conclusion $\mathbf{C}[\mathbb{N}_r, \mathbb{J}_r]$ is an orthogonal matrix therefore, $\mathbf{\Pi}_C[\mathbb{J}_r, \mathbb{N}_r]$ is the orthogonal matrix $\mathbf{C}[\mathbb{N}_r, \mathbb{J}_r]'$. Moreover, we had $\mathbf{T}_p^{-1} \mathcal{D}_r \mathbf{T}_p^{-1} = \mathbf{\Pi}_C \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C$ is a $p \times p$ diagonal matrix with exactly r nonzero diagonal elements. Hence, $\mathbf{\Pi}_C[\mathbb{N}_r, \mathbb{J}_r] \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C[\mathbb{N}_r, \mathbb{J}_r]$ is an $r \times r$ positive definite diagonal matrix with $\mathbf{\Lambda}_{\mathbb{I}_r}$ having distinct diagonal elements, and $\mathbf{\Pi}_C[\mathbb{N}_r, \mathbb{J}_r]$ being orthogonal. Therefore, $\mathbf{\Pi}_C[\mathbb{J}_r, \mathbb{N}_r]$ (as well as $\mathbf{C}[\mathbb{N}_r, \mathbb{J}_r]$) should be a square permutation matrix. Putting back the zero columns, we conclude that \mathbf{C} should be such that $\mathbf{\Pi}_C := (\mathbf{S}_p \circ (\mathbf{C}'\mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}'$ is a rectangular permutation matrix and $\mathbf{C} \mathbf{\Pi}_C = \mathbf{I}_r$. Note that it is possible to further analyze these conditions and determine the exact structure of \mathbf{C} . However, this is not needed in general for the critical point analysis of the next theorem except for the case where $r = p$ and \mathbf{C} is a square invertible matrix. In this case, square matrix $\mathbf{\Pi}_C$ is of full rank p , $\mathbb{J}_r = \mathbb{N}_p$ and therefore, $\mathbf{C}[\mathbb{N}_r, \mathbb{J}_r] = \mathbf{C}[\mathbb{N}_p, \mathbb{N}_p] = \mathbf{C}$. Hence, \mathbf{C} is any square permutation matrix $\mathbf{\Pi}$, $\mathbf{C}'\mathbf{C} = \mathbf{\Pi}'\mathbf{\Pi} = \mathbf{I}_p$ and $\mathbf{\Pi}_C := (\mathbf{S}_p \circ (\mathbf{C}'\mathbf{C}))^{-1} \mathbf{T}_p \mathbf{C}' = \mathbf{T}_p^{-1} \mathbf{T}_p \mathbf{\Pi}' = \mathbf{\Pi}'$, which verifies eq. (10) and eq. (11) for \mathbf{A} and \mathbf{B} when \mathbf{A} is of full rank p .

A.5 PROOF OF COROLLARY 1

1. We already show in the proof Theorem 1 that for critical (\mathbf{A}, \mathbf{B}) the matrix $\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}'$ is given by eq. (36) that is

$$\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}' = \mathbf{D}^{-1} \mathbf{\Pi}_C \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C \mathbf{D}^{-1}.$$

The matrix $\mathbf{\Pi}_C$ is a $p \times r$ rectangular permutation matrix so $\mathbf{\Pi}_C \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C$ is diagonal as well as $\mathbf{D}^{-1} \mathbf{\Pi}_C \mathbf{\Lambda}_{\mathbb{I}_r} \mathbf{\Pi}'_C \mathbf{D}^{-1}$. Therefore, $\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}'$ is diagonal. The diagonal matrix $\mathbf{\Lambda}_{\mathbb{I}_r}$ is of rank r therefore, $\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}'$ is of rank r .

2. Again by Theorem 1 critical (\mathbf{A}, \mathbf{B}) is of the form given by eq. (8) and eq. (9) with the proceeding conditions on the invariance \mathbf{C} . Therefore, the global map is

$$\begin{aligned} \mathbf{G} &= \mathbf{A} \mathbf{B} = \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D} \mathbf{D}^{-1} \mathbf{\Pi}_C \mathbf{U}'_{\mathbb{I}_r} \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1} \\ &= \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{\Pi}_C \mathbf{U}'_{\mathbb{I}_r} \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1} \xrightarrow{\mathbf{C} \mathbf{\Pi}_C = \mathbf{I}_r} \\ \mathbf{G} &= \mathbf{U}_{\mathbb{I}_r} \mathbf{U}'_{\mathbb{I}_r} \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1}. \end{aligned}$$

3. Based on Baldi & Hornik (1989) (\mathbf{A}, \mathbf{B}) define a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A} \mathbf{B} \mathbf{X}\|_F^2$ iff they satisfy

$$\mathbf{A}' \mathbf{A} \mathbf{B} \mathbf{\Sigma}_{xx} = \mathbf{A}' \mathbf{\Sigma}_{yx} \quad \text{and} \quad (48)$$

$$\mathbf{A} \mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}' = \mathbf{\Sigma}_{yx} \mathbf{B}'. \quad (49)$$

Again by assumption (\mathbf{A}, \mathbf{B}) define a critical point of $L(\mathbf{A}, \mathbf{B})$ so by Theorem 1 they are of the form given by eq. (8) and eq. (9) with the proceeding conditions on the invariance \mathbf{C} . Hence,

$$\mathbf{A}' \mathbf{A} \mathbf{B} \mathbf{\Sigma}_{xx} = \mathbf{D} \mathbf{C}' \underbrace{\mathbf{U}'_{\mathbb{I}_r} \mathbf{U}_{\mathbb{I}_r}}_{\mathbf{I}_r} \mathbf{C} \underbrace{\mathbf{D} \mathbf{D}^{-1}}_{\mathbf{I}_r} \mathbf{\Pi}_C \mathbf{U}'_{\mathbb{I}_r} \mathbf{\Sigma}_{yx} \underbrace{\mathbf{\Sigma}_{xx}^{-1} \mathbf{\Sigma}_{xx}}_{\mathbf{I}_r}$$

$$= DC' \underbrace{C\Pi_C}_{U'_{\mathbb{I}_r} \Sigma_{yx} \xrightarrow{C\Pi_C = I_r}} U'_{\mathbb{I}_r} \Sigma_{yx}$$

$$A'AB\Sigma_{xx} = DC'U'_{\mathbb{I}_r} \Sigma_{yx} = A' \Sigma_{yx}.$$

Hence, eq. (48) is satisfied. For the second equation we use the first property of this corollary that is $B\Sigma_{xx}B'$ is diagonal and satisfy eq. (7) of Proposition 2 that is

$$A(S_p \circ (B\Sigma_{xx}B')) = \Sigma_{yx}B'T_p \xrightarrow{B\Sigma_{xx}B' \text{ is diagonal}}$$

$$AT_pB\Sigma_{xx}B' = \Sigma_{yx}B'T_p \xrightarrow{B\Sigma_{xx}B' \text{ is diagonal}}$$

$$AB\Sigma_{xx}B'T_p = \Sigma_{yx}B'T_p \implies$$

$$AB\Sigma_{xx}B' = \Sigma_{yx}B'.$$

Hence, the second condition, eq. (49) is also satisfied. Therefore, any critical point of $L(\mathbf{A}, \mathbf{B})$ is a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B})$.

A.6 PROOF OF LEMMA 1

Proof. We have

$$L(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^p \|Y - AI_{i;p}BX\|_F^2 = \sum_{i=1}^p \langle Y - AI_{i;p}BX, Y - AI_{i;p}BX \rangle_F$$

$$= \sum_{i=1}^p (\langle Y, Y \rangle_F + \langle Y, -AI_{i;p}BX \rangle_F + \langle -AI_{i;p}BX, Y \rangle_F$$

$$+ \langle -AI_{i;p}BX, -AI_{i;p}BX \rangle_F)$$

$$= p\langle Y, Y \rangle_F - 2\langle Y, A \left(\sum_{i=1}^p I_{i;p} \right) BX \rangle_F + \sum_{i=1}^p \langle AI_{i;p}BX, AI_{i;p}BX \rangle_F \xrightarrow{\text{eq. (25)}}$$

$$= p\text{Tr}(YY') - 2\text{Tr}(AT_pBXY') + \sum_{i=1}^p \text{Tr}(X'B'I_{i;p}A'AI_{i;p}BX)$$

$$= p\text{Tr}(\Sigma_{yy}) - 2\text{Tr}(AT_pB\Sigma_{xy}) + \text{Tr} \left(XX'B' \sum_{i=1}^p (I_{i;p}A'AI_{i;p}) B \right) \xrightarrow{\text{eq. (26)}}$$

$$= p\text{Tr}(\Sigma_{yy}) - 2\text{Tr}(AT_pB\Sigma_{xy}) + \text{Tr}(B'(S_p \circ (A'A))B\Sigma_{xx}),$$

which is eq. (17). \square

A.7 PROOF OF THEOREM 2

Proof. The full rank matrices \mathbf{A}^* and \mathbf{B}^* given by eq. (18) and eq. (19) are clearly of the form given by Theorem 1 with $\mathbb{I}_p = \mathbb{N}_p := \{1, 2, \dots, p\}$, and $\Pi_p = I_p$. Hence, they define a critical point of $L(\mathbf{A}, \mathbf{B})$. We want to show that these are the only local minima, that is any other critical (\mathbf{A}, \mathbf{B}) is a saddle points. The proof is similar to the second partial derivative test. However, in this case the Hessian is a forth order tensor. Therefore, the second order Taylor approximation of the loss, derived in Lemma 4, is used directly. To prove the necessary condition, we show that at any other critical point (\mathbf{A}, \mathbf{B}) , where the first order derivatives are zero, there exists infinitesimal direction along which the second derivative of loss is negative. Next, for the sufficient condition we show that the any critical point of the form $(\mathbf{A}^*, \mathbf{B}^*)$ is a local and global minima.

The necessary condition:

Recall that $U_{\mathbb{I}_p}$ is the matrix of eigenvectors indexed by the p -index set \mathbb{I}_p and Π is a $p \times p$ permutation matrix. Since all the index sets \mathbb{I}_r , $r \leq p$ are assumed to be ordered, the only way to have $U_{\mathbb{N}_p} = U_{\mathbb{I}_p} \Pi$ is by having $\mathbb{I}_p = \mathbb{N}_p$ and $\Pi = I_p$. Let \mathbf{A} (with no zero column) and \mathbf{B} define an arbitrary critical point of $L(\mathbf{A}, \mathbf{B})$. Then Based on the previous theorem, either $\mathbf{A} = U_{\mathbb{I}_r} \mathbf{C}$ with $r < p$ or $\mathbf{A} = U_{\mathbb{I}_p} \Pi \mathbf{D}$ while in both cases $\mathbf{B} = \hat{\mathbf{B}}(\mathbf{A})$ given by eq. (6). If (\mathbf{A}, \mathbf{B}) is not of the form of $(\mathbf{A}^*, \mathbf{B}^*)$ then there are three possibilities either 1) $\mathbf{A} = U_{\mathbb{I}_r} \mathbf{C} \mathbf{D}$ with $r < p$, or 2)

$\mathbf{A} = \mathbf{U}_{\mathbb{I}_p} \mathbf{\Pi} \mathbf{D}$ with $\mathbb{I}_p \neq \mathbb{N}_p$ or 2) $\mathbf{A} = \mathbf{U}_{\mathbb{N}_p} \mathbf{\Pi} \mathbf{D}$ but $\mathbf{\Pi} \neq \mathbf{I}_p$. The first two cases corresponds to not having the ‘‘right’’ and/or ‘‘enough’’ eigenvectors, and the third corresponds to not having the ‘‘right’’ ordering. We introduce the following notation and investigate each case separately. Let $\varepsilon > 0$ and $\mathbf{U}_{i;j} \in \mathbb{R}^{n \times p}$ be a matrix of all zeros except the i^{th} column, which contains \mathbf{u}_j ; the eigenvector of $\mathbf{\Sigma}$ corresponding to the j^{th} largest eigenvalue. Therefore,

$$\mathbf{U}'_{i;j} \mathbf{\Sigma} \mathbf{U}_{i;j} = \mathbf{U}'_{i;j} \mathbf{U} \mathbf{\Lambda} \mathbf{U}' \mathbf{U}_{i;j} = \lambda_j \mathbf{E}_i, \quad (50)$$

where, $\mathbf{E}_i \in \mathbb{R}^{p \times p}$ is matrix of zeros except the i^{th} diagonal element that contains 1. In what follows, for each case we define a encoder direction $\mathbf{V} \in \mathbb{R}^{n \times p}$ with $\|\mathbf{V}\|_F = O(\varepsilon)$, and set the decoder direction $\mathbf{W} \in \mathbb{R}^{p \times n}$ as $\mathbf{W} = \bar{\mathbf{W}} := (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \mathbf{\Sigma}_{yx} \mathbf{\Sigma}_{xx}^{-1}$. Then we use eq. (30) and eq. (31) of Lemma 4, to show that the given direction (\mathbf{V}, \mathbf{W}) infinitesimally reduces the loss and hence, in every case the corresponding critical (\mathbf{A}, \mathbf{B}) is a saddle point.

1. For the case $\mathbf{A} = \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D}$, with $r < p$, note that based on the first item in Corollary 1, $\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}'$ is a $p \times p$ diagonal matrix of rank r so it has $p - r$ zero diagonal elements. Pick an $i \in \mathbb{N}_p$ such that $(\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}')_{ii}$ is zero and a $j \in \mathbb{N}_p \setminus \mathbb{I}_r$. Set $\mathbf{V} = \varepsilon \mathbf{U}_{i;j} \mathbf{D}$ and $\mathbf{W} = \bar{\mathbf{W}}$. Clearly,

$$\mathbf{V}' \mathbf{A} = \varepsilon \mathbf{D} \mathbf{U}'_{i;j} \mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D} = 0, \quad (51)$$

$$\begin{aligned} \mathbf{V}' \mathbf{V} \mathbf{T}_p \mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}' &= \varepsilon^2 \mathbf{D} \underbrace{\mathbf{U}'_{i;j} \mathbf{U}_{i;j}} \mathbf{D} \mathbf{T}_p \mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}', \\ &= \varepsilon^2 \mathbf{D} \mathbf{E}_i \mathbf{D} \mathbf{T}_p \mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}' = \varepsilon^2 \mathbf{D}^2 \mathbf{T}_p \mathbf{E}_i (\mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}') = 0 \text{ and} \end{aligned} \quad (52)$$

$$\mathbf{V}' \mathbf{\Sigma} \mathbf{V} = \varepsilon^2 \mathbf{D} \mathbf{U}'_{i;j} \mathbf{U} \mathbf{\Lambda} \mathbf{U}' \mathbf{U}_{i;j} \mathbf{D} = \varepsilon^2 \lambda_j \mathbf{D}^2 \mathbf{E}_i. \quad (53)$$

Notice, $\|\mathbf{V}\|_F, \|\mathbf{W}\|_F = O(\varepsilon)$, so based on eq. (30) of Lemma 4, we have

$$\begin{aligned} L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \\ &= \text{Tr}(\mathbf{V}' \mathbf{V} \mathbf{T}_p \mathbf{B} \mathbf{\Sigma}_{xx} \mathbf{B}') - \text{Tr}(\mathbf{V}' \mathbf{\Sigma} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p) \\ &+ 2 \text{Tr}(\mathbf{V}' \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \mathbf{\Sigma}_{xy} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} + (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p \mathbf{V}' \mathbf{\Sigma}_{yx} \mathbf{B}')) \\ &+ O(\varepsilon^3) \xrightarrow[\text{eq. (52)}]{\text{eq. (51)}} \\ L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \\ &= - \text{Tr}(\mathbf{V}' \mathbf{\Sigma} \mathbf{V} \mathbf{T}_p (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A}))^{-1} \mathbf{T}_p) + O(\varepsilon^3) \xrightarrow[\mathbf{A}' \mathbf{A} = \mathbf{D} \mathbf{C}' \mathbf{C} \mathbf{D}]{\text{eq. (53)}} \\ L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \\ &= - \varepsilon^2 \lambda_j \text{Tr} \left(\mathbf{D}^2 \mathbf{E}_i \mathbf{D}^{-1} \left(\left(\underbrace{\mathbf{T}_p^{-1} \mathbf{S}_p \mathbf{T}_p^{-1}} \right) \circ (\mathbf{C}' \mathbf{C}) \right)^{-1} \mathbf{D}^{-1} \right) + O(\varepsilon^3) = \\ &= - \varepsilon^2 \lambda_j \left(\left(\hat{\mathbf{S}}_p \circ (\mathbf{C}' \mathbf{C}) \right)^{-1} \right)_{ii} + O(\varepsilon^3). \end{aligned}$$

Therefore, since $\left(\hat{\mathbf{S}}_p \circ (\mathbf{C}' \mathbf{C}) \right)^{-1}$ is a positive definite matrix, as $\varepsilon \rightarrow 0$, we have $L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) \leq L(\mathbf{A}, \mathbf{B})$. Hence, any $(\mathbf{A}, \mathbf{B}) = (\mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D}, \hat{\mathbf{B}}(\mathbf{U}_{\mathbb{I}_r} \mathbf{C} \mathbf{D}))$ with $r < p$ is a saddle point.

2. Next, consider the case where $\mathbf{A} = \mathbf{U}_{\mathbb{I}_p} \mathbf{\Pi} \mathbf{D}$ with $\mathbb{I}_p \neq \mathbb{N}_p$. Then there exists at least one $j \in \mathbb{I}_p \setminus \mathbb{N}_p$ and $i \in \mathbb{N}_p \setminus \mathbb{I}_p$ such that $i < j$ (so $\lambda_i > \lambda_j$). Let σ be the permutation corresponding to the permutation matrix $\mathbf{\Pi}$. Also, let $\varepsilon > 0$ and $\mathbf{U}_{\sigma(j);i} \in \mathbb{R}^{n \times p}$ be a matrix of all zeros except the $\sigma(j)^{\text{th}}$ column, which contains \mathbf{u}_i ; the eigenvector of $\mathbf{\Sigma}$ corresponding to the i^{th} largest eigenvalue. Set $\mathbf{V} = \varepsilon \mathbf{U}_{\sigma(j);i} \mathbf{D}$ and $\mathbf{W} = \bar{\mathbf{W}}$. Then, since $i \notin \mathbb{I}_p$ we have

$$\mathbf{V}' \mathbf{U}_{\mathbb{I}_p} = \varepsilon \mathbf{D} \mathbf{U}'_{\sigma(j);i} \mathbf{U}_{\mathbb{I}_p} = 0, \quad (54)$$

$$\mathbf{V}' \mathbf{V} = \varepsilon^2 \mathbf{D} \mathbf{U}'_{\sigma(j);i} \mathbf{U}_{\sigma(j);i} \mathbf{D} = \varepsilon^2 \mathbf{D}^2 \mathbf{E}_{\sigma(j)}, \text{ and} \quad (55)$$

$$\mathbf{V}'\Sigma\mathbf{V} = \varepsilon^2 \mathbf{D}\mathbf{U}'_{\sigma(j);i}\mathbf{U}\Lambda\mathbf{U}'\mathbf{U}_{\sigma(j);i}\mathbf{D} = \varepsilon^2 \lambda_i \mathbf{D}^2 \mathbf{E}_{\sigma(j)}. \quad (56)$$

Since $\|\mathbf{V}\|_F, \|\mathbf{W}\|_F = O(\varepsilon)$, based on eq. (31) of Lemma 4, we have

$$\begin{aligned} L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\mathbf{V}'\mathbf{V}\Pi'\Lambda_{\mathbb{I}_p}\Pi\mathbf{T}_p\mathbf{D}^{-2}) - \text{Tr}(\mathbf{V}'\Sigma\mathbf{V}\mathbf{T}_p\mathbf{D}^{-2}) \\ &\quad + 2\text{Tr}\left(\mathbf{V}'\mathbf{U}_{\mathbb{I}_p}\Pi\mathbf{D}\left(\mathbf{S}_p \circ \left(\mathbf{D}^{-1}\Pi'\mathbf{U}'_{\mathbb{I}_p}\Sigma\mathbf{V}\mathbf{D}^{-2}\right)\right)\right) \\ &\quad + 2\text{Tr}\left(\mathbf{V}'\mathbf{U}_{\mathbb{I}_p}\Pi\mathbf{D}\left(\mathbf{S}_p \circ \left(\mathbf{D}^{-2}\mathbf{V}'\Sigma\mathbf{U}_{\mathbb{I}_p}\Pi\mathbf{D}^{-1}\right)\right)\right) \\ &\quad + O(\varepsilon^3) \xrightarrow[\text{eq. (55), eq. (56)}]{\text{eq. (54)}} \\ L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\varepsilon^2 \mathbf{D}^2 \mathbf{E}_{\sigma(j)} \Pi' \Lambda_{\mathbb{I}_p} \Pi \mathbf{T}_p \mathbf{D}^{-2}) \\ &\quad - \text{Tr}(\varepsilon^2 \lambda_i \mathbf{D}^2 \mathbf{E}_{\sigma(j)} \mathbf{T}_p \mathbf{D}^{-2}) + O(\varepsilon^3) \\ &= \varepsilon^2 \text{Tr}\left(\underbrace{\mathbf{E}_{\sigma(j)} \Pi' \Lambda_{\mathbb{I}_p} \Pi \mathbf{T}_p}_{\text{}}\right) - \varepsilon^2 \lambda_i \text{Tr}(\mathbf{E}_{\sigma(j)} \mathbf{T}_p) + O(\varepsilon^3) \\ &= \varepsilon^2 \lambda_j \text{Tr}(\mathbf{E}_{\sigma(j)} \mathbf{T}_p) - \varepsilon^2 \lambda_i \text{Tr}(\mathbf{E}_{\sigma(j)} \mathbf{T}_p) + O(\varepsilon^3) \\ &= -\varepsilon^2 (p - \sigma(j) + 1) (\lambda_i - \lambda_j) + O(\varepsilon^3). \end{aligned}$$

Note that in the above, the diagonal matrix $\Pi' \Lambda_{\mathbb{I}_p} \Pi$ has the same diagonal elements as $\Lambda_{\mathbb{I}_p}$ but they are permuted by σ . So $\mathbf{E}_{\sigma(j)} \Pi' \Lambda_{\mathbb{I}_p} \Pi$ selects $\sigma(j)$ th diagonal element of $\Pi' \Lambda_{\mathbb{I}_p} \Pi$ that is the j th diagonal element of $\Lambda_{\mathbb{I}_p}$, which is nothing but λ_j . Now, since $i < j$ so $\lambda_i > \lambda_j$ and $\sigma(j) \leq p$, as $\varepsilon \rightarrow 0$, we have $L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) \leq L(\mathbf{A}, \mathbf{B})$. Hence, any $(\mathbf{A}, \mathbf{B}) = (\mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}, \hat{\mathbf{B}}(\mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}))$ is a saddle point.

3. Finally consider the case where $\mathbf{A} = \mathbf{U}_{\mathbb{N}_p} \Pi \mathbf{D}$ with $\Pi \neq \mathbf{I}_p$. Since $\Pi \neq \mathbf{I}_p$, the permutation σ of the set \mathbb{N}_p , corresponding to the permutation matrix Π , has at least a cycle $(i_1 i_2 \cdots i_k)$, where $1 < i_1 < i_2 < \cdots < i_k < p$ and $2 \leq k \leq p$. Hence, Π can be decomposed as $\Pi = \Pi_{(i_1 i_2 \cdots i_k)} \hat{\Pi}$, where $\hat{\Pi}$ is the permutation matrix corresponding to other cycles of σ . The cycle $(i_1 i_2 \cdots i_k)$ can be decomposed into transpositions as $(i_1 i_2 \cdots i_k) = (i_k i_{k-1}) \cdots (i_k i_1)$, which in matrix form is $\Pi_{(i_1 i_2 \cdots i_k)} = \Pi_{(i_k i_{k-1})} \Pi_{(i_k i_2)} \cdots \Pi_{(i_k i_1)}$. Therefore, Π can be decomposed as $\Pi = \Pi_{(i_k i_1)} \tilde{\Pi}$, where $\tilde{\Pi} = \Pi_{(i_k i_2)} \cdots \Pi_{(i_k i_{k-1})} \hat{\Pi}$. Note that $\Pi_{(i_k i_1)}$, the permutation matrix corresponding to transposition $(i_k i_1)$ is a symmetric involutory matrix, i.e. $\Pi_{(i_k i_1)}^2 = \mathbf{I}_p$. Set $\mathbf{V} = \varepsilon(\mathbf{U}_{i_1; i_1} - \mathbf{U}_{i_k; i_k}) \tilde{\Pi} \mathbf{D}$ and $\mathbf{W} = \tilde{\mathbf{W}}$. Again we replace \mathbf{V} and \mathbf{W} in eq. (31) of Lemma 4. There are some tedious steps to simplify the equation, which is given in appendix A.7.1. The final result is as follows. With the given \mathbf{V} and \mathbf{W} , the third and fourth terms of the RHS of eq. (31) are canceled and the first two terms are simplified to

$$\text{Tr}(\mathbf{V}'\mathbf{V}\Pi'\Lambda_{\mathbb{N}_p}\Pi\mathbf{T}_p\mathbf{D}^{-2}) = \varepsilon^2 \lambda_{i_k} (p - i_1 + 1) + \varepsilon^2 \lambda_{i_1} (p - i_m + 1), \quad \text{and} \quad (57)$$

$$\text{Tr}(\mathbf{V}'\Sigma\mathbf{V}\mathbf{T}_p\mathbf{D}^{-2}) = \varepsilon^2 \lambda_{i_1} (p - i_1 + 1) + \varepsilon^2 \lambda_{i_k} (p - i_m + 1), \quad (58)$$

in which, $m = \max\{k - 1, 2\}$. This means that If the selected cycle is just a transposition $(i_1 i_2)$ then $i_m = i_2$. But if for the selected cycle $(i_1 i_2 \cdots i_k)$, k is greater than 2 then $i_m = i_{k-1}$. Using above equations, eq. (31) yields

$$\begin{aligned} L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\mathbf{V}'\mathbf{V}\Pi'\Lambda_{\mathbb{I}_p}\Pi\mathbf{T}_p\mathbf{D}^{-2}) - \text{Tr}(\mathbf{V}'\Sigma\mathbf{V}\mathbf{T}_p\mathbf{D}^{-2}) + O(\varepsilon^3) \\ &= \varepsilon^2 \lambda_{i_k} (p - i_1 + 1) + \varepsilon^2 \lambda_{i_1} (p - i_m + 1) \\ &\quad - \varepsilon^2 \lambda_{i_1} (p - i_1 + 1) - \varepsilon^2 \lambda_{i_k} (p - i_m + 1) + O(\varepsilon^3) \\ &= -\varepsilon^2 i_1 \lambda_{i_k} - \varepsilon^2 i_m \lambda_{i_1} + \varepsilon^2 i_1 \lambda_{i_1} + \varepsilon^2 i_m \lambda_{i_k} \\ &= -\varepsilon^2 ((\lambda_{i_1} - \lambda_{i_k})(i_m - i_1)) + O(\varepsilon^3). \quad (59) \end{aligned}$$

By the above definition of i_m , we have $i_m - i_1 > 0$ and since $i_1 < i_k$, $\lambda_{i_1} - \lambda_{i_k} > 0$. Hence, the first term in the above equation is negative and as $\varepsilon \rightarrow 0$, we have $L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) < 0$. Therefore, any any $(\mathbf{A}, \mathbf{B}) = (\mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}, \hat{\mathbf{B}}(\mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}))$ with $\Pi \neq \mathbf{I}_p$ is a saddle point.

The Sufficient condition:

From Lemma 1 we know that the loss $L(\mathbf{A}, \mathbf{B})$ can be written in the form of eq. (17). Use this equation to evaluate loss at $(\mathbf{A}^*, \mathbf{B}^*) = (\mathbf{U}_{\mathbb{N}_p} \mathbf{D}_p, \mathbf{D}_p^{-1} \mathbf{U}'_{\mathbb{N}_p} \Sigma_{yx} \Sigma_{xx}^{-1})$ as follows

$$\begin{aligned}
L(\mathbf{A}^*, \mathbf{B}^*) &= p \operatorname{Tr}(\Sigma_{yy}) - 2 \operatorname{Tr}(\mathbf{A}^* \mathbf{T}_p \mathbf{B}^* \Sigma_{xy}) + \operatorname{Tr}(\mathbf{B}^{*'} (\mathbf{S}_p \circ (\mathbf{A}^{*'} \mathbf{A}^*)) \mathbf{B}^* \Sigma_{xx}) \implies \\
L(\mathbf{A}^*, \mathbf{B}^*) &= p \operatorname{Tr}(\Sigma_{yy}) - 2 \operatorname{Tr}(\mathbf{U}_{\mathbb{N}_p} \mathbf{D}_p \mathbf{T}_p \mathbf{D}_p^{-1} \mathbf{U}'_{\mathbb{N}_p} \underbrace{\Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}}_{\Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}}) \\
&\quad + \operatorname{Tr}(\left(\mathbf{S}_p \circ \left(\mathbf{D}_p \underbrace{\mathbf{U}'_{\mathbb{N}_p} \mathbf{U}_{\mathbb{N}_p}}_{\Sigma_{xx}} \mathbf{D}_p \right) \right) \mathbf{D}_p^{-1} \mathbf{U}'_{\mathbb{N}_p} \underbrace{\Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xx} \Sigma_{xx}^{-1} \Sigma_{xy}}_{\Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}} \mathbf{U}_{\mathbb{N}_p} \mathbf{D}_p^{-1}) \implies \\
L(\mathbf{A}^*, \mathbf{B}^*) &= p \operatorname{Tr}(\Sigma_{yy}) - 2 \operatorname{Tr}(\mathbf{T}_p \mathbf{D}_p \mathbf{D}_p^{-1} \mathbf{U}'_{\mathbb{N}_p} \Sigma \mathbf{U}_{\mathbb{N}_p}) \\
&\quad + \operatorname{Tr}(\left(\mathbf{S}_p \circ (\mathbf{I}_p) \right) \mathbf{D}_p \mathbf{D}_p^{-1} \mathbf{U}'_{\mathbb{N}_p} \Sigma \mathbf{U}_{\mathbb{N}_p} \mathbf{D}_p^{-1} \mathbf{D}_p) \implies \\
L(\mathbf{A}^*, \mathbf{B}^*) &= p \operatorname{Tr}(\Sigma_{yy}) - 2 \operatorname{Tr}(\mathbf{T}_p \Lambda_{\mathbb{N}_p}) + \operatorname{Tr}(\mathbf{T}_p \Lambda_{\mathbb{N}_p}) \implies \\
L(\mathbf{A}^*, \mathbf{B}^*) &= p \operatorname{Tr}(\Sigma_{yy}) - \operatorname{Tr}(\mathbf{T}_p \Lambda_{\mathbb{N}_p}) = p \operatorname{Tr}(\Sigma_{yy}) - \sum_{i=1}^p (p - i + 1) \lambda_i,
\end{aligned}$$

which is eq. (20), as claimed. Notice that the above value is independent of the diagonal matrix \mathbf{D}_p . From the necessary condition we know that any critical point not in the form of $(\mathbf{A}^*, \mathbf{B}^*)$ is a saddle point. Hence, due to the convexity of the loss at least one $(\mathbf{A}^*, \mathbf{B}^*)$ is a global minimum but since the value of the loss at $(\mathbf{A}^*, \mathbf{B}^*)$ is independent of \mathbf{D}_p all these critical points yield the same value for the loss. Therefore, any critical point in the form of $(\mathbf{A}^*, \mathbf{B}^*)$ is a local and global minima. \square

A.7.1 SUPPLEMENTARY DETAILS OF THE PROOF OF THEOREM 2

To verify eq. (57), eq. (58), and eq. (59) in the proof of Theorem 2, we want to replace \mathbf{V} and \mathbf{W} in eq. (31) of Lemma 4 with $\mathbf{V} = \varepsilon(\mathbf{U}_{i_1:i_1} - \mathbf{U}_{i_k:i_k}) \tilde{\Pi} \mathbf{D}$ and $\mathbf{W} = \tilde{\mathbf{W}}$ and simplify. eq. (31) is

$$\begin{aligned}
L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= \operatorname{Tr}(\mathbf{V}' \mathbf{V} \Pi' \Lambda_{\mathbb{I}_p} \Pi \mathbf{T}_p \mathbf{D}^{-2}) - \operatorname{Tr}(\mathbf{V}' \Sigma \mathbf{V} \mathbf{T}_p \mathbf{D}^{-2}) \\
&\quad + 2 \operatorname{Tr}(\mathbf{V}' \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D} (\mathbf{S}_p \circ (\mathbf{D}^{-1} \Pi' \mathbf{U}'_{\mathbb{I}_p} \Sigma \mathbf{V} \mathbf{D}^{-2}))) \\
&\quad + 2 \operatorname{Tr}(\mathbf{V}' \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D} (\mathbf{S}_p \circ (\mathbf{D}^{-2} \mathbf{V}' \Sigma \mathbf{U}_{\mathbb{I}_p} \Pi \mathbf{D}^{-1}))) \\
&\quad + O(\varepsilon^3).
\end{aligned}$$

We investigate each term on the RHS separately. but before note that

$$\mathbf{E}_i \tilde{\Pi} \mathbf{T}_p \tilde{\Pi}' = \left(\tilde{\Pi} \mathbf{T}_p \tilde{\Pi}' \right)_{i,i} \mathbf{E}_i = (\mathbf{T}_p)_{\tilde{\sigma}^{-1}(i), \tilde{\sigma}^{-1}(i)} \mathbf{E}_i = (p - \tilde{\sigma}^{-1}(i) + 1) \mathbf{E}_i, \quad (60)$$

where, $\tilde{\sigma}$ and its function inverse $\tilde{\sigma}^{-1}$ are permutations corresponding to $\tilde{\Pi}$ and $\tilde{\Pi}'$ respectively. $\tilde{\Pi} \mathbf{T}_p \tilde{\Pi}'$ is a diagonal matrix where diagonal elements of \mathbf{T}_p are ordered based on $\tilde{\sigma}^{-1}$. Moreover, recall that we decomposed the permutation matrix Π in \mathbf{A} with a cycle $(i_1 i_2 \cdots i_k)$ as $\Pi = \Pi_{(i_1 i_k)} \underbrace{\Pi_{(i_k i_2)} \cdots \Pi_{(i_k i_{k-1})}}_{\tilde{\Pi}} \hat{\Pi} = \Pi_{(i_1 i_k)} \tilde{\Pi}$, where i_1, i_2, \cdots, i_k are fixed points of $\hat{\Pi}$. Therefore, with $\tilde{\sigma}$ being the permutation corresponding to $\tilde{\Pi}$ we have

$$\tilde{\sigma}(i_1) = i_1 \implies \tilde{\sigma}^{-1}(i_1) = i_1, \text{ and} \quad (61)$$

$$\tilde{\sigma}(i_{k-1}) = i_m \implies \tilde{\sigma}^{-1}(i_k) = i_m, \quad (62)$$

where, $m = \max\{k - 1, 2\}$. This means that If the selected cycle is just a transposition $(i_1 i_2)$ then $i_m = i_2$. But if for the selected cycle $(i_1 i_2 \cdots i_k)$, k is greater than 2 then $i_m = i_{k-1}$.

For the first term we have

$$\mathbf{V}' \mathbf{V} = \varepsilon^2 \mathbf{D} \tilde{\Pi}' (\mathbf{U}'_{i_1:i_1} - \mathbf{U}'_{i_k:i_k}) (\mathbf{U}_{i_1:i_1} - \mathbf{U}_{i_k:i_k}) \tilde{\Pi} \mathbf{D} \xrightarrow{\mathbf{U}'_{i_1:i_1} \mathbf{U}_{i_k:i_k} = 0}$$

$$\begin{aligned}
\mathbf{V}'\mathbf{V} &= \varepsilon^2 \mathbf{D}\tilde{\mathbf{\Pi}}'(U'_{i_1;i_1}U_{i_1;i_1} + U'_{i_k;i_k}U_{i_k;i_k})\tilde{\mathbf{\Pi}}\mathbf{D} \xrightarrow{U'_{i_1;i_1}U_{i_1;i_1}=\mathbf{E}_{i_1}} \\
&\xrightarrow{U'_{i_k;i_k}U_{i_k;i_k}=\mathbf{E}_{i_k}} \\
\mathbf{V}'\mathbf{V} &= \varepsilon^2 \mathbf{D}\tilde{\mathbf{\Pi}}'(\mathbf{E}_{i_1} + \mathbf{E}_{i_k})\tilde{\mathbf{\Pi}}\mathbf{D} \xrightarrow{\tilde{\mathbf{\Pi}}'(\mathbf{E}_{i_1} + \mathbf{E}_{i_k})\tilde{\mathbf{\Pi}} \text{ is diagonal}} \\
\mathbf{V}'\mathbf{V} &= \varepsilon^2 \tilde{\mathbf{\Pi}}'(\mathbf{E}_{i_1} + \mathbf{E}_{i_k})\tilde{\mathbf{\Pi}}\mathbf{D}^2 \implies \\
\text{Tr}(\mathbf{V}'\mathbf{V}\mathbf{\Pi}'\mathbf{\Lambda}_{\mathbb{N}_p}\mathbf{\Pi}\mathbf{T}_p\mathbf{D}^{-2}) &= \text{Tr}\left(\widetilde{\mathbf{V}'\mathbf{V}}\mathbf{D}^{-2}\tilde{\mathbf{\Pi}}'\mathbf{\Pi}_{(i_1 i_k)}\mathbf{\Lambda}_{\mathbb{N}_p}\mathbf{\Pi}_{(i_1 i_k)}\tilde{\mathbf{\Pi}}\mathbf{T}_p\right) \\
&= \text{Tr}\left(\varepsilon^2 \tilde{\mathbf{\Pi}}'(\mathbf{E}_{i_1} + \mathbf{E}_{i_k})\underbrace{\tilde{\mathbf{\Pi}}\mathbf{D}^2\mathbf{D}^{-2}\tilde{\mathbf{\Pi}}'}_{\mathbf{I}_p}\mathbf{\Pi}_{(i_1 i_k)}\mathbf{\Lambda}_{\mathbb{N}_p}\mathbf{\Pi}_{(i_1 i_k)}\tilde{\mathbf{\Pi}}\mathbf{T}_p\right) \\
&= \varepsilon^2 \text{Tr}\left((\mathbf{E}_{i_1} + \mathbf{E}_{i_k})\mathbf{\Pi}_{(i_1 i_k)}\mathbf{\Lambda}_{\mathbb{N}_p}\mathbf{\Pi}_{(i_1 i_k)}\tilde{\mathbf{\Pi}}\mathbf{T}_p\tilde{\mathbf{\Pi}}'\right) \\
&= \varepsilon^2 \text{Tr}\left(\lambda_{i_k}\mathbf{E}_{i_1}\tilde{\mathbf{\Pi}}\mathbf{T}_p\tilde{\mathbf{\Pi}}' + \lambda_{i_1}\mathbf{E}_{i_k}\tilde{\mathbf{\Pi}}\mathbf{T}_p\tilde{\mathbf{\Pi}}'\right) \xrightarrow{\text{eq. (60)}} \\
\text{Tr}(\mathbf{V}'\mathbf{V}\mathbf{\Pi}'\mathbf{\Lambda}_{\mathbb{N}_p}\mathbf{\Pi}\mathbf{T}_p\mathbf{D}^{-2}) &= \varepsilon^2 \lambda_{i_k}(p - \tilde{\sigma}^{-1}(i_1) + 1)\mathbf{E}_{i_1} + \varepsilon^2 \lambda_{i_1}(p - \tilde{\sigma}^{-1}(i_k) + 1)\mathbf{E}_{i_k} \xrightarrow{\text{eq. (61)}} \\
&\xrightarrow{\text{eq. (62)}} \\
\text{Tr}(\mathbf{V}'\mathbf{V}\mathbf{\Pi}'\mathbf{\Lambda}_{\mathbb{N}_p}\mathbf{\Pi}\mathbf{T}_p\mathbf{D}^{-2}) &= \varepsilon^2 \lambda_{i_k}(p - i_1 + 1)\mathbf{E}_{i_1} + \varepsilon^2 \lambda_{i_1}(p - i_m + 1)\mathbf{E}_{i_k}, \\
\end{aligned}$$

which is eq. (57) as claimed.

For the second term we have

$$\begin{aligned}
\mathbf{V}'\Sigma\mathbf{V} &= \varepsilon^2 \mathbf{D}\tilde{\mathbf{\Pi}}'(U'_{i_1;i_1} - U'_{i_k;i_k})\mathbf{U}\mathbf{\Lambda}\mathbf{U}'(U_{i_1;i_1} - U_{i_k;i_k})\tilde{\mathbf{\Pi}}\mathbf{D} \\
&= \varepsilon^2 \mathbf{D}\tilde{\mathbf{\Pi}}'\underbrace{(U'_{i_1;i_1}\mathbf{U}\mathbf{\Lambda}\mathbf{U}'U_{i_1;i_1})}_{\lambda_{i_1}\mathbf{E}_{i_1}} - \underbrace{U'_{i_1;i_1}\mathbf{U}\mathbf{\Lambda}\mathbf{U}'U_{i_k;i_k}}_0 \\
&\quad - \underbrace{U'_{i_k;i_k}\mathbf{U}\mathbf{\Lambda}\mathbf{U}'U_{i_1;i_1}}_0 + \underbrace{U'_{i_k;i_k}\mathbf{U}\mathbf{\Lambda}\mathbf{U}'U_{i_k;i_k}}_{\lambda_{i_k}\mathbf{E}_{i_k}}\tilde{\mathbf{\Pi}}\mathbf{D} \\
&= \varepsilon^2 \tilde{\mathbf{\Pi}}'(\lambda_{i_1}\mathbf{E}_{i_1} + \lambda_{i_k}\mathbf{E}_{i_k})\tilde{\mathbf{\Pi}}\mathbf{D}^2 \implies \\
\text{Tr}(\mathbf{V}'\Sigma\mathbf{V}\mathbf{T}_p\mathbf{D}^{-2}) &= \text{Tr}\left(\varepsilon^2 \tilde{\mathbf{\Pi}}'(\lambda_{i_1}\mathbf{E}_{i_1} + \lambda_{i_k}\mathbf{E}_{i_k})\tilde{\mathbf{\Pi}}\mathbf{D}^2\mathbf{T}_p\mathbf{D}^{-2}\right) \\
&= \varepsilon^2 \text{Tr}\left(\lambda_{i_1}\mathbf{E}_{i_1}\tilde{\mathbf{\Pi}}\mathbf{T}_p\tilde{\mathbf{\Pi}}' + \lambda_{i_k}\mathbf{E}_{i_k}\tilde{\mathbf{\Pi}}\mathbf{T}_p\tilde{\mathbf{\Pi}}'\right) \xrightarrow{\text{eq. (60)}} \\
\text{Tr}(\mathbf{V}'\Sigma\mathbf{V}\mathbf{T}_p\mathbf{D}^{-2}) &= \varepsilon^2 \lambda_{i_1}(p - \tilde{\sigma}^{-1}(i_1) + 1) + \varepsilon^2 \lambda_{i_k}(p - \tilde{\sigma}^{-1}(i_k) + 1) \xrightarrow{\text{eq. (61)}} \\
&\xrightarrow{\text{eq. (62)}} \\
\text{Tr}(\mathbf{V}'\Sigma\mathbf{V}\mathbf{T}_p\mathbf{D}^{-2}) &= \varepsilon^2 \lambda_{i_1}(p - i_1 + 1) + \varepsilon^2 \lambda_{i_k}(p - i_m + 1), \\
\end{aligned}$$

which is eq. (58) as claimed.

Finally, we have to show that the third and the fourth terms of the eq. (31) are canceled. First, observe that

$$\begin{aligned}
&\text{Tr}\left(\mathbf{V}'\mathbf{U}_{\mathbb{N}_p}\mathbf{\Pi}\mathbf{D}\left(\mathbf{S}_p \circ \left(\mathbf{D}^{-1}\mathbf{\Pi}'\mathbf{U}'_{\mathbb{N}_p}\Sigma\mathbf{V}\mathbf{D}^{-2}\right)\right)\right) = \\
&\text{Tr}\left(\varepsilon\mathbf{D}\tilde{\mathbf{\Pi}}'(U'_{i_1;i_1} - U'_{i_k;i_k})\mathbf{U}_{\mathbb{N}_p}\mathbf{\Pi}\left(\mathbf{S}_p \circ \left(\mathbf{\Pi}'\mathbf{U}'_{\mathbb{N}_p}\Sigma\mathbf{V}\mathbf{D}^{-2}\right)\right)\right) = \\
&\varepsilon \text{Tr}\left(\tilde{\mathbf{\Pi}}'(\mathbf{E}_{i_1} - \mathbf{E}_{i_k})\mathbf{\Pi}\left(\mathbf{S}_p \circ \left(\mathbf{\Pi}'\mathbf{U}'_{\mathbb{N}_p}\Sigma\mathbf{V}\mathbf{D}^{-2}\right)\right)\mathbf{D}\right) = \\
&\varepsilon^2 \text{Tr}\left(\tilde{\mathbf{\Pi}}'(\mathbf{E}_{i_1} - \mathbf{E}_{i_k})\mathbf{\Pi}\left(\mathbf{S}_p \circ \left(\mathbf{\Pi}'(\lambda_{i_1}\mathbf{E}_{i_1} - \lambda_{i_k}\mathbf{E}_{i_k})\tilde{\mathbf{\Pi}}\right)\right)\right) = \\
&\varepsilon^2 \text{Tr}\left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k})\left(\left(\mathbf{\Pi}\mathbf{S}_p\tilde{\mathbf{\Pi}}'\right) \circ \left(\mathbf{\Pi}\mathbf{\Pi}'(\lambda_{i_1}\mathbf{E}_{i_1} - \lambda_{i_k}\mathbf{E}_{i_k})\tilde{\mathbf{\Pi}}\tilde{\mathbf{\Pi}}'\right)\right)\right) = \\
&\varepsilon^2 \text{Tr}\left(\left(\mathbf{\Pi}\mathbf{S}_p\tilde{\mathbf{\Pi}}'\right) \circ \left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k})(\lambda_{i_1}\mathbf{E}_{i_1} - \lambda_{i_k}\mathbf{E}_{i_k})\right)\right) = \\
&\varepsilon^2 \text{Tr}\left(\left(\mathbf{\Pi}\mathbf{S}_p\tilde{\mathbf{\Pi}}'\right) \circ (\lambda_{i_1}\mathbf{E}_{i_1} + \lambda_{i_k}\mathbf{E}_{i_k})\right), \text{ and} \\
&\text{Tr}\left(\mathbf{V}'\mathbf{U}_{\mathbb{N}_p}\mathbf{\Pi}\mathbf{D}\left(\mathbf{S}_p \circ \left(\mathbf{D}^{-2}\mathbf{V}'\Sigma\mathbf{U}_{\mathbb{N}_p}\mathbf{\Pi}\mathbf{D}^{-1}\right)\right)\right) =
\end{aligned}$$

$$\begin{aligned}
& \text{Tr} \left(\varepsilon D \tilde{\Pi}' (U'_{i_1; i_1} - U'_{i_k; i_k}) U_{\mathbb{N}_p} \Pi (S_p \circ (D^{-1} V' \Sigma U_{\mathbb{N}_p} \Pi D^{-1})) \right) = \\
& \quad \varepsilon \text{Tr} \left(\tilde{\Pi}' (\mathbf{E}_{i_1} - \mathbf{E}_{i_k}) \Pi (S_p \circ (D^{-1} V' \Sigma U_{\mathbb{N}_p} \Pi)) \right) = \\
& \quad \varepsilon^2 \text{Tr} \left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k}) \Pi \left(S_p \circ \left(\tilde{\Pi}' (\lambda_{i_1} \mathbf{E}_{i_1} - \lambda_{i_k} \mathbf{E}_{i_k}) \Pi \right) \right) \tilde{\Pi}' \right) = \\
& \quad \varepsilon^2 \text{Tr} \left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k}) \left(\left(\Pi S_p \tilde{\Pi}' \right) \circ \left(\Pi \tilde{\Pi}' (\lambda_{i_1} \mathbf{E}_{i_1} - \lambda_{i_k} \mathbf{E}_{i_k}) \Pi \tilde{\Pi}' \right) \right) \right) = \\
& \quad \varepsilon^2 \text{Tr} \left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k}) \left(\left(\Pi S_p \tilde{\Pi}' \right) \circ \left(\Pi_{(i_1 i_k)} (\lambda_{i_1} \mathbf{E}_{i_1} - \lambda_{i_k} \mathbf{E}_{i_k}) \Pi_{(i_1 i_k)} \right) \right) \right) = \\
& \quad \varepsilon^2 \text{Tr} \left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k}) \left(\left(\Pi S_p \tilde{\Pi}' \right) \circ \left((\lambda_{i_1} \mathbf{E}_{i_k} - \lambda_{i_k} \mathbf{E}_{i_1}) \right) \right) \right) = \\
& \quad \varepsilon^2 \text{Tr} \left(\left(\Pi S_p \tilde{\Pi}' \right) \circ \left((\mathbf{E}_{i_1} - \mathbf{E}_{i_k}) (\lambda_{i_1} \mathbf{E}_{i_k} - \lambda_{i_k} \mathbf{E}_{i_1}) \right) \right) = \\
& \quad -\varepsilon^2 \text{Tr} \left(\left(\Pi S_p \tilde{\Pi}' \right) \circ (\lambda_{i_1} \mathbf{E}_{i_k} + \lambda_{i_k} \mathbf{E}_{i_1}) \right) = \\
& \quad -\varepsilon^2 \text{Tr} \left(\left(\Pi S_p \tilde{\Pi}' \right) \circ (\lambda_{i_1} \mathbf{E}_{i_k} + \lambda_{i_k} \mathbf{E}_{i_1}) \right).
\end{aligned}$$

Now, note that in both cases the matrices that are multiplied elementwise with $\Pi S_p \tilde{\Pi}'$ are diagonal and hence, we only need to look at diagonal elements of $\Pi S_p \tilde{\Pi}'$. Moreover,

$$\Pi S_p \tilde{\Pi}' = \Pi_{(i_1 i_k)} \Pi_{(i_k i_2)} \cdots \Pi_{(i_k i_{k-1})} \hat{\Pi} S_p \hat{\Pi}' \Pi_{(i_k i_{k-1})} \cdots \Pi_{(i_k i_2)},$$

where, $i_1 \cdots i_k$ are fixed points of permutation corresponding to $\hat{\Pi}$ so $\hat{\Pi} S_p \hat{\Pi}'$ has the same values at diagonal positions i_1 and i_k as the original matrix S_p . The only permutation that is only on the left side is $\Pi_{(i_1 i_k)}$ which exchanges the i_1 and i_k rows of S_p . Since S_p is such that the elements at each row before the diagonal element are the same and $i_k > i_1$, we have the i_1 and i_k diagonal elements of $\Pi S_p \tilde{\Pi}'$ have the same value. Let that value be denoted as s . Then the sum of the above two equations yields $m(\lambda_{i_1} + \lambda_{i_k}) - m(\lambda_{i_1} + \lambda_{i_k}) = 0$, as claimed.

B DERIVATIVES OF THE LOSS FUNCTION

B.1 FIRST AND SECOND ORDER FRÉCHET DERIVATIVE

In order to derive and analyze the critical points of the cost function which is a real-valued function of matrices we use the first and second order Fréchet derivatives as described in chapter 4 of Zeidler (1995). For a function $f : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ the first order Fréchet derivative at the point $\mathbf{A} \in \mathbb{R}^{n \times m}$ is a linear functional $df(\mathbf{A}) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ such that

$$\lim_{\mathbf{V} \rightarrow 0} \frac{|f(\mathbf{A} + \mathbf{V}) - f(\mathbf{A}) - df(\mathbf{A})\mathbf{V}|}{\|\mathbf{V}\|_F} = 0,$$

where we used the shorthand $df(\mathbf{A})\mathbf{V} \equiv (df(\mathbf{A}))(\mathbf{V})$. Similarly, the 2nd derivative is a bilinear functional $d^2f(\mathbf{A}) : \mathbb{R}^{n \times m} \times \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ such that

$$\lim_{\mathbf{V} \rightarrow 0} \frac{|df(\mathbf{A} + \mathbf{V})\mathbf{K} - df(\mathbf{A})\mathbf{K} - d^2f(\mathbf{A})\mathbf{V}\mathbf{K}|}{\|\mathbf{V}\|_F} = 0,$$

for all $\|\mathbf{K}\|_F \leq 1$, where again $d^2f(\mathbf{A})\mathbf{V}\mathbf{K} \equiv (d^2f(\mathbf{A}))(\mathbf{V}, \mathbf{K})$. The generalized Taylor formula then becomes:

$$f(\mathbf{A} + \mathbf{V}) = f(\mathbf{A}) + df(\mathbf{A})\mathbf{V} + \frac{1}{2}d^2f(\mathbf{A})\mathbf{V}^2 + o(\|\mathbf{V}\|^2),$$

Moreover, we derive functions $\nabla f : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{n \times m}$ and $\mathbf{H}(\mathbf{A}) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{n \times m}$ such that $df(\mathbf{A})\mathbf{V} = \langle \nabla f(\mathbf{A}), \mathbf{V} \rangle_F$ and $d^2f(\mathbf{A})\mathbf{V}^2 = \langle \mathbf{H}(\mathbf{A})\mathbf{V}, \mathbf{V} \rangle_F$, where again $\mathbf{H}(\mathbf{A})\mathbf{V} \equiv \mathbf{H}(\mathbf{A})(\mathbf{V})$. Then clearly, $\mathbf{A} \in \mathbb{R}^{n \times m}$ is a critical point of f iff $\nabla f(\mathbf{A}) = 0$ and for such \mathbf{A} s the sign of the bilinear form $\langle \mathbf{H}(\mathbf{A})\mathbf{V}, \mathbf{V} \rangle$ over directions \mathbf{V} determines the type of the critical point.

Extending the generalized Taylor theorem of Zeidler (1995), the second order Taylor expansion for the loss $L(\mathbf{A}, \mathbf{B})$ is then given by

$$\begin{aligned} L(\mathbf{A} + \mathbf{V}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= d_{\mathbf{A}}L(\mathbf{A}, \mathbf{B})\mathbf{V} + d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} + \frac{1}{2}d_{\mathbf{A}}^2L(\mathbf{A}, \mathbf{B})\mathbf{V}^2 \\ &\quad + d_{\mathbf{A}\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{V}\mathbf{W} + \frac{1}{2}d_{\mathbf{B}}^2L(\mathbf{A}, \mathbf{B})\mathbf{W}^2 + R_{\mathbf{V}, \mathbf{W}}(\mathbf{A}, \mathbf{B}), \end{aligned} \quad (63)$$

where, if $\|\mathbf{V}\|_F, \|\mathbf{W}\|_F = O(\varepsilon)$ then $\|R(\mathbf{V}, \mathbf{W})\| = O(\varepsilon^3)$. Clearly, as at critical points where $d_{\mathbf{A}}L(\mathbf{A}, \mathbf{B})\mathbf{V} + d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} = 0$, as $\varepsilon \rightarrow 0$ we have $R_{\mathbf{V}, \mathbf{W}}(\mathbf{A}, \mathbf{B}) \rightarrow 0$ and the sign of the sum of the second order partial Fréchet derivatives determines the type of the critical point very much similar to second partial derivative test for two variable functions. However, here for local minima we have to show the sign is positive in all directions and for saddle points have to show the sign is positive in some directions and negative at least in on direction. Finally, note that the smoothness of the loss entails that Fréchet derivative and directional derivative (Gateaux) both exist and (foregoing some subtleties in definition) are the same.

B.2 FIRST AND SECOND ORDER DERIVATIVE OF THE LOSS WRT TO \mathbf{B}

Lemma 5. *The first and second (partial Fréchet) derivative of the loss $L(\mathbf{A}, \mathbf{B})$ wrt to \mathbf{B} is derived as follows.*

$$d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} = -2 \text{Tr}(\mathbf{W}'(\mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx})) \quad (64)$$

$$= -2 \langle \mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx}, \mathbf{W} \rangle_F. \quad (65)$$

$$d_{\mathbf{B}}^2L(\mathbf{A}, \mathbf{B})\mathbf{W}^2 = 2 \langle (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{W} \Sigma_{xx}, \mathbf{W} \rangle_F = 2 \text{Tr}(\mathbf{W}'(\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{W} \Sigma_{xx}). \quad (66)$$

Proof. Directly compute

$$L(\mathbf{A}, \mathbf{B} + \mathbf{W}) = \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A} \mathbf{I}_{i:p}(\mathbf{B} + \mathbf{W}) \mathbf{X}\|_F^2$$

$$\begin{aligned}
&= \sum_{i=1}^p \langle Y - \mathbf{A} \mathbf{I}_{i;p} (\mathbf{B} + \mathbf{W}) \mathbf{X}, Y - \mathbf{A} \mathbf{I}_{i;p} (\mathbf{B} + \mathbf{W}) \mathbf{X} \rangle_F \\
&= \sum_{i=1}^p \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F + \sum_{i=1}^p \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, -\mathbf{A} \mathbf{I}_{i;p} \mathbf{W} \mathbf{X} \rangle_F \\
&+ \sum_{i=1}^p \langle -\mathbf{A} \mathbf{I}_{i;p} \mathbf{W} \mathbf{X}, Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F + \sum_{i=1}^p \langle -\mathbf{A} \mathbf{I}_{i;p} \mathbf{W} \mathbf{X}, -\mathbf{A} \mathbf{I}_{i;p} \mathbf{W} \mathbf{X} \rangle_F \\
&= L(\mathbf{A}, \mathbf{B}) - \sum_{i=1}^p 2 \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{A} \mathbf{I}_{i;p} \mathbf{W} \mathbf{X} \rangle + O(\|\mathbf{W}\|_F^2) \implies
\end{aligned}$$

$$\begin{aligned}
L(\mathbf{A}, \mathbf{B} + \mathbf{W}) - L(\mathbf{A}, \mathbf{B}) &= -2 \sum_{i=1}^p \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{A} \mathbf{I}_{i;p} \mathbf{W} \mathbf{X} \rangle_F + O(\|\mathbf{W}\|_F^2) \xrightarrow{\mathbf{W} \rightarrow 0} \\
d_{\mathbf{B}} L(\mathbf{A}, \mathbf{B}) \mathbf{W} &= -2 \sum_{i=1}^p \text{Tr}(\mathbf{X}' \mathbf{W}' \mathbf{I}_{i;p} \mathbf{A}' (Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X})) \\
&= -2 \text{Tr} \left(\mathbf{W}' \left(\left(\sum_{i=1}^p \mathbf{I}_{i;p} \right) \mathbf{A}' \mathbf{Y} \mathbf{X}' - \left(\sum_{i=1}^p \mathbf{I}_{i;p} \mathbf{A}' \mathbf{A} \mathbf{I}_{i;p} \right) \mathbf{B} \mathbf{X} \mathbf{X}' \right) \right) \\
&= -2 \text{Tr} (\mathbf{W}' (\mathbf{T}_p \mathbf{A}' \mathbf{Y} \mathbf{X}' - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \mathbf{X} \mathbf{X}')),
\end{aligned}$$

which can be written as the given form. For the second derivative wrt \mathbf{B} we have

$$\begin{aligned}
d_{\mathbf{B}} L(\mathbf{A}, \mathbf{B}) \mathbf{W} &= -2 \langle \mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx}, \mathbf{W} \rangle_F \implies \\
d_{\mathbf{B}} L(\mathbf{A}, \mathbf{B} + \bar{\mathbf{W}}) \mathbf{W} &= -2 \langle \mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) (\mathbf{B} + \bar{\mathbf{W}}) \Sigma_{xx}, \mathbf{W} \rangle_F \\
&= -2 \langle \mathbf{T}_p \mathbf{A}' \Sigma_{yx} - (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \mathbf{B} \Sigma_{xx}, \mathbf{W} \rangle_F \\
&+ 2 \langle (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \bar{\mathbf{W}} \Sigma_{xx}, \mathbf{W} \rangle_F \implies \\
d_{\mathbf{B}} L(\mathbf{A}, \mathbf{B} + \bar{\mathbf{W}}) \mathbf{W} - d_{\mathbf{B}} L(\mathbf{A}, \mathbf{B}) \mathbf{W} &= 2 \langle (\mathbf{S}_p \circ (\mathbf{A}' \mathbf{A})) \bar{\mathbf{W}} \Sigma_{xx}, \mathbf{W} \rangle_F,
\end{aligned}$$

which by having $\bar{\mathbf{W}} \rightarrow 0$ results in the second order partial derivative. \square

B.3 FIRST AND SECOND ORDER DERIVATIVE OF THE LOSS WRT TO \mathbf{A}

Lemma 6. *The first and second (partial Fréchet) derivative of the loss $L(\mathbf{A}, \mathbf{B})$ wrt to \mathbf{A} is derived as follows.*

$$d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} = -2 \langle \Sigma_{yx} \mathbf{B}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F, \quad (67)$$

$$d_{\mathbf{A}\mathbf{B}}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{V} \mathbf{W} = -2 \langle \Sigma_{yx} \mathbf{W}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{W}')) - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{W} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F, \quad (68)$$

$$d_{\mathbf{A}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{V}^2 = 2 \langle \mathbf{V} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F. \quad (69)$$

Proof. Directly compute

$$\begin{aligned}
L(\mathbf{A} + \mathbf{V}, \mathbf{B}) &= \sum_{i=1}^p \langle Y - (\mathbf{A} + \mathbf{V}) \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, Y - (\mathbf{A} + \mathbf{V}) \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F \\
&= \sum_{i=1}^p \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F - \sum_{i=1}^p \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F \\
&+ \sum_{i=1}^p \langle -\mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F + \sum_{i=1}^p \langle -\mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, -\mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F \\
&= L(\mathbf{A}, \mathbf{B}) - \sum_{i=1}^p 2 \langle Y - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F + \sum_{i=1}^p \langle \mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F
\end{aligned}$$

$$\begin{aligned}
L(\mathbf{A} + \mathbf{V}, \mathbf{B}) - L(\mathbf{A}, \mathbf{B}) &= - \sum_{i=1}^p 2 \langle \mathbf{Y} - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F + O(\|\mathbf{V}\|_F^2) \xrightarrow{\mathbf{V} \rightarrow 0} \\
d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} &= - \sum_{i=1}^p 2 \langle \mathbf{Y} - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}, \mathbf{V} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X} \rangle_F \\
&= -2 \text{Tr}(\mathbf{V}' (\sum_{yx} \mathbf{B}' \sum_{i=1}^p \mathbf{I}_{i;p} - \mathbf{A} \sum_{i=1}^p \mathbf{I}_{i;p} \mathbf{B} \Sigma_{xx} \mathbf{B}' \mathbf{I}_{i;p})) \implies \\
d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} &= -2 \langle \sum_{yx} \mathbf{B}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \implies \\
d_{\mathbf{A}} L(\mathbf{A} + \bar{\mathbf{V}}, \mathbf{B}) \mathbf{V} &= -2 \langle \sum_{yx} \mathbf{B}' \mathbf{T}_p - (\mathbf{A} + \bar{\mathbf{V}}) (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \\
d_{\mathbf{A}} L(\mathbf{A} + \bar{\mathbf{V}}, \mathbf{B}) \mathbf{V} - d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} &= 2 \langle \bar{\mathbf{V}} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \xrightarrow{\bar{\mathbf{V}} \rightarrow 0} \\
d_{\mathbf{A}^2}^2 L(\mathbf{A}, \mathbf{B})(\mathbf{V}, \bar{\mathbf{V}}) &= 2 \langle \bar{\mathbf{V}} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \implies \\
d_{\mathbf{A}^2}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{V}^2 &= 2 \langle \mathbf{V} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \\
d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B} + \mathbf{W}) \mathbf{V} &= -2 \langle \sum_{yx} (\mathbf{B} + \mathbf{W})' \mathbf{T}_p, \mathbf{V} \rangle_F \\
&\quad -2 \langle -\mathbf{A} (\mathbf{S}_p \circ ((\mathbf{B} + \mathbf{W}) \Sigma_{xx} (\mathbf{B} + \mathbf{W})')), \mathbf{V} \rangle_F \\
&\quad -2 \langle \sum_{yx} \mathbf{B}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \\
&= d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} - 2 \langle \sum_{yx} \mathbf{W}' \mathbf{T}_p, \mathbf{V} \rangle_F \\
&\quad -2 \langle -\mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{W}')) - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{W} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F + O(\|\mathbf{W}\|_F^2) \implies \\
d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B} + \mathbf{W}) \mathbf{V} - d_{\mathbf{A}} L(\mathbf{A}, \mathbf{B}) \mathbf{V} &= -2 \langle \sum_{yx} \mathbf{W}' \mathbf{T}_p, \mathbf{V} \rangle_F \\
&\quad -2 \langle -\mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{W}')) - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{W} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F \\
&\quad + O(\|\mathbf{W}\|_F^2) \xrightarrow{\mathbf{W} \rightarrow 0} \\
d_{\mathbf{A} \mathbf{B}}^2 L(\mathbf{A}, \mathbf{B}) \mathbf{V} \mathbf{W} &= -2 \langle \sum_{yx} \mathbf{W}' \mathbf{T}_p - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{B} \Sigma_{xx} \mathbf{W}')) - \mathbf{A} (\mathbf{S}_p \circ (\mathbf{W} \Sigma_{xx} \mathbf{B}')), \mathbf{V} \rangle_F.
\end{aligned}$$

□