

Eliminating the Invariance on the Loss Landscape of Linear Autoencoders

Reza Oftadeh, Jiayi Shen, Atlas Wang, Dylan Shell

Texas A&M University
Department of Computer Science and Engineering

ICML 2020

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Overview

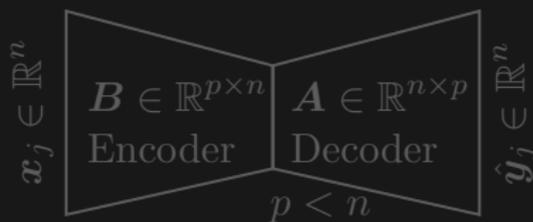
- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - **Invariant local minima become saddle points.**
 - Computational complexity is of the same order of MSE loss.

Overview

- ▶ Linear Autoencoder (LAE) with Mean Square Error (MSE).
The classical results:
 - Loss surface has been analytically characterized.
 - All local minima are global minima.
 - The columns of the optimal decoder does not identify the principal directions but only their low dimensional subspace (the so-called *invariance* problem).
- ▶ We present a new loss function for LAE:
 - Analytically characterize the loss landscape.
 - All local minima are global minima.
 - The columns of the optimal decoder span the *principal directions*.
 - Invariant local minima become saddle points.
 - Computational complexity is of the same order of MSE loss.

Setup

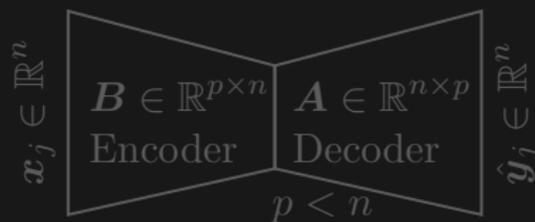
- ▶ **Data:** m sample points of dimension n :
 - Input: $\mathbf{x}_j \in \mathbb{R}^n$, Output: $\mathbf{y}_j \in \mathbb{R}^n$ for $j = 1, \dots, m$.
 - In matrix form: $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$.
- ▶ **LAE:** A neural network with linear activation functions and single hidden layer of width $p < n$.



- The weights: The encoder matrix \mathbf{B} , and the decoder matrix \mathbf{A} .
- The global map is $\hat{\mathbf{y}}_j = \mathbf{A}\mathbf{B}\mathbf{x}_j$ or $\hat{\mathbf{Y}} = \mathbf{A}\mathbf{B}\mathbf{X}$.

Setup

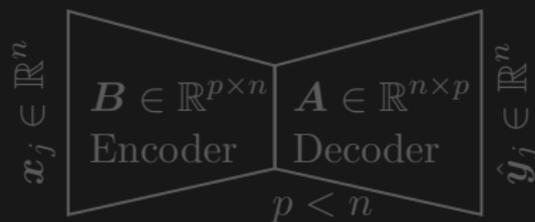
- ▶ Data: m sample points of dimension n :
 - Input: $\mathbf{x}_j \in \mathbb{R}^n$, Output: $\mathbf{y}_j \in \mathbb{R}^n$ for $j = 1, \dots, m$.
 - In matrix form: $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$.
- ▶ LAE: A neural network with linear activation functions and single hidden layer of width $p < n$.



- The weights: The encoder matrix \mathbf{B} , and the decoder matrix \mathbf{A} .
- The global map is $\hat{\mathbf{y}}_j = \mathbf{A}\mathbf{B}\mathbf{x}_j$ or $\hat{\mathbf{Y}} = \mathbf{A}\mathbf{B}\mathbf{X}$.

Setup

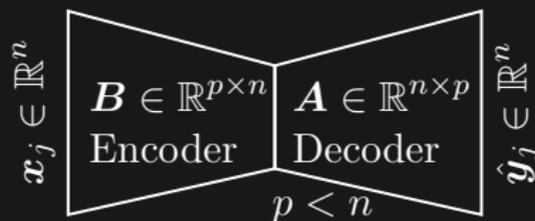
- ▶ Data: m sample points of dimension n :
 - Input: $\mathbf{x}_j \in \mathbb{R}^n$, Output: $\mathbf{y}_j \in \mathbb{R}^n$ for $j = 1, \dots, m$.
 - In matrix form: $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$.
- ▶ LAE: A neural network with linear activation functions and single hidden layer of width $p < n$.



- The weights: The encoder matrix \mathbf{B} , and the decoder matrix \mathbf{A} .
- The global map is $\hat{\mathbf{y}}_j = \mathbf{A}\mathbf{B}\mathbf{x}_j$ or $\hat{\mathbf{Y}} = \mathbf{A}\mathbf{B}\mathbf{X}$.

Setup

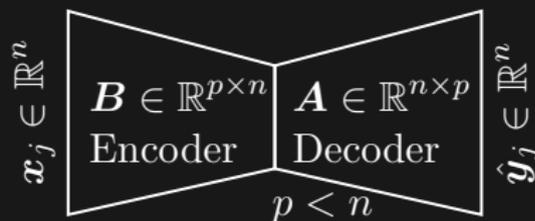
- ▶ Data: m sample points of dimension n :
 - Input: $\mathbf{x}_j \in \mathbb{R}^n$, Output: $\mathbf{y}_j \in \mathbb{R}^n$ for $j = 1, \dots, m$.
 - In matrix form: $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$.
- ▶ LAE: A neural network with linear activation functions and single hidden layer of width $p < n$.



- The weights: The encoder matrix \mathbf{B} , and the decoder matrix \mathbf{A} .
- The global map is $\hat{\mathbf{y}}_j = \mathbf{A}\mathbf{B}\mathbf{x}_j$ or $\hat{\mathbf{Y}} = \mathbf{A}\mathbf{B}\mathbf{X}$.

Setup

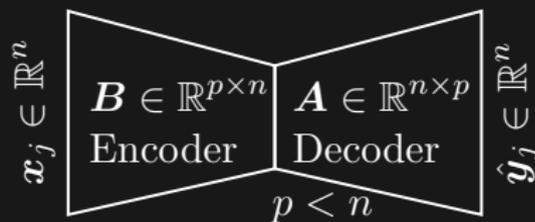
- ▶ Data: m sample points of dimension n :
 - Input: $\mathbf{x}_j \in \mathbb{R}^n$, Output: $\mathbf{y}_j \in \mathbb{R}^n$ for $j = 1, \dots, m$.
 - In matrix form: $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$.
- ▶ LAE: A neural network with linear activation functions and single hidden layer of width $p < n$.



- The weights: The encoder matrix \mathbf{B} , and the decoder matrix \mathbf{A} .
- The global map is $\hat{\mathbf{y}}_j = \mathbf{A}\mathbf{B}\mathbf{x}_j$ or $\hat{\mathbf{Y}} = \mathbf{A}\mathbf{B}\mathbf{X}$.

Setup

- ▶ Data: m sample points of dimension n :
 - Input: $\mathbf{x}_j \in \mathbb{R}^n$, Output: $\mathbf{y}_j \in \mathbb{R}^n$ for $j = 1, \dots, m$.
 - In matrix form: $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times m}$.
- ▶ LAE: A neural network with linear activation functions and single hidden layer of width $p < n$.



- The weights: The encoder matrix \mathbf{B} , and the decoder matrix \mathbf{A} .
- The global map is $\hat{\mathbf{y}}_j = \mathbf{A}\mathbf{B}\mathbf{x}_j$ or $\hat{\mathbf{Y}} = \mathbf{A}\mathbf{B}\mathbf{X}$.

The loss functions

- ▶ The MSE loss: $\tilde{L}(\mathbf{A}, \mathbf{B}) := \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$.
 - If $(\mathbf{A}^*, \mathbf{B}^*)$ is a local minimum of \tilde{L} then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*)$ is another local minima:

$$\tilde{L}(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*) = \|\mathbf{Y} - \mathbf{A}^* \mathbf{C} \mathbf{C}^{-1} \mathbf{B}^* \mathbf{X}\|_F^2 = \tilde{L}(\mathbf{A}^*, \mathbf{B}^*).$$

- ▶ The proposed loss: $L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}\|_F^2$,
where, $\mathbf{I}_{i;p} = \text{diag}(\underbrace{1, \dots, 1}_i, 0, \dots, 0) \in \mathbb{R}^{p \times p}$.

The loss functions

- ▶ The MSE loss: $\tilde{L}(\mathbf{A}, \mathbf{B}) := \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$.
 - If $(\mathbf{A}^*, \mathbf{B}^*)$ is a local minimum of \tilde{L} then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*)$ is another local minima:

$$\tilde{L}(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*) = \|\mathbf{Y} - \mathbf{A}^* \mathbf{C} \mathbf{C}^{-1} \mathbf{B}^* \mathbf{X}\|_F^2 = \tilde{L}(\mathbf{A}^*, \mathbf{B}^*).$$

- ▶ The proposed loss: $L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}\|_F^2$,
where, $\mathbf{I}_{i;p} = \text{diag}(\underbrace{1, \dots, 1}_i, 0, \dots, 0) \in \mathbb{R}^{p \times p}$.

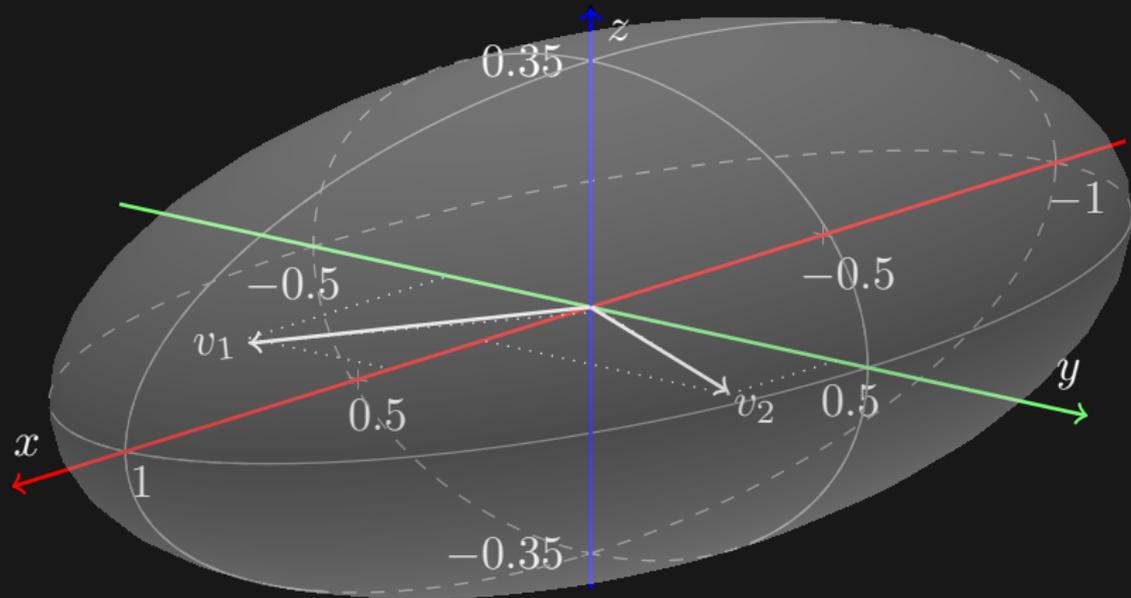
The loss functions

- ▶ The MSE loss: $\tilde{L}(\mathbf{A}, \mathbf{B}) := \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$.
 - If $(\mathbf{A}^*, \mathbf{B}^*)$ is a local minimum of \tilde{L} then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*)$ is another local minima:

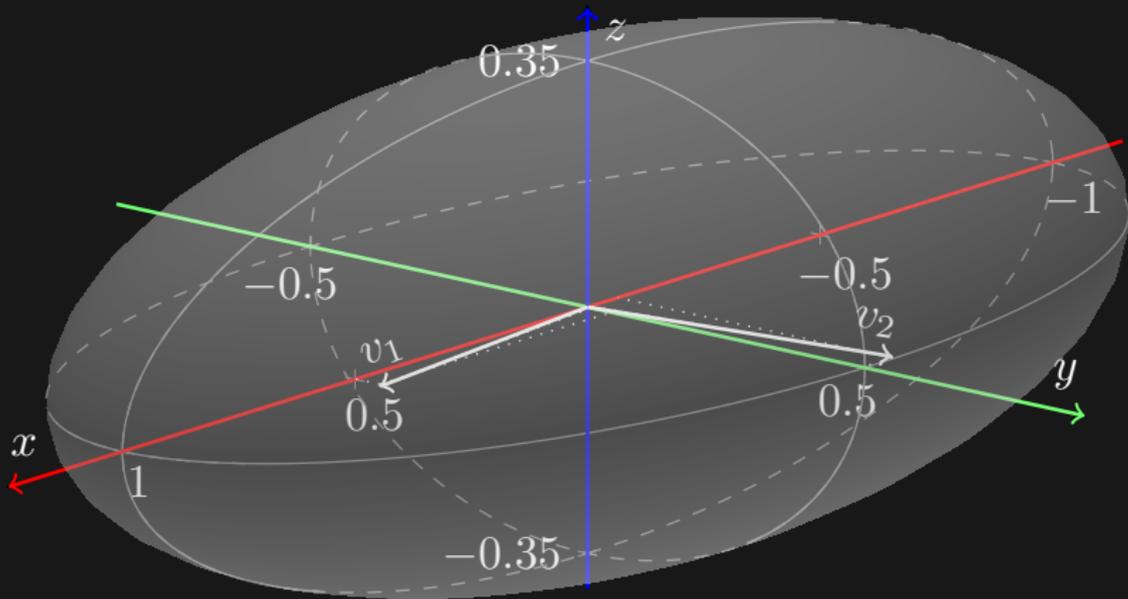
$$\tilde{L}(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*) = \|\mathbf{Y} - \mathbf{A}^* \mathbf{C} \mathbf{C}^{-1} \mathbf{B}^* \mathbf{X}\|_F^2 = \tilde{L}(\mathbf{A}^*, \mathbf{B}^*).$$

- ▶ The proposed loss: $L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A} \mathbf{I}_{i;p} \mathbf{B} \mathbf{X}\|_F^2$,
where, $\mathbf{I}_{i;p} = \text{diag}(\underbrace{1, \dots, 1}_i, 0, \dots, 0) \in \mathbb{R}^{p \times p}$.

A Visualization: MSE Loss



A Visualization: Proposed Loss



The loss functions

► The MSE loss: $\tilde{L}(\mathbf{A}, \mathbf{B}) := \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$.

- If $(\mathbf{A}^*, \mathbf{B}^*)$ is a local minimum of \tilde{L} then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^*\mathbf{C}, \mathbf{C}^{-1}\mathbf{B}^*)$ is another local minima:

$$\tilde{L}(\mathbf{A}^*\mathbf{C}, \mathbf{C}^{-1}\mathbf{B}^*) = \|\mathbf{Y} - \mathbf{A}^*\mathbf{C}\mathbf{C}^{-1}\mathbf{B}^*\mathbf{X}\|_F^2 = \tilde{L}(\mathbf{A}^*, \mathbf{B}^*).$$

► The proposed loss: $L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}\mathbf{X}\|_F^2$,
where, $\mathbf{I}_{i;p} = \text{diag}(\underbrace{1, \dots, 1}_i, 0, \dots, 0) \in \mathbb{R}^{p \times p}$.

- Intuition: (Sequential) As an example look at $p = 3$, where

$$\mathbf{I}_{1;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{I}_{2;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{I}_{3;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

- But does this work simultaneously? And is it computationally feasible (p can be large)?
- Well, *it does and it is!* But before getting into details let's discuss some implications ...

The loss functions

► The MSE loss: $\tilde{L}(\mathbf{A}, \mathbf{B}) := \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$.

- If $(\mathbf{A}^*, \mathbf{B}^*)$ is a local minimum of \tilde{L} then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^*\mathbf{C}, \mathbf{C}^{-1}\mathbf{B}^*)$ is another local minima:

$$\tilde{L}(\mathbf{A}^*\mathbf{C}, \mathbf{C}^{-1}\mathbf{B}^*) = \|\mathbf{Y} - \mathbf{A}^*\mathbf{C}\mathbf{C}^{-1}\mathbf{B}^*\mathbf{X}\|_F^2 = \tilde{L}(\mathbf{A}^*, \mathbf{B}^*).$$

► The proposed loss: $L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}\mathbf{X}\|_F^2$,
where, $\mathbf{I}_{i;p} = \text{diag}(\underbrace{1, \dots, 1}_i, 0, \dots, 0) \in \mathbb{R}^{p \times p}$.

- Intuition: (Sequential) As an example look at $p = 3$, where

$$\mathbf{I}_{1;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{I}_{2;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{I}_{3;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

- But does this work simultaneously? And is it computationally feasible (p can be large)?
- Well, *it does and it is!* But before getting into details let's discuss some implications ...

The loss functions

► The MSE loss: $\tilde{L}(\mathbf{A}, \mathbf{B}) := \|\mathbf{Y} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$.

- If $(\mathbf{A}^*, \mathbf{B}^*)$ is a local minimum of \tilde{L} then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^*\mathbf{C}, \mathbf{C}^{-1}\mathbf{B}^*)$ is another local minima:

$$\tilde{L}(\mathbf{A}^*\mathbf{C}, \mathbf{C}^{-1}\mathbf{B}^*) = \|\mathbf{Y} - \mathbf{A}^*\mathbf{C}\mathbf{C}^{-1}\mathbf{B}^*\mathbf{X}\|_F^2 = \tilde{L}(\mathbf{A}^*, \mathbf{B}^*).$$

► The proposed loss: $L(\mathbf{A}, \mathbf{B}) := \sum_{i=1}^p \|\mathbf{Y} - \mathbf{A}\mathbf{I}_{i;p}\mathbf{B}\mathbf{X}\|_F^2$,
where, $\mathbf{I}_{i;p} = \text{diag}(\underbrace{1, \dots, 1}_i, 0, \dots, 0) \in \mathbb{R}^{p \times p}$.

- Intuition: (Sequential) As an example look at $p = 3$, where

$$\mathbf{I}_{1;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{I}_{2;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{I}_{3;3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

- But does this work simultaneously? And is it computationally feasible (p can be large)?
- Well, *it does and it is!* But before getting into details let's discuss some implications ...

Implications

- ▶ Let $(\mathbf{A}^*, \mathbf{B}^*)$ be the local minimum of MSE loss, where the columns of \mathbf{A}^* are the largest eigenvectors of the sample covariance matrix, then for any invertible $\mathbf{C} \in \mathbb{R}^{p \times p}$, $(\mathbf{A}^* \mathbf{C}, \mathbf{C}^{-1} \mathbf{B}^*)$ is another local minima.
 - Numerically, on the same dataset different runs with different initializations lead to different optimal points.
 - Almost surely none will represent the principal directions.
- ▶ *The only local minimum of the loss L is $(\mathbf{A}^*, \mathbf{B}^*)$, up to the normalization of the columns.*
 - The loss L enables low rank decomposition as a single optimization block that can be incorporated as part of a larger pipeline.
 - Potentially enabling LAEs to compete with other approaches for low rank decomposition.

Critical Points

- The critical point equations of \tilde{L} and L .

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

where,

- \mathbf{A}' is the transpose of \mathbf{A} .
- $\boldsymbol{\Sigma}_{xx} = \mathbf{X}\mathbf{X}'$, $\boldsymbol{\Sigma}_{yx} = \mathbf{Y}\mathbf{X}'$ are covariance matrices.
- \circ is the (element-wise) Hadamard product, and
- \mathbf{T}_p , and \mathbf{S}_p are

$$\mathbf{T}_p = \text{diag}(p, p-1, \dots, 1),$$

$$\mathbf{S}_p = \begin{bmatrix} p & p-1 & \dots & 1 \\ p-1 & p-1 & \dots & 1 \\ \vdots & \vdots & \ddots & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } \mathbf{S}_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

Critical Points

- The critical point equations of \tilde{L} and L .

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

where,

- \mathbf{A}' is the transpose of \mathbf{A} .
- $\boldsymbol{\Sigma}_{xx} = \mathbf{X}\mathbf{X}'$, $\boldsymbol{\Sigma}_{yx} = \mathbf{Y}\mathbf{X}'$ are covariance matrices.
- \circ is the (element-wise) Hadamard product, and
- \mathbf{T}_p , and \mathbf{S}_p are

$$\mathbf{T}_p = \text{diag}(p, p-1, \dots, 1),$$

$$\mathbf{S}_p = \begin{bmatrix} p & p-1 & \cdots & 1 \\ p-1 & p-1 & \cdots & 1 \\ \vdots & \vdots & \ddots & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } \mathbf{S}_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

Critical Points

- ▶ The critical point equations of \tilde{L} and L .

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

where,

- \mathbf{A}' is the transpose of \mathbf{A} .
- $\boldsymbol{\Sigma}_{xx} = \mathbf{X}\mathbf{X}'$, $\boldsymbol{\Sigma}_{yx} = \mathbf{Y}\mathbf{X}'$ are covariance matrices.
- \circ is the (element-wise) Hadamard product, and
- \mathbf{T}_p , and \mathbf{S}_p are

$$\mathbf{T}_p = \text{diag}(p, p-1, \dots, 1),$$

$$\mathbf{S}_p = \begin{bmatrix} p & p-1 & \dots & 1 \\ p-1 & p-1 & \dots & 1 \\ \vdots & \vdots & \ddots & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } \mathbf{S}_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

Critical Points

- The critical point equations of \tilde{L} and L .

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

where,

- \mathbf{A}' is the transpose of \mathbf{A} .
- $\boldsymbol{\Sigma}_{xx} = \mathbf{X}\mathbf{X}'$, $\boldsymbol{\Sigma}_{yx} = \mathbf{Y}\mathbf{X}'$ are covariance matrices.
- \circ is the (element-wise) Hadamard product, and
- \mathbf{T}_p , and \mathbf{S}_p are

$$\mathbf{T}_p = \text{diag}(p, p-1, \dots, 1),$$

$$\mathbf{S}_p = \begin{bmatrix} p & p-1 & \dots & 1 \\ p-1 & p-1 & \dots & 1 \\ \vdots & \vdots & \ddots & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } \mathbf{S}_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

Critical Points

- The critical point equations of \tilde{L} and L .

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

where,

- \mathbf{A}' is the transpose of \mathbf{A} .
- $\boldsymbol{\Sigma}_{xx} = \mathbf{X}\mathbf{X}'$, $\boldsymbol{\Sigma}_{yx} = \mathbf{Y}\mathbf{X}'$ are covariance matrices.
- \circ is the (element-wise) Hadamard product, and
- \mathbf{T}_p , and \mathbf{S}_p are

$$\mathbf{T}_p = \text{diag}(p, p-1, \dots, 1),$$

$$\mathbf{S}_p = \begin{bmatrix} p & p-1 & \dots & 1 \\ p-1 & p-1 & \dots & 1 \\ \vdots & \vdots & \ddots & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \text{ e.g. } \mathbf{S}_4 = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 3 & 2 & 1 \\ 2 & 2 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

Results I

- ▶ Every critical point of $L(\mathbf{A}, \mathbf{B})$ is a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B})$, but not the other way around.
- ▶ Local minima of L , and \tilde{L} :

$$\begin{array}{l} \text{For } \tilde{L}(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{C}_p, \\ \mathbf{B}^* = \mathbf{C}_p^{-1} \mathbf{U}'_{1:p} \Sigma_{yx} \Sigma_{xx}^{-1}, \end{array} \left| \begin{array}{l} \text{For } L(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{D}_p, \\ \mathbf{B}^* = \mathbf{D}_p^{-1} \mathbf{U}'_{1:p} \Sigma_{yx} \Sigma_{xx}^{-1}, \end{array} \right.$$

- The i^{th} column of $\mathbf{U}_{1:p}$ is a unit eigenvector of $\Sigma := \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$ corresponding the i^{th} largest eigenvalue.
 - \mathbf{D}_p is a diagonal matrix with nonzero diagonal elements, and $\mathbf{C}_p \in \text{GL}_p(\mathbb{R})$.
- ▶ The characterization of the loss landscape:
 - The structure of full rank saddle points.
 - The structure of low rank saddle points (rather involved!).

Results I

- ▶ Every critical point of $L(\mathbf{A}, \mathbf{B})$ is a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B})$, but not the other way around.
- ▶ Local minima of L , and \tilde{L} :

$$\begin{array}{l} \text{For } \tilde{L}(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{C}_p, \\ \mathbf{B}^* = \mathbf{C}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array} \left| \begin{array}{l} \text{For } L(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{D}_p, \\ \mathbf{B}^* = \mathbf{D}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array} \right.$$

- 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.
 - \mathbf{D}_p is a diagonal matrix with nonzero diagonal elements, and $\mathbf{C}_p \in \text{GL}_p(\mathbb{R})$.
- ▶ The characterization of the loss landscape:
 - The structure of full rank saddle points.
 - The structure of low rank saddle points (rather involved!).

Results I

- ▶ Every critical point of $L(\mathbf{A}, \mathbf{B})$ is a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B})$, but not the other way around.
- ▶ Local minima of L , and \tilde{L} :

$$\begin{array}{l|l} \text{For } \tilde{L}(\mathbf{A}, \mathbf{B}): & \text{For } L(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{C}_p, & \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{D}_p, \\ \mathbf{B}^* = \mathbf{C}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, & \mathbf{B}^* = \mathbf{D}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array}$$

- 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.
 - \mathbf{D}_p is a diagonal matrix with nonzero diagonal elements, and $\mathbf{C}_p \in \text{GL}_p(\mathbb{R})$.
- ▶ The characterization of the loss landscape:
 - The structure of full rank saddle points.
 - The structure of low rank saddle points (rather involved!).

Results I

- ▶ Every critical point of $L(\mathbf{A}, \mathbf{B})$ is a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B})$, but not the other way around.
- ▶ Local minima of L , and \tilde{L} :

$$\begin{array}{l} \text{For } \tilde{L}(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{C}_p, \\ \mathbf{B}^* = \mathbf{C}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array} \left| \begin{array}{l} \text{For } L(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{D}_p, \\ \mathbf{B}^* = \mathbf{D}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array} \right.$$

- 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.
 - \mathbf{D}_p is a diagonal matrix with nonzero diagonal elements, and $\mathbf{C}_p \in \text{GL}_p(\mathbb{R})$.
- ▶ The characterization of the loss landscape:
 - The structure of full rank saddle points.
 - The structure of low rank saddle points (rather involved!).

Results I

- ▶ Every critical point of $L(\mathbf{A}, \mathbf{B})$ is a critical point of $\tilde{L}(\mathbf{A}, \mathbf{B})$, but not the other way around.
- ▶ Local minima of L , and \tilde{L} :

$$\begin{array}{l} \text{For } \tilde{L}(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{C}_p, \\ \mathbf{B}^* = \mathbf{C}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array} \left| \begin{array}{l} \text{For } L(\mathbf{A}, \mathbf{B}): \\ \mathbf{A}^* = \mathbf{U}_{1:p} \mathbf{D}_p, \\ \mathbf{B}^* = \mathbf{D}_p^{-1} \mathbf{U}'_{1:p} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \end{array} \right.$$

- 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.
 - \mathbf{D}_p is a diagonal matrix with nonzero diagonal elements, and $\mathbf{C}_p \in \text{GL}_p(\mathbb{R})$.
- ▶ The characterization of the loss landscape:
 - The structure of full rank saddle points.
 - The structure of low rank saddle points (rather involved!).

Results II

- ▶ The MSE loss \tilde{L} and our loss L can be written as

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

- ▶ The analytical gradients are:

$$\begin{aligned}d_B \tilde{L}(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2\langle \mathbf{A}'\boldsymbol{\Sigma}_{yx} - \mathbf{A}'\mathbf{A}\mathbf{B}\boldsymbol{\Sigma}_{xx}, \mathbf{W} \rangle_F, \\ d_B L(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2\langle \mathbf{T}_p\mathbf{A}'\boldsymbol{\Sigma}_{yx} - (\mathbf{S}_{p^\circ}(\mathbf{A}'\mathbf{A}))\mathbf{B}\boldsymbol{\Sigma}_{xx}, \mathbf{W} \rangle_F,\end{aligned}$$

in direction of $\mathbf{W} \in \mathbb{R}^{p \times n}$. The gradient for \mathbf{A} is similar.

- ▶ Finally, since the loss function is explicitly provided, any optimization method that works with MSE loss is usable with the proposed loss.

Results II

- ▶ The MSE loss \tilde{L} and our loss L can be written as

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

- ▶ The analytical gradients are:

$$\begin{aligned}d_{\mathbf{B}}\tilde{L}(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2\langle\mathbf{A}'\boldsymbol{\Sigma}_{yx} - \mathbf{A}'\mathbf{A}\mathbf{B}\boldsymbol{\Sigma}_{xx}, \mathbf{W}\rangle_F, \\ d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2\langle\mathbf{T}_p\mathbf{A}'\boldsymbol{\Sigma}_{yx} - (\mathbf{S}_{p\circ}(\mathbf{A}'\mathbf{A}))\mathbf{B}\boldsymbol{\Sigma}_{xx}, \mathbf{W}\rangle_F,\end{aligned}$$

in direction of $\mathbf{W} \in \mathbb{R}^{p \times n}$. The gradient for \mathbf{A} is similar.

- ▶ Finally, since the loss function is explicitly provided, any optimization method that works with MSE loss is usable with the proposed loss.

Results II

- ▶ The MSE loss \tilde{L} and our loss L can be written as

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

- ▶ The analytical gradients are:

$$\begin{aligned}d_{\mathbf{B}}\tilde{L}(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2\langle\mathbf{A}'\boldsymbol{\Sigma}_{yx} - \mathbf{A}'\mathbf{A}\mathbf{B}\boldsymbol{\Sigma}_{xx}, \mathbf{W}\rangle_F, \\ d_{\mathbf{B}}L(\mathbf{A}, \mathbf{B})\mathbf{W} &= -2\langle\mathbf{T}_p\mathbf{A}'\boldsymbol{\Sigma}_{yx} - (\mathbf{S}_{p\circ}(\mathbf{A}'\mathbf{A}))\mathbf{B}\boldsymbol{\Sigma}_{xx}, \mathbf{W}\rangle_F,\end{aligned}$$

in direction of $\mathbf{W} \in \mathbb{R}^{p \times n}$. The gradient for \mathbf{A} is similar.

- ▶ Finally, since the loss function is explicitly provided, any optimization method that works with MSE loss is usable with the proposed loss.