

Introduction to Machine Learning

Lab 5: Matrix Factorization and Linear Dimensionality Reduction

Hongteng Xu

April 3, 2022

1 Motivation

- Implement PCA and achieve data whitening based on it
- Try to explore the algorithm of robust PCA (RPCA) based on the guidance. Learn how to reformulate problems approximately, and do your second alternating optimization algorithm (Do you remember which is your first alternating optimization algorithm in this course?)
- Learn how to change data statistics by machine learning-guided data manipulation.

2 Tasks

Please read Lecture 6 carefully before doing this lab work.

1. Implement a PCA function, which takes the data matrix and the number of principal components you want as its input.
2. Achieve data whitening based on your PCA method.
3. Implement a robust PCA (RPCA) algorithm via alternating optimization. Given a data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$, the original RPCA problem is

$$\min_{\mathbf{X}} \|\mathbf{X}_{noise} - \mathbf{X}\|_1, \quad s.t. \text{rank}(\mathbf{X}) \leq L. \quad (1)$$

Assuming the target $\mathbf{X} = \mathbf{L} + \mathbf{S}$, where \mathbf{L} is low-rank and \mathbf{S} is sparse, we can solve (1) by optimizing the following two subproblems iteratively. At the k -th iteration:

$$\begin{aligned} \text{P1: } \mathbf{L}^{(k)} &= \arg \min_{\mathbf{L}} \|\mathbf{X}_{noise} - \mathbf{L} - \mathbf{S}^{(k-1)}\|_F^2, \quad s.t. \text{rank}(\mathbf{L}) \leq L. \\ \text{P2: } \mathbf{S}^{(k)} &= \arg \min_{\mathbf{S}} \|\mathbf{X}_{noise} - \mathbf{L}^{(k)} - \mathbf{S}\|_F^2, \quad s.t. \|\mathbf{S}\|_0 \leq \tau ND, \end{aligned} \quad (2)$$

where $\|\mathbf{S}\|_0$ counts the number of nonzero elements in \mathbf{S} , and $\tau \in (0, 1)$ controls the ratio of nonzero elements.

4. **Data manipulation:** Given a set of zero-mean data, i.e., $\mathbf{X} \in \mathbb{R}^{N \times D}$, whose covariance matrix is $\mathbf{\Gamma}_{\mathbf{X}} \in \mathbb{R}^{D \times D}$, we have a chance to add two outliers, i.e., \mathbf{x}_1 and \mathbf{x}_2 into \mathbf{X} and formulate a new data matrix $\widehat{\mathbf{X}} \in \mathbb{R}^{(N+2) \times D}$. The two outliers satisfy that $\|\mathbf{x}_1\|_2 = \|\mathbf{x}_2\|_2 = 1$. Try to ensure $\widehat{\mathbf{X}}$ to be a zero-mean data matrix and maximize the difference between its covariance matrix $\mathbf{\Gamma}_{\widehat{\mathbf{X}}}$ and the original $\mathbf{\Gamma}_{\mathbf{X}}$, i.e., $\max \|\mathbf{\Gamma}_{\widehat{\mathbf{X}}} - \mathbf{\Gamma}_{\mathbf{X}}\|_F^2$