

ELEC70037
Topics in Large Dimensional Data Processing
Exercise

Dr Wei Dai

November 25, 2024

EEE Department, Imperial College London

Contents

Contents	i
I. TEACHING ORGANISAION	1
1. Teaching Organisation	2
1.1. Overview	2
1.1.1. Invited Lectures	2
1.2. Coursework Assignments	2
1.2.1. Schedule	2
1.2.2. Marking	3
1.3. Software for Coursework	4
1.3.1. Julia for Programming	4
2. Your Feedback Matters	5
2.1. You Said, We Did	5
2.1.1. Lecture Notes	5
2.1.2. Assessments	5
2.2. Recent Changes	5
2.3. Your Feedback	6
2.4. Our Commitment to You	6
3. Important Notes	7
3.1. Plagiarism	7
3.2. Notes for Coursework Submission	7
3.2.1. Handling Solutions	7
II. COURSEWORK 1	9
4. Linear Algebra: Basics	10
4.1. Trace and Inner Product of Matrices	10
4.2. Adjoint Operator	10
5. Least Squares	13
5.1. The Netflix Problem	13
6. MRI01 - Least Squares	14
6.1. Linear Operator	14
6.2. Least Squares	15
III. COURSEWORK 2	17
7. Greedy Algorithms	18
7.1. Estimation Accuracy	18
7.2. Greedy Algorithms	18
8. Dictionary Learning	19
8.1. The MOD Method	19
8.2. The K-SVD Method	20

IV. COURSEWORK 3	21
9. Proximal Operator	22
10. MRI CS Recovery	24
10.1. Sparsity under Wavelet Transform	24
10.2. MRI CS Recovery: Wavelet and DCT	25
V. COURSEWORK 4	26
11. MRI CS Recovery: ADMM	27
12. Blind Deconvolution: Convex Relaxation and ADMM	28
13. Neural Network Training: Block Coordinate Descent	30

Part I.

TEACHING ORGANISAION

Teaching Organisation

1.

1.1. Overview

- ▶ Lectures
 - Mondays 07/10 - 25/11/2024, 14:00-15:50, Room 305
 - Mondays 02/12 - 09/12/2024, 14:00-15:50, Room 508A&B
- ▶ Office hours
 - Mondays 07/10 - 02/12/2024, 12:00-12:55, Room 503.
- ▶ Communications: Should you have any questions or inquiries:
 - Attend the designated office hours.
 - Post them on the Q&A channel.
 - Refrain from using my college email.
- ▶ Assessment
 - 4 coursework assignments: each counts 25%

1.1.1. Invited Lectures

To enrich students' learning experience with diverse perspectives, particularly in the areas of optimization and machine learning, the lectures on Monday, 9th December (the final Monday of this term) will feature invited seminar lectures presented by guest colleagues. These sessions will offer insights into advanced topics:

- ▶ Advanced Sparsity and Learning Approach for Imaging
Dr Chen Qin,
14:00-14:50, 09/12/2024, Room 508A&B
- ▶ Introduction to Distributed Optimisation and Learning
Dr Stefan Vlaski,
15:00-15:50, 09/12/2024, Room 508A&B

1.2. Coursework Assignments

1.2.1. Schedule

Coursework Overview:

Individual Data + Team Work + Individual Assessment

The procedure for the coursework assignments is outlined below.

- ▶ Each student will receive a unique data file provided by the GTAs.
- ▶ Students collaborate in groups while working on the coursework. The group information is input in an excel file.
- ▶ By the due time, each student must submit their own individual coursework, clearly specifying their contributions to the group effort. Note that submissions cannot be altered after the deadline.

Table 1.1.: Coursework Schedule

	Tasks	Due Date/Time
Coursework 1	Assignment	(Week 3) 14/10
	Due	(Week 4) 25/10 11:59
	Feedback	(Week 6) 04/11
Coursework 2	Assignment	(Week 5) 28/10
	Due	(Week 6) 08/11 11:59
	Feedback	(Week 8) 18/11
Coursework 3	Assignment	(Week 7) 11/11
	Due	(Week 8) 22/11 11:59
	Feedback	(Week 10) 02/12
Coursework 4	Assignment	(Week 9) 25/11
	Due	(Week 10) 06/12
	Feedback (Marks only)	(Week 11) 13/12

- GTAs will assess individual submissions. Let $M_{i,0}$ denote the raw mark received by the student i . An adjusted mark, calculated using Equation (1.1), will then be provided to the student.
- Some coursework questions will be selected for in-depth discussion. Students chosen to lead these discussions will receive bonus marks to encourage active participation.

1.2.2. Marking

The marking formula for the first three coursework is given by

$$M_i = \min \left(\left(\sum_{k \in \mathcal{G}_l} M_{k,0} \right) * C_i * W_l, M_{i,0} \right), \quad (1.1)$$

where

- M_i is the adjusted mark received by the student i ,
- $M_{k,0}$ is the raw mark obtained by the student k ,
- $\sum_{k \in \mathcal{G}_l} M_{k,0}$ is the total raw mark obtained by the group \mathcal{G}_l ,
- C_i is the contribution of the student i in the group,¹
- W_l is the weighting coefficient for the group \mathcal{G}_l , depending on the size of the group,
- and we apply the min function to ensure the adjusted mark does not exceeds the raw mark.²

1: Note that

$$\sum_{k \in \mathcal{G}_l} C_k = 1.00 = 100\%.$$

2: It is typical that $M_i = M_{i,0} \times W_l$, i.e., the adjusted mark is simply the raw mark multiplied by the weighting coefficient.

The weighting coefficients for groups are detailed in Table 1.2. There are several considerations behind the design.

- We encourage both individual work and team work within the group.
- The default size of the group is 2 or 3.
- Students from the same group may share the same codes.
- **Cheating offences and plagiarism** are taken very seriously and are dealt with according to the College's [Academic Misconduct Policy and Procedure](#).

Size of the group	1	2	3	4	5	≥ 6
Weighting coefficient W	1.00	0.97	0.94	0.86	0.70	0.50

Table 1.2.: The weighting coefficient W for the group is decided by the size of the group.

Remark 1.2.1 (Bonus marks) GTAs will select students to present their solutions to specific questions in-depth. Students can earn bonus marks by sharing and explaining their coursework solutions during the feedback/discussion sessions.

The bonus marks awarded will be determined by the GTAs and the module leader, based on the quality of the solutions and the clarity of the explanation. The awarded marks will be communicated to the participating student.

1.3. Software for Coursework

The coursework requires Julia programming. The coursework submission files include a Jupyter notebook file and a data file in JLD2 format. We recommend VSCode as the default editor.

The software that you need to install/have

- ▶ Jupyter Notebook
 - Download and install the package management software [Anaconda](#). By default, it will install Python and Jupyter Notebook to your system.
- ▶ Julia programming language
 - Download and install the current stable version of [Julia](#).
 - Check whether Julia has been installed properly: Run Julia interactive session (a.k.a. REPL) by following the [instructions](#).
- ▶ VSCode
 - Download and install [VSCode](#).
 - We need to install some extensions for VSCode. See [here](#) for how to install VSCode extensions.
 - * Anaconda should have automatically installed VSCode Extensions [Python](#) and [Pylance](#) for you. If not, install them.
 - * Install extension [Julia](#) for Julia programming in VS-Code. See [here](#) and [here](#) for more details.

1.3.1. Julia for Programming

- ▶ [A collection of tutorials.](#)
- ▶ [Introduction to Scientific Programming and Machine Learning with Julia.](#)
- ▶ [A tutorial in pdf format.](#)

Your Feedback Matters

2.

2.1. You Said, We Did

Over the past few years, we've listened carefully to your feedback and implemented substantial changes to enhance your learning experience.

2.1.1. Lecture Notes

- ▶ **Transition to Book Format:** We've transformed our lecture notes from slides into a comprehensive book format. This allows for a more cohesive and self-contained learning experience.
- ▶ **LaTeX Template:** We've adopted the LaTeX template `kaobook`, which gives you more space for note-taking.
- ▶ **Increased Examples:** We've incorporated more examples and illustrations to reinforce key concepts and make the material more engaging.

2.1.2. Assessments

- ▶ **Shift to Coursework:** We've transitioned from traditional paper-based exams to coursework assignments. This approach allows us to focus more on the practical application of theoretical knowledge.
- ▶ **Skill Development:** Coursework assignments are designed to help you hone your analytical and programming skills.
- ▶ **Peer Marking Removed:** We've removed peer marking to ensure that your performance is accurately reflected in your final grade.
- ▶ **Individual Recognition:** Our coursework marking process is designed to fairly assess and appreciate your unique contributions and efforts.

2.2. Recent Changes

Based on your feedback from the last academic year, we've implemented the following changes:

- ▶ **Reduced Coursework Load:** We've reduced the overall workload of coursework assignments by removing a learn-and-tell presentation coursework and spreading the questions across four assignments instead of three.
- ▶ **Coursework submission deadlines.** Coursework deadlines are now on Fridays, providing you with a more relaxed weekend and allowing the department to efficiently handle extension requests.
- ▶ **Enhanced Feedback:** We've increased the amount of information provided with your coursework marks. In addition to your individual scores, we'll share statistical information about the class marks.

- ▶ **In-Lecture Problem Discussion:** To foster deeper understanding and active engagement, we've incorporated in-lecture problem discussions and solving sessions.

2.3. Your Feedback

Please take the time to complete the **MEQ** questionnaires at the end of the term. Your continued feedback is invaluable in helping us improve the learning experience for this module.

2.4. Our Commitment to You

- ▶ **Accessibility and Support:** We're dedicated to providing you with the support you need to succeed. We offer regular **office hours** (Mondays, 12:00-12:55, Room 503) where you can discuss any questions or difficulties you may encounter. Additionally, we maintain a **Q&A channel** on Teams for more immediate assistance.
- ▶ **Timely Feedback:** We understand the importance of timely feedback. We aim to complete coursework marking within 10 working days of the submission deadline. You'll receive your marked assignments along with statistical information about the class performance.
- ▶ **Engaging Learning Experiences:** We're committed to creating an engaging and supportive learning environment. **In-lecture problem discussions and solving sessions** provide opportunities for you to actively participate, ask questions, and deepen your understanding of the material.

3.1. Plagiarism

PLAGIARISM/COLLUSION DECLARATION

Coursework submitted for assessment must be the original work of you and your group. Assignments are subjected to regular checks for plagiarism and/or collusion. Plagiarism is the presentation of another person's thoughts or words (those outside your group) as if they were your own. Collusion involves obtaining help from someone outside your group to complete your work. In preparing your coursework, you should not seek help, or copy from any other person or source, including the Internet, without proper and explicit acknowledgement.

There is a procedure in place for you to declare individual contributions within your group for coursework. You must declare the contributions fairly and accurately.

You must not disclose your solutions or insights related to coursework with anyone else, including future students or the Internet.

By acknowledging the the statements above, you are declaring that both this and all subsequent pieces of coursework are, and will remain, the original work of you and your group.

- ▶ Submissions will not be accepted without the aforementioned declaration.
- ▶ Members of a group are deemed to have **collective responsibility** for the integrity for work submitted and are liable for any penalty imposed, proportionate to their contributions.

3.2. Notes for Coursework Submission

- ▶ Each registered student will get a data file. The data in the data file can be unique.
- ▶ Each registered student needs to submit the solutions related to their own data, no matter whether they are in groups or not.

3.2.1. Handling Solutions

In our coursework, we use the following convention.

- ▶ If the solution to a question is a unique integer, you need to assign an integer value to your solution variable.
- ▶ If the solution to a question does not exist, you need to create a 1-D array of length zero.
- ▶ If the solution to a question is not unique, you need to create a 1-D array, of which the length is the number of distinct solutions, and then specify the values in the **ascending** order.

Examples:

```
1      x = 3;  
2      y = Array{Int64,1}(undef,0);  
3      z = Array{Int64,1}(undef,3);  
4      z[1] = 3; z[2] = 4; z[3] = 5;
```

Part II.

COURSEWORK 1

4.1. Trace and Inner Product of Matrices

Give a square matrix $A \in \mathbb{R}^{n \times n}$, its trace is defined as the sum of the diagonal elements:

$$\text{trace}(A) := \sum_{i=1}^n A_{i,i}.$$

Use the parameters in your data file to answer the following questions.

1. Let $A, B \in \mathbb{R}^{m \times n}$.

- a) Identify the dimension of AB^T . [2]
- b) Identify the dimension of $A^T B$. [2]
- c) Is the following statement true or false?

$$\text{trace}(AB^T) = \text{trace}(A^T B).$$

[1]

2. For any two matrices $A, B \in \mathbb{R}^{m \times n}$, define the **inner product** between them:

$$\langle A, B \rangle := \text{trace}(A^T B).$$

Is the following statement true or false?

$$\langle A, B \rangle = \langle \text{vect}(A), \text{vect}(B) \rangle.$$

[1]

3. Let $A \in \mathbb{R}^{k \times l}$ and $B \in \mathbb{R}^{m \times n}$.

- a) Specify the dimensions of U and V such that

$$\langle UAV^T, B \rangle$$

is well-defined. [4]

- b) Is the following statement true or false?

$$\langle UAV^T, B \rangle = \langle A, U^T B V \rangle.$$

[1]

4.2. Adjoint Operator

Let \mathcal{A} be a linear operator mapping from \mathbb{R}^n to \mathbb{R}^m . It always can be represented by a matrix $A \in \mathbb{R}^{m \times n}$, that is,

$$\begin{aligned} \mathcal{A} : \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ x &\mapsto y = Ax. \end{aligned}$$

Definition 4.2.1 The adjoint of the linear operator \mathcal{A} , denoted by \mathcal{A}^* , is defined via

$$\langle \mathcal{A}(x), y \rangle = \langle x, \mathcal{A}^*(y) \rangle$$

for all $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$.

It can be verified that

$$\mathcal{A}^*(y) = A^T y.$$

If $\mathcal{A} : \mathbb{C}^n \rightarrow \mathbb{C}^m$, then \mathcal{A} is generally represented by a matrix $A \in \mathbb{C}^{m \times n}$. That is,

$$\mathcal{A}^*(y) = A^H y,$$

where the superscript H denotes ‘conjugate transpose’.

The following exercise questions involve identification of linear operators and their adjoint operators. One way to check the correctness of your adjoint operator is to see whether the equality

$$\langle \mathcal{A}(x), \mathcal{A}(y) \rangle = \langle x, \mathcal{A}^*(\mathcal{A}(y)) \rangle$$

holds.

1. Consider the subsampling operator \mathcal{P}_Ω where $\Omega = \{i_1, i_2, \dots, i_m\}$ is a subset of $[1 : n] := \{1, 2, \dots, n\}$ ($m \leq n$), defined as

$$\begin{aligned} \mathcal{P}_\Omega : \mathbb{C}^n &\rightarrow \mathbb{C}^m \\ x &\mapsto y \text{ with } y_i = x_{i_i}. \end{aligned}$$

- a) Write a function to find the corresponding matrix representations of \mathcal{P}_Ω and \mathcal{P}_Ω^* .
Use the data in the data file for a test. [4]
- b) Write a function to evaluate $y = \mathcal{P}_\Omega(x)$ without using the matrix representation of the linear map.
Use the data in the data file for a test. [2]
- c) Write a function to evaluate $z = \mathcal{P}_\Omega^*(y)$ without using the matrix representation of the linear map.
Use the data in the data file for a test. [2]

2. Consider the Hankel operator given by

$$\mathcal{H}_{p,q} : \mathbb{C}^{p+q-1} \rightarrow \mathbb{C}^{p \times q}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{p+q-1} \end{bmatrix} \mapsto H = \begin{bmatrix} x_1 & x_2 & x_3 & \cdots & x_q \\ x_2 & x_3 & x_4 & \cdots & x_{q+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_p & x_{p+1} & x_{p+2} & \cdots & x_{p+q-1} \end{bmatrix}.$$

Or equivalently,

$$H_{m,n} = x_{m+n-1}.$$

- a) Write a function to find the corresponding matrix representations of $\mathcal{H}_{p,q}$ and $\mathcal{H}_{p,q}^*$.

- Use the data in the data file for a test. [4]
- b) Write a function to evaluate $\mathbf{H} = \mathcal{H}_{p,q}(\mathbf{x})$ without using the matrix representation of the linear map.
Use the data in the data file for a test. [2]
- c) Write a function to evaluate $\mathbf{z} = \mathcal{H}_{p,q}^*(\mathbf{H})$ without using the matrix representation of the linear map.
Use the data in the data file for a test. [2]
- d) Evaluate $\mathbf{z} ./ \mathbf{x}$ where the symbol $./$ denotes element division. [1]

5.1. The Netflix Problem

1. The task is to recover an unknown matrix X from its partial observations given by

$$\mathbf{y} = \mathcal{P}_\Omega(X) + \mathbf{w},$$

where \mathbf{w} denotes noise.

The alternating minimization method to solve the Netflix problem is stated below. Suppose the prior information that the rank of X is at most r . Then one would find the matrices $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times r}$ such that $X = \mathbf{U}\mathbf{V}^\top$. In the k -th iteration of the alternating minimization, one solves the following two least squares sequentially:

$$\begin{aligned} \mathbf{V}^{k+1} &= \arg \min_{\mathbf{V}} \left\| \mathbf{y} - \mathcal{P}_\Omega(\mathbf{U}^k \mathbf{V}^\top) \right\|_2^2, \\ \mathbf{U}^{k+1} &= \arg \min_{\mathbf{U}} \left\| \mathbf{y} - \mathcal{P}_\Omega(\mathbf{U} \mathbf{V}^{k+1, \top}) \right\|_2^2. \end{aligned}$$

The partial observations \mathbf{y} and the index sets of the observed entries Ω are given in the data file. Suppose that $\text{rank}(X) = 10$.

- a) Write a function to update \mathbf{V} for a given \mathbf{U} . Use the \mathbf{U}_0 in the data file for a test. [2]
- b) Write a function to update \mathbf{U} for a given \mathbf{V} . Use the \mathbf{V} from the previous sub-question for a test. [2]
- c) Write a function for the alternating minimization process. Use the \mathbf{U}_0 in the data file as the initial point. Run a test. [2]

6.1. Linear Operator

This question requires background in

- ▶ 2D Discrete Fourier Transform. See Wikipedia page https://en.wikipedia.org/wiki/Discrete_Fourier_transform, Sections 'Definition', 'Properties > The unitary DFT', and 'Multidimensional DFT'.
- ▶ `fftshift`: See <https://www.matecdev.com/posts/julia-fft.html#fftshift-and-fftfreq>.

Let $X \in \mathbb{R}^{m \times n}$ be an MRI image. The MRI data compressed sensing can be mathematically expressed as

$$y = \underbrace{\mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}}_{\mathcal{A}}(X)$$

where

- ▶ \mathcal{F}_{2D} denotes the 2D Discrete Fourier Transform,
- ▶ \mathcal{S} denotes the 2D `fftshift` operator (which is the straightforward extension of 1D `fftshift`), and
- ▶ \mathcal{P}_Ω is the subsampling operator.

In the following, you are allowed to use `FFTW.fft`, `FFTW.ifft`, and `FFTW.fftshift`.

Make sure that your function works for both odd and even m and n .

1. Write a function to evaluate

$$y = \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}(X) = \mathcal{A}(X)$$

without using matrix representation of the linear map \mathcal{A} .

The outputs of your function must include $O_1 = \mathcal{F}_{2D}(X)$, $O_2 = \mathcal{S} \circ \mathcal{F}_{2D}(X)$, and $O_3 = \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}(X)$.

Use the data in the data file for a test. [6]

2. Write a function to evaluate

$$Z = (\mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D})^*(y) = \mathcal{A}^*(y).$$

without using matrix representation of the linear map.

The outputs of your function must include $O_1 = \mathcal{P}_\Omega^*(y)$, $O_2 = \mathcal{S}^* \circ \mathcal{P}_\Omega^*(y)$, and $O_3 = \mathcal{F}_{2D}^* \circ \mathcal{S}^* \circ \mathcal{P}_\Omega^*(y)$.

Use the data in the data file for a test. [6]

6.2. Least Squares

1. The goal of this question is to find the pseudo-inverse of the linear operator

$$\mathcal{A} = \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}$$

without using matrix representation of the linear map.

- a) Let $\mathbf{Y}_1 = \mathcal{F}_{2D}(\mathbf{X})$. Compute $\mathbf{Z}_1 = \mathcal{F}_{2D}^\dagger(\mathbf{Y}_1)$. [1]
- b) Let $\mathbf{Y}_2 = \mathcal{S}(\mathbf{Y}_1)$. Compute $\mathbf{Z}_2 = \mathcal{S}^\dagger(\mathbf{Y}_2)$. [1]
- c) Let $\mathbf{y}_3 = \mathcal{P}_\Omega(\mathbf{Y}_2)$. Compute $\mathbf{Z}_3 = \mathcal{P}_\Omega^\dagger(\mathbf{y}_3)$. [1]
- d) Write a function to compute $\mathbf{Z} = \mathcal{A}^\dagger(\mathbf{y})$, where $\mathbf{y} = \mathcal{A}(\mathbf{X})$. [1]

Remark 6.2.1 Further insights into the pseudo-inverse can be gained by relating your solution to the singular value decomposition (SVD) of the matrix representations of the linear operators involved.

2. Consider the following recovery of undersampled MRI data:

$$\begin{aligned} \hat{\mathbf{X}} &= \arg \min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}(\mathbf{X})\|_2^2 + \frac{\alpha}{2} \|\mathbf{X}\|_F^2 \\ &= \arg \min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \frac{\alpha}{2} \|\mathbf{X}\|_F^2 \end{aligned}$$

The above formulation can be written as

$$\left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathcal{B}(\mathbf{X}) \right\|_2^2,$$

where $\mathbf{0}$ denotes a zero vector of proper dimension.

- a) Find the dimension of the above zero vector. [1]
- b) Write a function to compute

$$\mathbf{b} = \mathcal{B}(\mathbf{X}).$$

without using matrix representations of linear operators. [2]

- c) Write a function to compute

$$\mathbf{Z} = \mathcal{B}^* \left(\begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} \right)$$

without using matrix representations of linear operators. [2]

- d) It is clear that

$$\hat{\mathbf{X}} = (\mathcal{B}^* \mathcal{B})^{-1} \mathcal{B}^* \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} = \mathcal{B}^\dagger \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}.$$

Write a function to compute $\hat{\mathbf{X}}$. [5]

(Hint: Understanding of the SVD of the matrix representation of $\mathcal{B}^* \mathcal{B}$ is helpful in solving this problem. Let \mathbf{B} be the matrix representation of \mathcal{B} . The matrix $\mathbf{B}^H \mathbf{B}$ corresponds to the matrix representation of $\mathcal{B}^* \mathcal{B}$, and its SVD exhibits a convenient structure. Although deriving the explicit SVD of $\mathcal{B}^* \mathcal{B}$ is beyond the scope of this exercise, this insight will help

you to efficiently compute $(\mathcal{B}^*\mathcal{B})^{-1}$.)

Remark 6.2.2 The questions in this chapter adopt a matrix-free approach. As every finite-dimensional linear operator can be represented as a matrix, the standard method for solving least squares problems involves converting the problem into a matrix/vector form and then solving it using matrix (pseudo-)inverse. However, this approach is often computationally expensive and doesn't scale well. Instead, the questions are designed to explore the structure of the linear operators without explicit matrix inverse, allowing for much more efficient computations.

Part III.

COURSEWORK 2

Greedy Algorithms

7.

7.1. Estimation Accuracy

1. Alice, Bob, and Charlie's jobs are to estimate unknown x from their respective measurement vector y given by

$$y = Ax_0 + w,$$

by solving the minimization problem

$$\hat{x} = \arg \min_z \frac{1}{2} \|y - Az\|^2.$$

Daniel knows the ground truth x_0 and has the power to choose the noise vector w with $\|w\|_2 = 1$ in generating y .

- a) Daniel is a friend of Alice. He would like to choose a noise vector w so that the estimation error

$$\|\hat{x} - x_0\|^2$$

is minimised.

Find the best w and compute the corresponding estimation error for Alice. [2]

- b) Daniel is a foe of Bob. He would like to choose a noise vector w so that the estimation error

$$\|\hat{x} - x_0\|^2$$

is maximised.

Find the best w and compute the corresponding estimation error for Bob. [2]

- c) For Charlie, Daniel randomly generates a noise vector w with i.i.d. Gaussian entries and then normalises it so that $\|w\|_2 = 1$ (i.e., w has unit l_2 norm and statistically invariant under rotations). Find the expectation of the estimation error, i.e.,

$$\mathbb{E} \left[\|\hat{x} - x_0\|^2 \right]$$

that Charlie may expect.

[2]

7.2. Greedy Algorithms

1. Implement the OMP algorithm. Test its performance using the data given in the data file. [4]
2. Implement the SP algorithm. Test its performance using the data given in the data file. [4]
3. Implement the HIT algorithm. Test its performance using the data given in the data file. [4]

Dictionary Learning 8.

The dictionary learning problem refers to the task of learning a set of atoms that can be used to efficiently represent data. The mathematical formulation of the dictionary learning problem is as follows. Given a set of data consisting of N data points

$$\mathbf{Y} = [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \cdots \quad \mathbf{y}_N] \in \mathbb{R}^{M \times N},$$

dictionary learning seeks to solve the optimisation problem:

$$\begin{aligned} \min_{\mathbf{D}, \mathbf{C}} \quad & \|\mathbf{Y} - \mathbf{DC}\|_F^2 \\ \text{subject to} \quad & \|\mathbf{d}_k\|_2 = 1, \quad 1 \leq k \leq K \\ & \|\mathbf{c}_n\|_0 \leq S, \quad 1 \leq n \leq N, \end{aligned}$$

where \mathbf{D} is the dictionary matrix, with each column \mathbf{d}_k being a dictionary atom, the constraint $\|\mathbf{d}_k\|_2 = 1$ is to normalise the ℓ_2 -norm of each atom, \mathbf{C} is the coefficient matrix, with each column \mathbf{c}_n representing the sparse coefficients for \mathbf{y}_n , and the constraint $\|\mathbf{c}_n\|_0 \leq S$ means that each data point is represented using only up to S many dictionary atoms.

8.1. The MOD Method

The MOD (Method of Optimal Directions) is a simple algorithm for solving the dictionary learning problem through alternating minimisation. Each iteration of the MOD algorithm consists of two main steps: sparse coding and dictionary update. In the sparse coding step, the current dictionary \mathbf{D} is used to update the sparse coefficients \mathbf{C} by solving the following problem:

$$\begin{aligned} \min_{\mathbf{C}} \quad & \|\mathbf{Y} - \mathbf{DC}\|_F^2 \\ \text{subject to} \quad & \|\mathbf{c}_n\|_0 \leq S, \quad 1 \leq n \leq N. \end{aligned}$$

In the dictionary update step, the entire dictionary is updated by solving a least squares problem:

$$\min_{\mathbf{D}} \|\mathbf{Y} - \mathbf{DC}\|_F^2$$

followed by normalising each column of the dictionary to have unit ℓ_2 -norm.

The data given in the data file includes the dataset \mathbf{Y} , and the initial values of the dictionary and sparse coefficients, denoted by \mathbf{D}_0 and \mathbf{C}_0 respectively. In the j -th iteration of the MOD algorithm, the process starts with \mathbf{C}_{j-1} and \mathbf{D}_{j-1} . After the sparse coding step, \mathbf{C}_{j-1} is updated into \mathbf{C}_j , and following the dictionary update step, \mathbf{D}_{j-1} is updated to \mathbf{D}_j .

1. Implement the sparse coding step as a function `SC_OMP` using the OMP algorithm. Test it by using \mathbf{Y} and \mathbf{D}_0 as inputs, and generate the output \mathbf{C}_{OMP} . [3]

Engan, Kjersti, Sven Ole Aase, and J. Hakon Husoy. "Method of optimal directions for frame design." In IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), vol. 5, pp. 2443-2446. IEEE, 1999.

2. Implement the sparse coding step as a function `SC_SP` using the Subspace Pursuit (SP) algorithm. Unlike the standard SP algorithm, here C_{j-1} is used as the initialisation of the SP process in the j -th iteration of the MOD algorithm. Test it by using Y , D_0 , and C_0 as inputs, and generating the output C_{SP} . [3]
3. Implement the dictionary update step as a function `DU_MOD`. Test it by using Y and C_0 as inputs, and generating the output D_{MOD} . [3]
4. Now implement the full MOD algorithm for dictionary learning by combining `SC_SP` and `DU_MOD`. Run 50 iterations. Save C_{50} and D_{50} into the data file. [3]
5. Let

$$f(C, D) = \|Y - DC\|_F^2.$$

Compute and save $f(C_j, D_j)$, and the norm of the gradient of f with respect to C , i.e., $\|\nabla_C f(C_j, D_j)\|$, $j = 1, 2, \dots, 50$. [3]

8.2. The K-SVD Method

K-SVD (K-means Singular Value Decomposition) is a widely-used algorithm designed to solve the dictionary learning problem. Like the MOD method, each iteration of K-SVD consists of a sparse coding step and a dictionary update step. The sparse coding step in K-SVD is the same as in the MOD algorithm. However, the dictionary update step in K-SVD is significantly different. Specifically, K-SVD updates dictionary atoms d_k one by one. For each column $k = 1, 2, \dots, K$ in D_{j-1} , the update follows the following steps:

- (a) Identify the group of data samples that use atom d_k

$$\omega_k = \{n : C_{k,n} \neq 0, 1 \leq n \leq N\}.$$

- (b) Compute the overall representation error matrix, E_k , by

$$E_k = Y - \sum_{p \neq k} d_p C_{p,:}.$$

- (c) Restrict E_k to the columns corresponding to ω_k , yielding $E_{k,\omega}$.
- (d) Perform SVD on $E_{k,\omega}$: $E_{k,\omega} = U \Sigma V^T$. Update the k -th column of the dictionary d_k to be the first column of U , and the coefficient vector C_{k,ω_k} to be the first column of V multiplied by $\Sigma_{1,1}$.

It is important to note that during the dictionary update step, both the dictionary and the sparse coefficients are updated, while the sparsity pattern remains unchanged.

1. Implement the dictionary update step of K-SVD as a function `DU_KSVD`. Test it by using Y , C_0 , and D_0 as inputs, and generating the output C_{KSVD} and D_{KSVD} . [3]
2. Implement the full K-SVD algorithm by combining `SC_SP` and `DU_KSVD`. Run 50 iterations. Save C_{50} and D_{50} into the data file. [3]
3. Compute and save $f(C_j, D_j)$, and the norm of the gradient of f with respect to C , i.e., $\|\nabla_C f(C_j, D_j)\|$, $j = 1, 2, \dots, 50$. [3]

Aharon, Michal, Michael Elad, and Alfred Bruckstein. "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation." *IEEE Transactions on signal processing*, no. 11, vol. 54, (2006): 4311-4322.

Part IV.

COURSEWORK 3

Proximal Operator

9.

1. Use the data in the data file to compute the solutions of the following optimization problems.

a)

$$\min_x \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x} + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

[2]

b)

$$\min_x \delta(\mathbf{a} \leq \mathbf{x} \leq \mathbf{b}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2,$$

where $\mathbf{a} \leq \mathbf{x} \leq \mathbf{b}$ denotes element-wise inequalities.

[2]

c)

$$\min_x \delta(\|\mathbf{x}\|_0 \leq k) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

[2]

d)

$$\min_x \|\mathbf{x}\|_0 + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

[2]

e)

$$\min_{\mathbf{X}} \delta(\text{rank}(\mathbf{X}) \leq r) + \frac{1}{2\gamma} \|\mathbf{X} - \mathbf{Z}\|_F^2.$$

[2]

f)

$$\min_{\mathbf{X}} \text{rank}(\mathbf{X}) + \frac{1}{2\gamma} \|\mathbf{X} - \mathbf{Z}\|_F^2.$$

[2]

g)

$$\min_{\mathbf{X}} \delta(\mathbf{X} \geq 0) + \frac{1}{2\gamma} \|\mathbf{X} - \mathbf{Z}\|_F^2,$$

where $\mathbf{X} \geq 0$ denotes that \mathbf{X} is a positive semi-definite matrix, i.e., all the eigenvalues of the symmetric matrix \mathbf{X} are non-negative.

[2]

h)

$$\min_x \|\mathbf{x}\|_2 + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

[2]

i)

$$\min_x \delta(\|\mathbf{x}\|_2 = 1) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

[2]

2. (Proximal operators related to ℓ_1 -norm) Use the data in the data file to compute the solutions of the following optimization problems.

a)

$$\min_x \lambda \|\mathbf{x}\|_1 + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

[2]

b)

$$\min_x \lambda \|\mathbf{x}\|_1 + \frac{1}{2\gamma} \|\text{diag}(\mathbf{a})\mathbf{x} - \mathbf{z}\|_2^2,$$

where $\text{diag}(\mathbf{a})$ turns the vector \mathbf{a} into a diagonal matrix with the diagonal vector \mathbf{a} .

[2]

c)

$$\min_x \lambda \|\mathbf{x}\|_1 + \frac{1}{2\gamma} \|\mathbf{U}\mathbf{x} - \mathbf{z}\|_2^2,$$

where the matrix \mathbf{U} is a unitary matrix, i.e., $\mathbf{U}\mathbf{U}^\top$ and $\mathbf{U}^\top\mathbf{U} = \mathbf{I}$.

[2]

d)

$$\min_x \lambda \|\mathbf{U}\mathbf{x}\|_1 + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|_2^2,$$

where the matrix \mathbf{U} is a unitary matrix, i.e., $\mathbf{U}\mathbf{U}^\top$ and $\mathbf{U}^\top\mathbf{U} = \mathbf{I}$.

[2]

10.1. Sparsity under Wavelet Transform

Let X represent an unknown MRI image. A reasonable assumption is that the wavelet transform coefficients of X are sparse.

We consider the following MRI compressed sensing recovery formulation:

$$\begin{aligned}\hat{X} &= \arg \min_X f(X) \\ &:= \arg \min_X \lambda \|\mathcal{W}(X)\|_1 + \frac{1}{2} \|\mathbf{y} - \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}(X)\|_2^2 \\ &= \arg \min_X \underbrace{\lambda \|\mathcal{W}(X)\|_1}_{h(X)} + \underbrace{\frac{1}{2} \|\mathbf{y} - \mathcal{A}(X)\|_2^2}_{q(X)},\end{aligned}$$

where $\|\cdot\|_1$ is the sum of the absolute values of the entries, \mathcal{W} denotes the 2D discrete Daubechies wavelet transform (which can be represented by an orthonormal matrix) (See `Wavelets.jl` for Julia functions `dwt` and `idwt`), the operators $\mathcal{P}_\Omega, \mathcal{S}, \mathcal{F}_{2D}$ have been defined in coursework 1, and for simplicity we have used the notation

$$\mathcal{A} = \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}.$$

We apply the proximal gradient method to solve the optimization problem above. Let X^k denote the value of X after the k -th iteration, which also serves as the initial value for the $k + 1$ -th iteration.

1. Use the X^0 given in the data file to compute the gradient

$$\nabla_X q(X^0) = \nabla_X \left(\frac{1}{2} \|\mathbf{y} - \mathcal{A}(X^0)\|_2^2 \right).$$

[3]

2. Let $X \in \mathbb{R}^{n \times n}$ (the value of n can be inferred from the data that you are given). Let L be the Lipschitz constant of $\nabla_X q(X)$. Find the value of L . [2]
3. Set the step size of the proximal gradient algorithm as $\tau = \frac{1}{L}$ where L is the Lipschitz constant of the gradient of $q(X)$. Find the X_1 after the first iteration of the proximal gradient method. [2]
4. To monitor the convergence of the algorithm, we compute $\|g^k\|_2$ at the end of each iteration k , where $g^k \in \partial f(X^k)$. Calculate $\|g^k\|_2$, $k = 1, 2, \dots, 10$, and save these values as a vector in the data file. [2]
5. Run the proximal gradient method that you implemented until it converges, i.e., $\|g^k\|_2 \leq 10^{-3} \times \|X^k\|_F$. Draw your recovery using the entry magnitudes. [2]

10.2. MRI CS Recovery: Wavelet and DCT

Let X be an unknown MRI image. We consider the following MRI compressed sensing recovery formulation:

$$\begin{aligned}\hat{X} &= \arg \min_{X, Z_1, Z_2} f(X, Z_1, Z_2) \\ &= \arg \min_{X, Z_1, Z_2} \underbrace{\lambda \|Z_1\|_1 + \lambda \|Z_2\|_1}_{h(Z_1, Z_2)} \\ &\quad + \underbrace{\frac{\alpha}{2} \|Z_1 - \mathcal{W}(X)\|_2^2 + \frac{\alpha}{2} \|Z_2 - \mathcal{D}(X)\|_2^2 + \frac{1}{2} \|y - \mathcal{A}(X)\|_2^2}_{q(X, Z_1, Z_2)},\end{aligned}$$

where where \mathcal{D} denotes the 2D discrete cosine transform (which can be represented by an orthonormal matrix) (See `FFTW.jl` for Julia functions `dct` and `idct`), and \mathcal{W} and \mathcal{A} are defined as in the previous question.

We use alternating minimization method to solve the above optimization problem.

1. Use the data in the data file including Z_1^0, Z_2^0 to find X^1 . [3]
2. Use the obtained X^1 to find Z_1^1 and Z_2^1 . [3]
3. Evaluate

$$\|g^k\|_2, \quad \text{where } g^k \in \partial_{(X, Z_1, Z_2)} f(X^k, Z_1^k, Z_2^k),$$

for $k = 1, 2, \dots, 10$. [2]

4. Implement and run the alternating minimization method until it converges, i.e., $\|g^k\|_2 \leq 10^{-3} \times \|[X^k, Z_1^k, Z_2^k]\|_F$. Draw your recovery using the entry magnitudes. [2]

Part V.

COURSEWORK 4

MRI CS Recovery: ADMM

Let $X \in \mathbb{R}^{n \times n}$ be an unknown MRI image. We consider the following MRI compressed sensing recovery formulation:

$$\hat{X} = \arg \min_X \frac{1}{2} \|y - \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}(X)\|_2^2 + \lambda_1 \|\mathcal{W}(X)\|_1 + \lambda_2 \|\mathcal{D}(X)\|_1,$$

where \mathcal{P}_Ω is the under-sampling operator, \mathcal{S} is the shift operator, \mathcal{F}_{2D} is the two dimensional FFT transform, \mathcal{W} is the two-dimensional discrete (Daubechies) wavelet transform (Julia package `Wavelets`), and \mathcal{D} is the two-dimensional discrete cosine transform (Julia package `FFTW`).

The above optimization is convex and can be solved using ADMM, we reformulate the above unstrained optimization problem into the following constrained counterpart:

$$\begin{aligned} \min_{X, Z_{1:2}} \quad & \frac{1}{2} \|y - \mathcal{A}(X)\|_2^2 + \lambda_1 \|Z_1\|_1 + \lambda_2 \|Z_2\|_1 \\ \text{s.t.} \quad & Z_1 = \mathcal{W}(X), \quad Z_2 = \mathcal{D}(X), \end{aligned}$$

where $\mathcal{A} = \mathcal{P}_\Omega \circ \mathcal{S} \circ \mathcal{F}_{2D}$.

The corresponding augmented Lagrangian is then given by

$$\begin{aligned} L_\rho = & \frac{1}{2} \|y - \mathcal{A}(X)\|_2^2 + \lambda_1 \|Z_1\|_1 + \lambda_2 \|Z_2\|_1 \\ & + \langle U_1, Z_1 - \mathcal{W}(X) \rangle + \frac{\rho}{2} \|Z_1 - \mathcal{W}(X)\|_F^2 \\ & + \langle U_2, Z_2 - \mathcal{D}(X) \rangle + \frac{\rho}{2} \|Z_2 - \mathcal{D}(X)\|_F^2. \end{aligned} \quad (11.1)$$

The data file provide the initial values for all variables and the values of the key parameters.

1. Write a function to update X when all other variables are fixed. Save the value of X^1 into the data file. [3]
 2. Write a function to update $Z_{1:2}$ when all other variables are fixed. Save the values of $Z_{1:2}^1$ into the data file. [4]
 3. Write a function to update $U_{1:2}$ when all other variables are fixed. Save the values of $U_{1:2}^1$ into the data file. [2]
 4. Compute the primal residual at the end of the first iteration R^1 . [2]
 5. Compute the dual residual at the end of the first iteration S^1 . [2]
 6. Now write a function to implement the ADMM algorithm including the stopping criteria mentioned in the lecture notes. Here we use the error tolerance constants $\epsilon_{\text{abs}} = \epsilon_{\text{rel}} = 10^{-3}$. The maximum number of iteration is 200. Save your final result \hat{X} into the data file. [2]
- Draw your final result \hat{X} as an image in the Jupyter notebook. [1]

The formulation uses the combination of wavelet and DCT transforms. It has been observed in the literature that wavelet transform captures point-like features, and DCT enhances homogeneous texture components.

Blind Deconvolution: Convex Relaxation and ADMM

12.

This chapter addresses the blind deconvolution problem, which involves recovering a signal from its blurred version without prior knowledge of the blurring function.

Specifically, let

$$\mathbf{y} = \mathbf{s} \circledast \mathbf{h},$$

where $\mathbf{h} \in \mathbb{R}^m$, $\mathbf{s} \in \mathbb{R}^n$, and $\mathbf{y} \in \mathbb{R}^{n+m-1}$. By definition of convolution, we have

$$y[t] = \sum_{\tau=0}^{n-1} s[\tau]h[t-\tau].$$

Assume that \mathbf{s} is sparse under DCT transform, meaning $\mathcal{D}(\mathbf{s})$ is sparse. A convex relaxation of the blind deconvolution can be formulated as

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda_1 \|\mathbf{X}\|_* + \lambda_2 \|\mathcal{D}(\mathbf{X})\|_{2,1},$$

where $\mathbf{X} = \mathbf{s}\mathbf{h}^\top \in \mathbb{R}^{n \times m}$ is a “matrix lifting” of the convolution, \mathcal{A} denotes the linear operator corresponding to the convolution (see lecture notes for further details), $\|\cdot\|_*$ denotes the nuclear norm, $\mathcal{D}(\mathbf{X}) \in \mathbb{R}^{n \times m}$ applies the DCT transform to each column of \mathbf{X} , and

$$\|\mathbf{B}\|_{2,1} = \sum_{l=1}^n \|\mathbf{B}_{:,l}\|_2$$

is the summation of ℓ_2 -norms of the columns of a matrix, promoting column sparsity.

The above optimization problem can be solved using ADMM. We reformulate the above unconstrained optimization problem into a constrained one:

$$\begin{aligned} \min_{\mathbf{X}, \mathbf{Z}_{1,2}} \quad & \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda_1 \|\mathbf{Z}_1\|_* + \lambda_2 \|\mathbf{Z}_2^\top\|_{2,1} \\ \text{s.t.} \quad & \mathbf{X} = \mathbf{Z}_1, \mathcal{D}(\mathbf{X}) = \mathbf{Z}_2. \end{aligned}$$

The corresponding augmented Lagrangian is then given by

$$\begin{aligned} L_\rho = & \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda_1 \|\mathbf{Z}_1\|_* + \lambda_2 \|\mathbf{Z}_2^\top\|_{2,1} \\ & + \langle \mathbf{u}_1, \mathbf{Z}_1 - \mathbf{X} \rangle + \frac{\rho}{2} \|\mathbf{Z}_1 - \mathbf{X}\|_F^2 \\ & + \langle \mathbf{u}_2, \mathbf{Z}_2 - \mathcal{D}(\mathbf{X}) \rangle + \frac{\rho}{2} \|\mathbf{Z}_2 - \mathcal{D}(\mathbf{X})\|_F^2. \end{aligned}$$

The data file provide the initial values for all variables and the values of the key parameters.

1. Write a function to update \mathbf{X} when all other variables are fixed. Save the value of \mathbf{X}^1 into the data file. [3]

2. Write a function to update $\mathbf{Z}_{1:2}$ when all other variables are fixed. Save the values of $\mathbf{Z}_{1:2}^1$ into the data file. [4]
3. Write a function to update $\mathbf{U}_{1:2}$ when all other variables are fixed. Save the values of $\mathbf{U}_{1:2}^1$ into the data file. [2]
4. Compute the primal residual at the end of the first iteration \mathbf{R}^1 . [2]
5. Compute the dual residual at the end of the first iteration \mathbf{S}^1 . [2]
6. Write a function to implement the ADMM algorithm including the stopping criteria mentioned in the lecture notes. Here we use the error tolerance constants $\epsilon_{\text{abs}} = \epsilon_{\text{rel}} = 10^{-3}$. The maximum number of iteration is 200. Save your final result $\hat{\mathbf{X}}$ into the data file. [2]
7. Apply SVD to extract $\hat{\mathbf{s}}$ and $\hat{\mathbf{h}}$ from $\hat{\mathbf{X}}$ such that $\|\mathbf{h}\|_2 = 1$ and $h[1] > 0$. [1]

Neural Network Training: Block Coordinate Descent

13.

For demonstration purposes, we will examine the training of a simplified neural network with a single hidden layer, using proximal algorithms instead of the standard gradient backpropagation.

This two-layer network can be described as

$$f(x; W_{1:2}) = h_2(W_2 h_1(W_1 x)),$$

where, for simplicity, we assume both the activation functions h_1 and h_2 to be the Relu function:

$$h_1(x) = h_2(x) = \begin{cases} x, & \text{if } x > 0, \\ 0, & \text{otherwise,} \end{cases}$$

and the activation function is applied element-wise to the input vectors and matrices. Denote the loss function as l , the training problem can be posed as

$$\min_{W_{1:2}} l(f(x; W_{1:2}), y).$$

For simplicity, we assume a quadratic loss function, giving:

$$\min_{W_{1:2}} \frac{1}{2} \|y - f(x; W_{1:2})\|_2^2.$$

We relax the above optimisation problem by introducing auxiliary variables and applying quadratic penalties. This leads to the following unconstrained optimisation problem:

$$\begin{aligned} \min_{x_1, z_{1:2}, W_{1:2}} & \frac{1}{2} \|y - h_2(z_2)\|_2^2 + \frac{\alpha}{2} \|z_2 - W_2 x_1\|_2^2 \\ & + \frac{\alpha}{2} \|x_1 - h_1(z_1)\|_2^2 + \frac{\alpha}{2} \|z_1 - W_1 x\|_2^2. \end{aligned}$$

When extended to multiple training samples, this becomes

$$\begin{aligned} \min_{X_1, Z_{1:2}, W_{1:2}} & f_{\text{NN}}(W_{1:2}, X_1, Z_{1:2}; X, Y) \\ & := \frac{1}{2} \|Y - h_2(Z_2)\|_F^2 + \frac{\alpha}{2} \|Z_2 - W_2 X_1\|_F^2 \\ & + \frac{\alpha}{2} \|X_1 - h_1(Z_1)\|_F^2 + \frac{\alpha}{2} \|Z_1 - W_1 X\|_F^2. \end{aligned} \quad (13.1)$$

We will apply block coordinate descent to solve the above optimisation problem. The idea is alternating minimisation, i.e., we minimise the objective function with respect to one of the variables $W_{1:2}$, X_1 , and $Z_{1:2}$, by fixing the others.

We also provide the closed form solution for the proximal operator of the Relu function. Let h denote the Relu function. Consider the following

one-dimensional minimisation problem:

$$x^\star = \arg \min_x \frac{1}{2} (h(x) - a)^2 + \frac{\gamma}{2} (x - b)^2.$$

The optimal solution is given by

$$\begin{aligned} & \text{prox}_{\frac{1}{2\gamma}(h(\cdot)-a)^2}(b) \\ &= \begin{cases} \frac{a+\gamma b}{1+\gamma}, & \text{if } a + \gamma b \geq 0, \ b \geq 0, \\ & \text{or if } -\left(\sqrt{\gamma(\gamma+1)} - \gamma\right)a \leq \gamma b < 0, \\ b, & \text{if } -a \leq \gamma b \leq -\left(\sqrt{\gamma(\gamma+1)} - \gamma\right)a < 0, \\ \min\{b, 0\}, & \text{if } a + \gamma b < 0. \end{cases} \end{aligned}$$

The data file provide the initial values for all variables and the values of the key parameters.

1. Write a function to update $W_{1:2}$ when all other variables are fixed.
Save the value of $W_{1:2}^1$ into the data file. [2]
2. Write a function to update X_1 when all other variables are fixed.
Save the values of X_1^1 into the data file. [1]
3. Write a function to update $Z_{1:2}$ when all other variables are fixed.
Save the values of $Z_{1:2}^1$ into the data file. [4]
4. Compute

$$g^1 \in \partial_{W_{1:2}, X_1, Z_{1:2}} f_{\text{NN}}(W_{1:2}^1, X_1^1, Z_{1:2}^1; X, Y).$$

[4]

5. Implement the BCD algorithm. Run your implementation until either $\|g^k\|_F \leq \epsilon \|\bar{X}^k\|_F$, where $\epsilon = 10^{-2}$ and \bar{X}^k includes all the optimisation variables at the end of k -th iteration, or the 100-th iteration has finished.

Show the following four curves:

- $f_{\text{NN}}(\bar{X}^k; X, Y)$ versus k ,
- $\|g^k\|_F / \|\bar{X}^k\|_F$ versus k ,
- $l(f(X; W_{1:2}), Y)$ versus k ,
- $l(f(X_{\text{test}}; W_{1:2}), Y_{\text{test}})$ versus k .

[4]