

HERIOT-WATT UNIVERSITY

MASTERS THESIS

Bayesian Reconsruction and Regression over Networks

Author:

John SMITH

Supervisor:

Dr. James SMITH

*A thesis submitted in fulfilment of the requirements
for the degree of MSc.*

in the

School of Mathematical and Computer Sciences

February 2023



Declaration of Authorship

I, John SMITH, declare that this thesis titled, 'Bayesian Reconsruction and Regression over Networks' and the work presented in it is my own. I confirm that this work submitted for assessment is my own and is expressed in my own words. Any uses made within it of the works of other authors in any form (e.g., ideas, equations, figures, text, tables, programs) are properly acknowledged at any point of their use. A list of the references employed is included.

Signed:

Date:

“Thanks to my solid academic training, today I can write hundreds of words on virtually any topic without possessing a shred of information, which is how I got a good job in journalism.”

Dave Barry

Abstract

The Thesis Abstract is written here (and usually kept to just this page).

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor :)

Contents

Declaration of Authorship	i
Abstract	iii
Acknowledgements	iv
Contents	v
List of Figures	viii
List of Tables	ix
Abbreviations	x
Symbols	xi
Identities	xiv
1 Introduction	1
1.1 Background and Definitions	1
1.2 Thesis overview	2
2 Outline and Fundamentals	3
2.1 Graph Signal Processing	3
2.1.1 A broad overview of the field	3
2.1.2 The graph Laplacian	3
2.1.3 Graph filters	3
2.2 Regression and Reconstruction	3
2.2.1 Graph Signal Reconstruction	3
2.2.2 Kernel Graph Regression	4
2.2.3 Regression with Network Cohesion	4
3 Kernel Generalized Least Squares Regression for Network Data	5
3.1 Kernel Graph Regression with Missing Values	5
3.2 GLS Kernel Graph Regression	5
3.2.1 A Gauss-Markov estimator	5

3.2.2	AR(1) processes	5
3.2.3	Experiments	5
4	Regression and Reconstruction on Cartesian Product Graphs	6
4.1	Graph Products	6
4.1.1	Basic definitions	7
4.1.2	The spectral properties of graph products	8
4.1.3	GSP with Cartesian product graphs	10
4.2	Graph Signal Reconstruction on Cartesian Product Graphs	13
4.2.1	Problem statement	15
4.2.2	A stationary iterative method	18
4.2.3	A conjugate gradient method	21
4.3	Convergence properties	24
4.3.1	Convergence of the SIM	24
4.3.2	Convergence of the CGM	27
4.4	Image processing experiments	28
4.5	Kernel Graph Regression with Unrestricted Missing Data Patterns	31
4.5.1	Cartesian product graphs and KGR	31
4.5.2	Convergence properties	31
4.6	Regression with Network Cohesion	31
4.6.1	Regression with node-level covariates	31
4.6.2	Convergence properties	31
5	Multi-Dimensional Cartesian Product Graphs	32
5.1	Fast computation with d -dimensional Kronecker products	32
5.2	Signal reconstruction	32
5.3	Kernel Graph Regression	32
5.4	Regression with Network Cohesion	32
6	Signal Uncertainty: Estimation and Sampling	33
6.1	Introduction	33
6.2	Posterior Estimation	33
6.2.1	Log-variance prediction	33
6.2.2	Estimation models	33
6.2.3	Query strategies	33
6.2.4	Comparison and analysis	33
6.3	Posterior Sampling	33
6.3.1	Perturbation optimization	33
6.4	Estimation vs Sampling	33
6.4.1	Experiments	33
7	Working with Binary-Valued Graph Signals	34
7.1	Logistic Graph Signal Reconstruction	34
7.2	Logistic Kernel Graph Regression	34
7.3	Logistic Regression with Network Cohesion	34
7.4	Approximate Sampling via the Laplace Approximation	34
8	Conclusions	35

8.1 Main Section 1	35
A Proofs	36

List of Figures

1.1	A graphical depiction of a graph signal. Here, the nodes are represented by circles, the edges as dotted lines, and the value of the signal at each node is represented by the height of its associated bar.	2
4.1	Graphical depiction of the standard graph products	9
4.2	A time-vertex Cartesian product graph	13
4.3	A time-vertex Cartesian product graph	14
4.4	The output from the first experiment is depicted. In the top left quadrant, the input images are shown across a range of noise levels and missing pixel percentages, with missing pixels chosen uniformly at random. Below that, in the lower left quadrant, the corresponding reconstructed images are shown. The right half of the plot is the same, except here entire columns and rows of pixels are removed at random.	28
4.5	The total runtime in seconds for the SIM and CGM compared to a naive Gaussian elimination approach is shown as a function of the total number of nodes using a quad-core Intel i7-7700HQ CPU.	29
4.6	The number of iterations required to reach a certain level of precision is shown experimentally, along with the theoretical upper bound, for the SIM and CGM	29

List of Tables

2.1	Isotropic graph filter functions	3
4.1	The adjacency and Laplacian matrices for the standard graph products .	8
4.2	Spectral decomposition of product graphs	10
4.3	Anisotropic graph filter functions	12

Abbreviations

GSP	Graph Signal Processing
GFT	Graph Fourier Transform
IGFT	Inverse Graph Fourier Transform
GSR	Graph Signal Reconstruction
KGR	Kernel Graph Regression
RNC	Regression with Network Cohesion
GLS	Generalised Least Squares
DCT	Discrete Cosine Transform
FCT	Fast Cosine Transform
FFT	Fast Fourier Transform
PSD	Positive Semi-Definite
SIM	Stationary Iterative Method
CGM	Conjugate Gradient Method
SNR	Signal to Noise Ratio

Symbols

Unless otherwise specified, the following naming conventions apply.

Integer constants

N	The number of nodes in a graph
T	The number of time points considered
M	The number of explanatory variables
Q	The number of queries

Integer variables

n	The index of a specific node in a graph
t	The index of a specific time point
m	The index of a specific explanatory variable
q	The index of a specific query
i, j, k	Generic indexing variables

Scalar variables

α	An autocorrelation regularisation parameter
β	A hyperparameter characterising a graph filter
γ	A precision parameter
λ	An eigenvalue <i>or</i> ridge regression penalty parameter
μ	The mean of a random variable
θ	AR(1) autocorrelation parameter
σ^2	The variance of a random variable

Matrices

A	The graph adjacency matrix
D	A diagonal matrix
E	The prediction residuals
F	A predicted graph signal
G	A spectral scaling matrix
H	A graph filter <i>or</i> Hessian matrix
I_N	The $(N \times N)$ identity matrix
J_N	An $(N \times N)$ matrix of ones
K	A kernel (Gram) matrix
L	The graph Laplacian
S	A binary selection matrix
U	Laplacian eigenvector matrix
V	Kernel eigenvector matrix
X	Data matrix of explanatory variables
Y	(Partially) observed graph signal
Λ	A diagonal eigenvalue matrix
Σ	A covariance matrix
Φ, Ψ	Generic eigenvector matrices
Ω	Log marginal variance matrix

Vectors/tensors

$\mathbf{1}_N$	A length- N vector of ones
\mathbf{e}	The prediction residuals
\mathbf{e}_i	The i -th unit basis vector
\mathbf{f}	The predicted graph signal
\mathbf{s}	A binary selection vector/tensor
\mathbf{x}	A vector of explanatory variables
\mathbf{y}	The observed graph signal
$\boldsymbol{\alpha}$	A flexible intercept vector/tensor
$\boldsymbol{\beta}$	A graph filter parameter vector <i>or</i> vector of regression coefficients
$\boldsymbol{\theta}$	A aggregated coefficient vector $[\boldsymbol{\alpha}^\top, \boldsymbol{\beta}^\top]^\top$

Functions

$g(\cdot)$	A graph filter function
$p(\text{statement})$	The probability that a statement is true
$\pi(\cdot)$	A probability density function
$\xi(\cdot)$	Optimisation target function
$\kappa(\cdot, \cdot)$	A kernel function

Operations

$(\cdot)^\top$	Transpose of a matrix/vector
$\ \cdot\ _F$	The Frobenius norm
$\text{tr}(\cdot)$	The trace of a square matrix
$\text{vec}(\cdot)$	Convert a matrix to a vector in column-major order
$\text{vec}_{\text{RM}}(\cdot)$	Convert a matrix to a vector in row-major order
$\text{mat}(\cdot)$	Convert a vector to a matrix in column-major order
$\text{mat}_{\text{RM}}(\cdot)$	Convert a vector to a matrix in row-major order
$\text{diag}(\cdot)$	Convert a vector to a diagonal matrix
$\text{diag}^{-1}(\cdot)$	Convert the diagonal of a matrix into a vector
\otimes	The Kronecker product
\oplus	The Kronecker sum
\circ	The Hadamard product

Graphs

\mathcal{G}	A graph
\mathcal{V}	A vertex/node set
\mathcal{E}	An edge set

Miscellaneous

$\hat{(\cdot)}$	The estimator of a matrix/vector/tensor
$O(\cdot)$	The runtime complexity
x_i	A vector element
\mathbf{X}_i	A matrix column
\mathbf{X}_{ij}	A matrix element

Identities

1	$\text{vec}(\mathbf{AXB})$	$(\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X})$
2	$\text{tr}(\mathbf{A}^\top \mathbf{B})$	$\text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$
3	$\mathbf{AC} \otimes \mathbf{BD}$	$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D})$
4	$(\mathbf{A} \otimes \mathbf{B})^{-1}$	$\mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$
5	$\text{tr}(\mathbf{X}^\top \mathbf{A} \mathbf{Y} \mathbf{B})$	$\text{vec}(\mathbf{X})^\top (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{Y})$
6	$\text{vec}(\mathbf{J} \circ \mathbf{Y})$	$\text{diag}(\text{vec}(\mathbf{J})) \text{vec}(\mathbf{Y})$
9	$\text{diag}^{-1}(\mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B})$	$(\mathbf{B}^\top \circ \mathbf{A}) \mathbf{x}$

For/Dedicated to/To my...

Chapter 1

Introduction

1.1 Background and Definitions

Graph Signal Processing (GSP) is a rapidly evolving field that sits at the intersection between spectral graph theory, statistics and data science [Shuman et al., 2013]. In this context, a graph is an abstract collection of objects in which any pair may be, in some sense, “related”. These objects are referred to as vertices (or nodes) and their connections as edges [Newman, 2018]. GSP is concerned with the mathematical analysis of signals that are defined over the nodes of a graph, referred to simply as *graph signals*.

A graph signal can be thought of as a value that is measured simultaneously at each node in a graph. In practice, it is represented as a vector where each element corresponds to a single node. For example, consider a social network where each node represents an individual and presence of an edge between two nodes indicates that the two individuals have met. An example of a graph signal in this context could be the age of each person in the network. Figure 1.1 shows a graphical depiction of a signal defined over a network.

Graphs and graph signals have proven a useful way to describe data across a broad range of applications owing to their flexibility and relative simplicity. They are able to summarise the of properties large, complex systems within a single easily-digestible structure. Much of the data

The GSP community, in particular, is focused on generalising tools designed for traditional signal processing tasks to irregular graph-structured domains.

[Ortega et al., 2018]

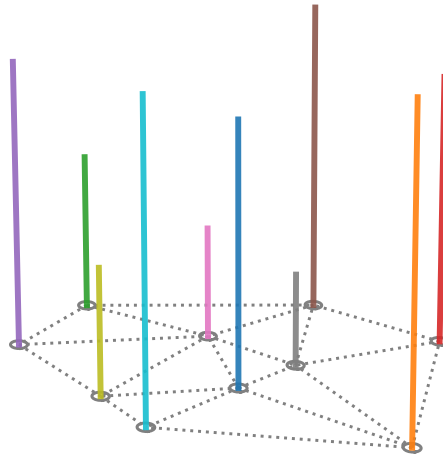


FIGURE 1.1: A graphical depiction of a graph signal. Here, the nodes are represented by circles, the edges as dotted lines, and the value of the signal at each node is represented by the height of its associated bar.

1.2 Thesis overview

Chapter 2

Outline and Fundamentals

2.1 Graph Signal Processing

2.1.1 A broad overview of the field

2.1.2 The graph Laplacian

2.1.3 Graph filters

2.2 Regression and Reconstruction

2.2.1 Graph Signal Reconstruction

Introduce the known work on GSR

Filter	$g(\lambda; \beta)$
1-hop random walk	$(1 + \beta\lambda)^{-1}$
Diffusion	$\exp(-\beta\lambda)$
ReLu	$\max(1 - \beta\lambda, 0)$
Sigmoid	$2(1 + \exp(\beta\lambda))^{-1}$
Bandlimited	1, if $\beta\lambda \leq 1$ else 0

TABLE 2.1: Isotropic graph filter functions

2.2.2 Kernel Graph Regression

Introduce the known work on KGR and GPoG

2.2.3 Regression with Network Cohesion

Introduce the known work on RNC

Chapter 3

Kernel Generalized Least Squares Regression for Network Data

3.1 Kernel Graph Regression with Missing Values

3.2 GLS Kernel Graph Regression

3.2.1 A Gauss-Markov estimator

3.2.2 AR(1) processes

3.2.3 Experiments

Chapter 4

Regression and Reconstruction on Cartesian Product Graphs

In this chapter we explore a class of reconstruction algorithms as applied to signals defined on the nodes of a Cartesian product graph. In particular, we pose the reconstruction task in terms of Bayesian inference of an underlying signal given a noisy partial observation, and investigate scalable methods for obtaining the posterior mean. We begin in section 4.1 by reviewing the concept of a graph product, and explain why we choose to look specifically at the Cartesian product. We also review how concepts from standard one-dimensional graph signal processing such as the GFT and spectral filtering can be extended to the two dimensional case. In section 4.2 we introduce the statistical model defining GSR on a Cartesian product graph, and derive two alternative methods for solving for the posterior mean. These comprise a Stationary Iterative Method (SIM) and a Conjugate Gradient Method (CGM). In each case, we show how graph spectral considerations can be leveraged to find efficient forms of said algorithms and derive the convergence properties. In section 4.4 we run several experiments on image data (which can be understood as a special case of a Cartesian product graph signal) to verify these properties.

The key outcomes of this chapter are as follows:

-

4.1 Graph Products

In this chapter, we will be primarily concerned with signal processing on *Cartesian product graphs*. This special class of graph finds applications in numerous areas, such

as video, hyper-spectral image processing and network time series problems. However, the Cartesian product is not the only way to consistently define a product between two graphs. In this section we formally introduce the concept of a graph product, examine some prominent examples, and explain why we choose to look specifically at the Cartesian graph product.

4.1.1 Basic definitions

In the general case, consider two undirected graphs $\mathcal{G}_A = (\mathcal{V}_A, \mathcal{E}_A)$ and $\mathcal{G}_B = (\mathcal{V}_B, \mathcal{E}_B)$ with vertex sets given by $\mathcal{V}_A = \{a \in \mathbb{N} \mid a \leq A\}$ and $\mathcal{V}_B = \{b \in \mathbb{N} \mid b \leq B\}$ respectively. (In this context we do not regard zero to be an element of the natural numbers). A new graph \mathcal{G} can be constructed by taking the product between \mathcal{G}_A and \mathcal{G}_B . This can be generically written as follows.

$$\mathcal{G} = \mathcal{G}_A \diamond \mathcal{G}_B = (\mathcal{V}, \mathcal{E}) \quad (4.1)$$

For all definitions of a graph product, the new vertex set \mathcal{V} is given by the Cartesian product of the vertex sets of the factor graphs, that is

$$\mathcal{V} = \mathcal{V}_A \times \mathcal{V}_B = \{(a, b) \in \mathbb{N}^2 \mid a \leq A \text{ and } b \leq B\} \quad (4.2)$$

Typically, vertices are arranged in lexicographic order, in the sense that $(a, b) \leq (a', b')$ iff $a < a'$ or $(a = a' \text{ and } b \leq b')$ [Harzheim, 2005]. Each consistent rule for constructing the new edge set \mathcal{E} corresponds to a different definition of a graph product. In general, there are eight possible conditions for deciding whether two nodes (a, b) and (a', b') are to be connected in the new graph.

1. $[a, a'] \in \mathcal{E}_A$ and $b = b'$
2. $[a, a'] \notin \mathcal{E}_A$ and $b = b'$
3. $[a, a'] \in \mathcal{E}_A$ and $[b, b'] \in \mathcal{E}_B$
4. $[a, a'] \notin \mathcal{E}_A$ and $[b, b'] \in \mathcal{E}_B$
5. $[a, a'] \in \mathcal{E}_A$ and $[b, b'] \notin \mathcal{E}_B$
6. $[a, a'] \notin \mathcal{E}_A$ and $[b, b'] \notin \mathcal{E}_B$
7. $a = a'$ and $[b, b'] \in \mathcal{E}_B$,
8. $a = a'$ and $[b, b'] \notin \mathcal{E}_B$

Each definition of a graph product corresponds to the union of a specific subset of these conditions, thus, there exist 256 different types of graph product [Barik et al., 2015]. Of these, the Cartesian product (conditions 1 or 7), the direct product (condition 3), the strong product (conditions 1, 3 or 7) and the lexicographic product (conditions 1, 3, 5 or 7) are referred to as the standard products and are well-studied [Imrich and Klavžar, 2000]. A graphical depiction of the standard graph products is shown in figure 4.1. In each of these four cases, the adjacency and Laplacian matrices of the product graph can be described in terms of matrices relating to the factor graphs [Barik et al., 2018, Fiedler, 1973]. This is shown in table 4.1.

	Adjacency matrix	Laplacian
Cartesian	$\mathbf{A}_A \oplus \mathbf{A}_B$	$\mathbf{L}_A \oplus \mathbf{L}_B$
Direct	$\mathbf{A}_A \otimes \mathbf{A}_B$	$\mathbf{D}_A \otimes \mathbf{L}_B + \mathbf{L}_A \otimes \mathbf{D}_B - \mathbf{L}_A \otimes \mathbf{L}_B$
Strong	$\mathbf{A}_A \otimes \mathbf{A}_B + \mathbf{A}_A \oplus \mathbf{A}_B$	$\mathbf{D}_A \otimes \mathbf{L}_B + \mathbf{L}_A \otimes \mathbf{D}_B - \mathbf{L}_A \otimes \mathbf{L}_B + \mathbf{L}_A \oplus \mathbf{L}_B$
Lexicographic	$\mathbf{I}_A \otimes \mathbf{A}_B + \mathbf{A}_A \otimes \mathbf{J}_A$	$\mathbf{I}_A \otimes \mathbf{L}_B + \mathbf{L}_A \otimes \mathbf{J}_B + \mathbf{D}_A \otimes (\mathcal{V}_B \mathbf{I}_B - \mathbf{J}_B)$

TABLE 4.1: The adjacency and Laplacian matrices for the standard graph products. Here, \mathbf{D}_A and \mathbf{D}_B are the diagonal degree matrices, i.e $\mathbf{D}_A = \text{diag}(\mathbf{A}_A \mathbf{1})$. \mathbf{I}_A and \mathbf{J}_A are the $(A \times A)$ identity matrix and matrix of ones respectively.

Given these definitions, it may seem that all the standard graph products are non-commutative in the sense that $\mathbf{A}_A \oplus \mathbf{A}_B \neq \mathbf{A}_B \oplus \mathbf{A}_A$ etc. However, the graphs $\mathcal{G}_A \diamond \mathcal{G}_B$ and $\mathcal{G}_B \diamond \mathcal{G}_A$ are in fact isomorphically identical in the case of the Cartesian, direct and strong products. This is not the case for the Lexicographic product [Imrich and Klavžar, 2000].

4.1.2 The spectral properties of graph products

In the field of graph signal processing, we are often concerned with analysing the properties of graphs via eigendecomposition of the graph Laplacian [Mieghem, 2010]. In the case of product graphs, it is greatly preferable if we are able to fully describe the spectrum of $\mathcal{G}_A \diamond \mathcal{G}_B$ in terms of the spectra of \mathcal{G}_A and \mathcal{G}_B alone. This is because direct decomposition of a dense \mathbf{L} has time-complexity $O(A^3 B^3)$, whereas decomposition of the factor Laplacians individually has complexity $O(A^3 + B^3)$. As the graphs under considerations become medium to large, this fact quickly makes direct decomposition of the product graph Laplacian intractable. However, in the general case, only the spectra of the Cartesian and lexicographic graph products can be described in this way [Barik et al., 2018]. In the case of the direct and strong product, it is possible to estimate

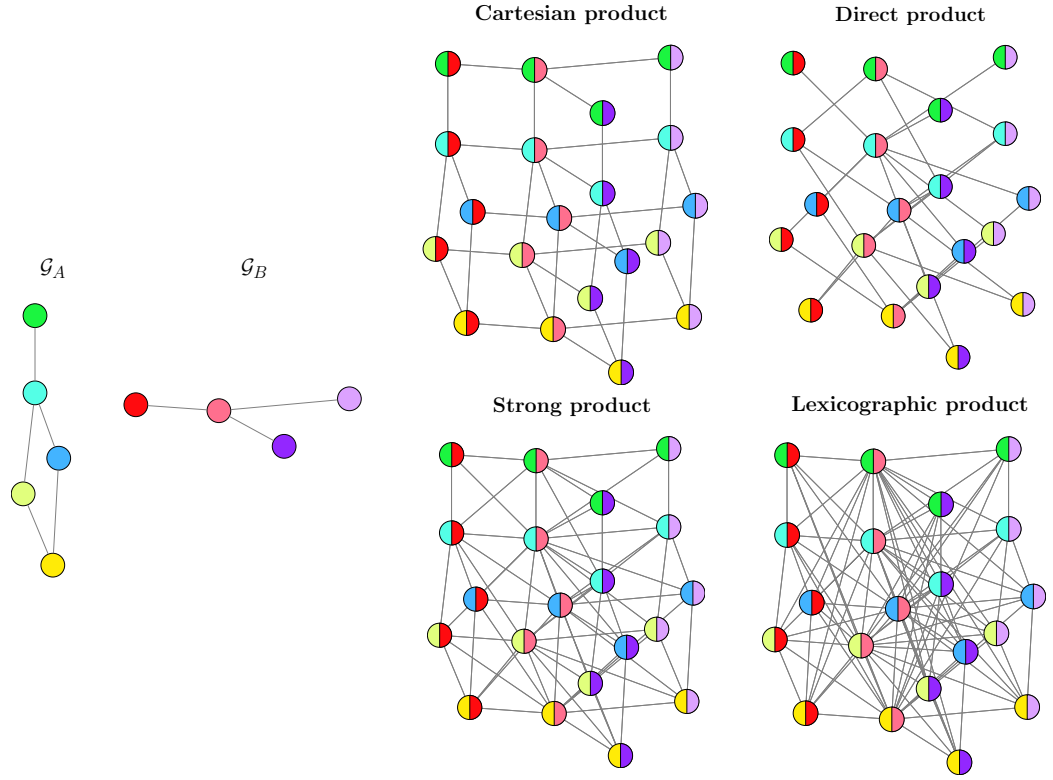


FIGURE 4.1: A graphical depiction of the four standard graph products

the spectra without performing the full decomposition (see [Sayama, 2016]). However, in general, the full eigendecomposition of the product graph Laplacian can only be described in terms of the factor eigendecompositions when both factor graphs are regular.

Consider the eigendecompositions of \mathbf{L}_A and \mathbf{L}_B .

$$\mathbf{L}_A = \mathbf{U}_A \mathbf{\Lambda}_A \mathbf{U}_A^\top, \quad \text{and} \quad \mathbf{L}_B = \mathbf{U}_B \mathbf{\Lambda}_B \mathbf{U}_B^\top \quad (4.3)$$

where \mathbf{U}_A and \mathbf{U}_B are the respective orthonormal eigenvector matrices, and $\mathbf{\Lambda}_A$ and $\mathbf{\Lambda}_B$ are the diagonal eigenvalue matrices given by

$$\mathbf{\Lambda}_A = \begin{bmatrix} \lambda_1^{(A)} & & & \\ & \lambda_2^{(A)} & & \\ & & \ddots & \\ & & & \lambda_A^{(A)} \end{bmatrix} \quad \text{and} \quad \mathbf{\Lambda}_B = \begin{bmatrix} \lambda_1^{(B)} & & & \\ & \lambda_2^{(B)} & & \\ & & \ddots & \\ & & & \lambda_B^{(B)} \end{bmatrix} \quad (4.4)$$

Given these definitions, table 4.2 gives information about the spectral decomposition of the standard graph products.

	Eigenvalues	Eigenvectors
Cartesian	$\lambda_a^{(A)} + \lambda_b^{(B)}$	$(\mathbf{U}_A)_a \otimes (\mathbf{U}_B)_b$
Direct [*]	$r_A \lambda_b^{(B)} + r_B \lambda_a^{(A)} - \lambda_a^{(A)} \lambda_b^{(B)}$	$(\mathbf{U}_A)_a \otimes (\mathbf{U}_B)_b$
Strong [*]	$(1 + r_A) \lambda_b^{(B)} + (1 + r_B) \lambda_a^{(A)} - \lambda_a^{(A)} \lambda_b^{(B)}$	$(\mathbf{U}_A)_a \otimes (\mathbf{U}_B)_b$
Lexicographic [†]	$B \lambda_a^{(A)}$	$(\mathbf{U}_A)_a \otimes \mathbf{1}_B$
	$\lambda_b^{(B)} + B \deg(a)$	$\mathbf{e}_a \otimes (\mathbf{U}_B)_b$

TABLE 4.2: Eigendecomposition of the Laplacian of the standard graph products. Here, a and b are understood to run from 1 to A and 1 to B respectively. ^{*} only for r_A and r_B -regular factor graphs. [†] note that the b runs from 2 to B in the lower row.

4.1.3 GSP with Cartesian product graphs

While both the direct and strong products do find uses in certain applications (for example, see [Kaveh and Alinejad, 2011]), they are both less common and more challenging to work with in a graph signal processing context due to their spectral properties described in the previous subsection. In practice, being limited to regular factor graphs means the majority of practical GSP applications are ruled out. The lexicographic product does not share this drawback, however it is also significantly less common than the Cartesian product in real-world applications. For this reason, in the following, we focus primarily on the Cartesian product.

Given the spectral decomposition of the Cartesian graph product stated in table 4.2, we can write the Laplacian eigendecomposition in matrix form as follows.

$$\mathbf{L} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top, \quad \text{where} \quad \mathbf{U} = \mathbf{U}_A \otimes \mathbf{U}_B \quad \text{and} \quad \mathbf{\Lambda} = \mathbf{\Lambda}_A \oplus \mathbf{\Lambda}_B \quad (4.5)$$

This motivates the following definitions for the Graph Fourier Transform (GFT) and its inverse (IGFT). Consider a signal defined over the nodes of a Cartesian product graph expressed as a matrix $\mathbf{Y} \in \mathbb{R}^{B \times A}$. We can perform the GFT as follows.

$$\text{GFT}(\mathbf{Y}) = \text{mat}\left((\mathbf{U}_A^\top \otimes \mathbf{U}_B^\top) \text{vec}(\mathbf{Y})\right) = \mathbf{U}_B^\top \mathbf{Y} \mathbf{U}_A \quad (4.6)$$

Correspondingly, we can define the IGFT acting on a matrix of spectral components $\mathbf{Z} \in \mathbb{R}^{B \times A}$ as follows.

$$\text{IGFT}(\mathbf{Z}) = \text{mat}\left((\mathbf{U}_A \otimes \mathbf{U}_B) \text{vec}(\mathbf{Z})\right) = \mathbf{U}_B \mathbf{Z} \mathbf{U}_A^\top \quad (4.7)$$

Product graph signals: representation and vectorisation

It is natural to assume that signals defined on the nodes of a Cartesian product graph $\mathcal{G}_A \square \mathcal{G}_B$ could be represented by matrices (order two tensors) of shape $(A \times B)$. Since product graph operators, such as the Laplacian $\mathbf{L}_A \oplus \mathbf{L}_B$, act on vectors of length AB , we must define a consistent function to map matrix graph signals $\in \mathbb{R}^{A \times B}$ to vector graph signals $\in \mathbb{R}^{AB}$. The standard mathematical operator for this purpose is the $\text{vec}(\cdot)$ function, along with its reverse operator $\text{mat}(\cdot)$. However, this is somewhat problematic since $\text{vec}(\cdot)$ is defined to act in *column-major* order, that is

$$\text{vec} \left(\begin{bmatrix} \mathbf{Y}_{(1,1)} & \mathbf{Y}_{(1,2)} & \cdots & \mathbf{Y}_{(1,B)} \\ \mathbf{Y}_{(2,1)} & \mathbf{Y}_{(2,2)} & \cdots & \mathbf{Y}_{(2,B)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Y}_{(A,1)} & \mathbf{Y}_{(A,2)} & \cdots & \mathbf{Y}_{(A,B)} \end{bmatrix} \right) = \begin{bmatrix} \mathbf{Y}_{(1,1)} \\ \mathbf{Y}_{(2,1)} \\ \vdots \\ \mathbf{Y}_{(A-1,1)} \\ \mathbf{Y}_{(A,1)} \\ \mathbf{Y}_{(1,2)} \\ \mathbf{Y}_{(2,2)} \\ \vdots \\ \mathbf{Y}_{(A-1,2)} \\ \mathbf{Y}_{(A,2)} \\ \vdots \\ \mathbf{Y}_{(1,B)} \\ \mathbf{Y}_{(2,B)} \\ \vdots \\ \mathbf{Y}_{(A-1,B)} \\ \mathbf{Y}_{(A,B)} \end{bmatrix}$$

As is visible, this does not result in a lexicographic ordering of the matrix elements when the graph signal has shape $(A \times B)$. Therefore, to avoid this issue and to be consistent with standard mathematical notation, we will assume that graph signals are represented by matrices of shape $(B \times A)$ when considering the product between two graphs $\mathcal{G}_A \square \mathcal{G}_B$. For graph signals of this shape, the first index represents traversal of the nodes in \mathcal{G}_B , and the second index represents traversal of the nodes in \mathcal{G}_A . This ensures that matrix elements are correctly mapped to vector elements when using the column-major $\text{vec}(\cdot)$ function.

Given these definitions, we can define a spectral operator (usually a low-pass filter) \mathbf{H} which acts on graph signals according to a spectral scaling function $g(\lambda; \beta)$ such as one of those defined in table 2.1. As with regular non-product graphs, the action of this operator can be understood as first transforming a signal into the frequency domain via the GFT, then scaling the spectral components according to some function, and finally transforming back into the vertex domain via the IGFT.

Filter	$g(\boldsymbol{\lambda}; \boldsymbol{\beta})$
1-hop random walk	$(1 + \boldsymbol{\beta}^\top \boldsymbol{\lambda})^{-1}$
Diffusion	$\exp(-\boldsymbol{\beta}^\top \boldsymbol{\lambda})$
ReLU	$\max(1 - \boldsymbol{\beta}^\top \boldsymbol{\lambda}, 0)$
Sigmoid	$2(1 + \exp(\boldsymbol{\beta}^\top \boldsymbol{\lambda}))^{-1}$
Bandlimited	1, if $\boldsymbol{\beta}^\top \boldsymbol{\lambda} \leq 1$ else 0

TABLE 4.3: Anisotropic graph filter functions

$$\begin{aligned}
\mathbf{H} &= g(\mathbf{L}_A \oplus \mathbf{L}_B) \\
&= (\mathbf{U}_A \otimes \mathbf{U}_B) g(\boldsymbol{\Lambda}_A \oplus \boldsymbol{\Lambda}_B) (\mathbf{U}_A^\top \otimes \mathbf{U}_B^\top) \\
&= (\mathbf{U}_A \otimes \mathbf{U}_B) \text{diag}(\text{vec}(\mathbf{G})) (\mathbf{U}_A^\top \otimes \mathbf{U}_B^\top)
\end{aligned} \tag{4.8}$$

The matrix $\mathbf{G} \in \mathbb{R}^{B \times A}$, which we refer to as the spectral scaling matrix, holds the value of the scaling function applied to the sum of pairs of eigenvalues, such that

$$\mathbf{G}_{ba} = g(\lambda_a^{(A)} + \lambda_b^{(B)}; \beta) \tag{4.9}$$

We observe that defining the filtering operation in this manner implies that the intensity is equal across both \mathcal{G}_A and \mathcal{G}_B . We refer to filters of this type as *isotropic*. This can be further generalised by considering an *anisotropic* graph filter, which offers independent control over the filter intensity in each of the two dimensions. In this case, we define \mathbf{G} as follows.

$$\mathbf{G}_{ba} = g\left(\begin{bmatrix} \lambda_a^{(A)} \\ \lambda_b^{(B)} \end{bmatrix}, \begin{bmatrix} \beta_a \\ \beta_b \end{bmatrix}\right) \tag{4.10}$$

where now $g(\boldsymbol{\lambda}; \boldsymbol{\beta})$ is chosen to be an anisotropic graph filter such as one of those listed in table 4.3. Note that the original parameter β is now replaced by a vector of parameters $\boldsymbol{\beta}$ which control the filter intensity in each dimension.

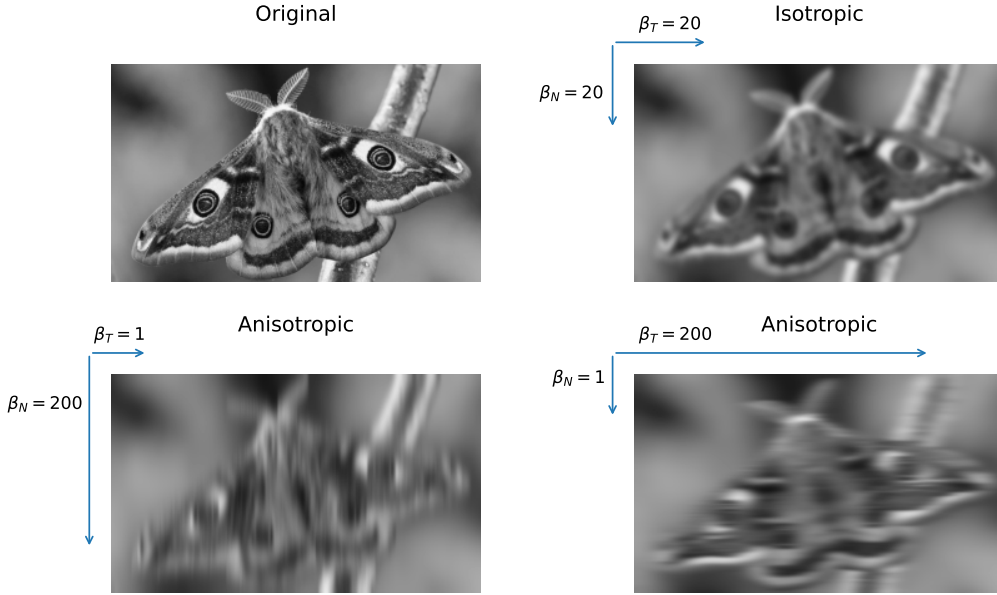


FIGURE 4.2: A graphical depiction of a time-vertex Cartesian product graph

4.2 Graph Signal Reconstruction on Cartesian Product Graphs

We now turn our attention to the task of signal reconstruction on Cartesian product graphs. In the following, we will replace the factor graph labels A and B with T and N respectively. The reason for this is that one application of particular interest is graph time-series problems, where we seek to model a network of N nodes across a series of T discrete time points. These so called “time-vertex” (T-V) problems have garnered significant interest recently in the context of GSP [Grassi et al., 2018, Isufi et al., 2017, Loukas and Foucard, 2016]. T-V signals can be understood as existing on the nodes of a Cartesian product graph $\mathcal{G}_T \square \mathcal{G}_N$. In particular, we can conceptualise T repeated measurements of a signal defined across the nodes of a N -node graph as a single measurement of a signal defined on the nodes of $\mathcal{G}_T \square \mathcal{G}_N$, where \mathcal{G}_T is a simple path graph.

On the Laplacian spectrum of the path graph

When considering time-vertex problems with uniformly spaced time intervals, \mathcal{G}_T will be described by a path graph with equal weights on each edge. This special case of a graph has vertices given by $\mathcal{V}_T = \{t \in \mathbb{N} | t \leq T\}$ and edges given by $\mathcal{E}_T = \{[t, t+1] | t < T\}$. The Laplacian matrix of the path graph is therefore given by

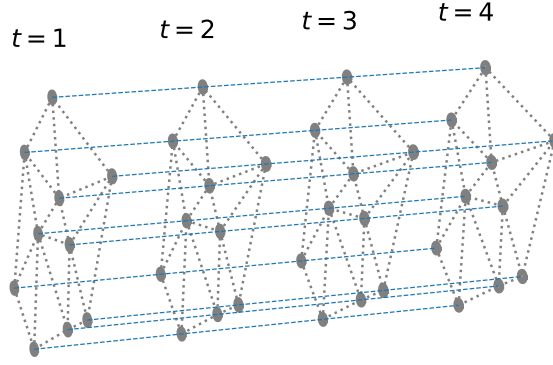


FIGURE 4.3: A graphical depiction of a time-vertex Cartesian product graph

$$\mathbf{L}_T = \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & & \ddots & & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 1 \end{bmatrix}$$

The eigenvalues and eigenvectors of this Laplacian are well-known and can be expressed in closed-form [Jiang, 2012]. In particular,

$$\lambda_t^{(T)} = 2 - 2 \cos \left(\frac{t-1}{T} \pi \right)$$

and

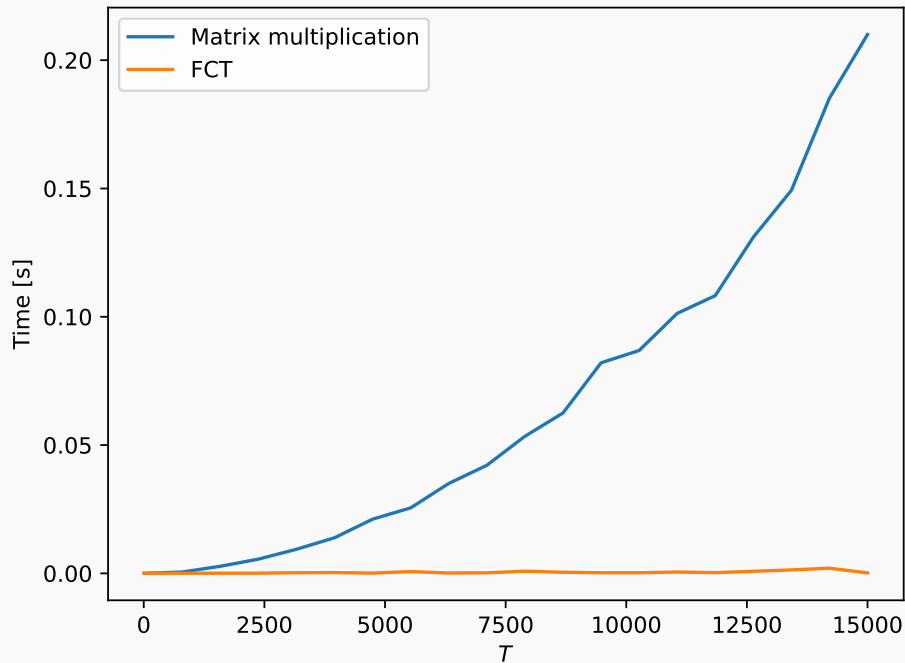
$$(\mathbf{U}_T)_{ij} = \cos \left(\frac{j-1}{T} \left(i - \frac{1}{2} \right) \pi \right)$$

where the columns of \mathbf{U} are appropriately normalised such that the magnitude of each eigenvector is one. Furthermore, this implies that the graph Fourier transform of a signal $\mathbf{y} \in \mathbb{R}^T$ is given by the orthogonal type-II Discrete Cosine Transform (DCT) [Ahmed et al., 1974]. This is of significance, as it means we can leverage Fast Cosine Transform (FCT) algorithms [Makhoul, 1980] which operate in a similar manner to the well-known Fast Fourier Transform (FFT) [Cooley and Tukey, 1965]. See chapter 4 of Rao and Yip [1990] for an overview of FCT algorithms.

In particular, this reduces both of the following procedures

$$\text{GFT}(\mathbf{y}) = \mathbf{U}^\top \mathbf{y} \quad \text{and} \quad \text{IGFT}(\mathbf{y}) = \mathbf{U} \mathbf{y}$$

from $O(T^2)$ operations to $O(T \log T)$ operations, which can be significant for large time-series problems. The figure below compares the time to compute the graph Fourier transform of a random signal using the matrix multiplication method vs the FCT implementation. In particular, we varied T from 10 to 15,000 in 20 equally spaced increments, and measured the mean time to compute $\mathbf{U}^\top \mathbf{y}$ across five independent trials using both the standard matrix multiplication and the Fast Cosine Transform method. As is visible, the difference becomes extremely pronounced as T grows large.



Note that, despite the observation that \mathcal{G}_T is often a path graph in the context of T-V problems, the methods introduced in this section are valid for the Cartesian product between arbitrary undirected factor graphs.

4.2.1 Problem statement

The goal of Graph Signal Reconstruction (GSR) is to estimate the value of a partially observed graph signal at nodes where no data was collected. In the context of GSR on a Cartesian product graph, the available data is an observed signal $\mathbf{Y} \in \mathbb{R}^{N \times T}$ where only a partial set $\mathcal{S} = \{(n_1, t_1), (n_2, t_2), \dots\}$ of the signal elements were recorded. All other missing elements of \mathbf{Y} are set to zero. Our model is based on the assumption that \mathbf{Y} is

a noisy partial observation of an underlying signal $\mathbf{F} \in \mathbb{R}^{N \times T}$, which is itself assumed to be smooth with respect to the graph topology.

We define the statistical model for the generation of an observation matrix \mathbf{Y} as

$$\mathbf{Y} = \mathbf{S} \circ (\mathbf{F} + \mathbf{E}) \quad (4.11)$$

where $\mathbf{S} \in [0, 1]^{N \times T}$ is referred to as the sensing matrix, and has entries given by

$$\mathbf{S}_{nt} = \begin{cases} 1 & \text{if } (n, t) \in \mathcal{S} \\ 0 & \text{otherwise} \end{cases} \quad (4.12)$$

The matrix \mathbf{E} represents the model error and is assumed to have an independent normal distribution with unit variance. Therefore, the probability distribution of \mathbf{Y} given the latent signal \mathbf{F} is

$$\text{vec}(\mathbf{Y}) | \mathbf{F} \sim \mathcal{N}(\text{vec}(\mathbf{S} \circ \mathbf{F}), \text{diag}(\text{vec}(\mathbf{S}))) \quad (4.13)$$

Note that the covariance matrix $\text{diag}(\text{vec}(\mathbf{S}))$ is semi-positive definite by construction. This naturally reflects the constraint that some elements of \mathbf{Y} are zero with probability 1. In order to estimate the latent signal \mathbf{F} , we must provide a prior distribution describing our belief about its likely profile ahead of time. In general, we expect \mathbf{F} to be smooth with respect to the topology of the graph. This can be expressed by setting the covariance matrix in its prior to be proportional to \mathbf{H}^2 , where \mathbf{H} is a graph filter as defined in equation (4.8). For now, in the absence of any further information, we assume that the prior mean for \mathbf{F} is zero across all elements.

$$\text{vec}(\mathbf{F}) \sim \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{H}^2) \quad (4.14)$$

Next, given an observation \mathbf{Y} , we use Bayes' rule to find the posterior distribution over \mathbf{F} . This is given by

$$\pi(\text{vec}(\mathbf{F}) | \mathbf{Y}) = \frac{\pi(\text{vec}(\mathbf{Y}) | \mathbf{F}) \pi(\mathbf{F})}{\pi(\mathbf{Y})}. \quad (4.15)$$

where we use the notation $\pi(\cdot)$ to denote a probability density function.

The posterior distribution for \mathbf{F} is given by

$$\text{vec}(\mathbf{F}) \mid \mathbf{Y} \sim \mathcal{N}(\boldsymbol{\Sigma} \text{vec}(\mathbf{Y}), \boldsymbol{\Sigma}) \quad (4.16)$$

where

$$\boldsymbol{\Sigma} = \left(\text{diag}(\text{vec}(\mathbf{S})) + \gamma \mathbf{H}^{-2} \right)^{-1} \quad (4.17)$$

A proof of this can be found in the appendix, theorem [A.1](#).

In this chapter, we are primarily interested in computing the posterior mean, which is the solution to the following linear system.

$$\text{vec}(\mathbf{F}) = \left(\text{diag}(\text{vec}(\mathbf{S})) + \gamma \mathbf{H}^{-2} \right)^{-1} \text{vec}(\mathbf{Y}) \quad (4.18)$$

We return to the question of sampling from the posterior and estimating the posterior covariance directly in chapter [6](#).

Two significant computational challenges arise when working with non-trivial graph signal reconstruction problems, where the number of vertices in the product graph is large. First, although the posterior mean point estimator given in eq. (4.18) has an exact closed-form solution, its evaluation requires solving an $NT \times NT$ system of equations, which is impractical for all but the smallest of problems. Second, since the eigenvalues of \mathbf{H} can be close to or exactly zero, \mathbf{H}^{-2} may be severely ill-conditioned and even undefined. This means the condition number of the coefficient matrix may not be finite, making basic iterative methods to numerically solve the linear system, such as steepest descent, slow or impossible. The models proposed in this paper aim to overcome these problems.

Since the coefficient matrix defining the system is of size $NT \times NT$, direct methods such as Gaussian elimination are assumed to be out of the question. In such cases, one often resorts to one of three possible solution approaches: stationary iterative methods; Krylov methods; and multigrid methods. Each are part of the family of iterative methods which are most commonly found in applications of sparse matrices, such as finite element methods [[Brenner et al., 2008](#)]. In the following, we propose a stationary iterative method and a Krylov method and compare their relative behaviour. In both cases, we show that each step of the iterative process can be completed in $O(N^2T + NT^2)$ operations, making a solution feasible. First, we present each of the methods in isolation. Then, the convergence behaviour of each is derived theoretically and verified numerically.

4.2.2 A stationary iterative method

In this section, we demonstrate a technique for obtaining the posterior mean by adopting a classic approach to solving linear systems, known as *matrix splitting*, which sits within the family of Stationary Iterative Methods (SIMs) [Saad, 2003]. The general splitting strategy is to break the coefficient matrix into the form $\mathbf{M} - \mathbf{N}$, such that

$$\text{vec}(\mathbf{F}) = (\mathbf{M} - \mathbf{N})^{-1} \text{vec}(\mathbf{Y}) \quad (4.19)$$

By noting that

$$\mathbf{M} \text{vec}(\mathbf{F}) = \mathbf{N} \text{vec}(\mathbf{F}) + \text{vec}(\mathbf{Y}) \quad (4.20)$$

$$\text{vec}(\mathbf{F}) = \mathbf{M}^{-1} \mathbf{N} \text{vec}(\mathbf{F}) + \mathbf{M}^{-1} \text{vec}(\mathbf{Y}) \quad (4.21)$$

we devise an iterative scheme given by

$$\text{vec}(\mathbf{F}_{k+1}) = \mathbf{M}^{-1} \mathbf{N} \text{vec}(\mathbf{F}_k) + \mathbf{M}^{-1} \text{vec}(\mathbf{Y}) \quad (4.22)$$

When \mathbf{M} is a simple matrix that is easy to invert, this update function can be vastly more efficient to compute. Common approaches to finding a suitable value for \mathbf{M} and \mathbf{N} include the Jacobi, Gauss-Seidel and successive over-relaxation methods, each of which represent a different strategy for splitting the coefficient matrix [Saad, 2003]. However, whilst these techniques are well-studied, they are not appropriate for use in the case of graph signal reconstruction. This is because, for each of these methods, the coefficient matrix is split according to its diagonal and off-diagonal elements in some way. Consequently, this would require the evaluation of \mathbf{H}^{-2} directly which, as we have discussed, may be large, severely ill-conditioned and possibly ill-defined.

Instead, we require a custom splitting that avoids direct evaluation of \mathbf{H}^{-2} , and allows the right hand side of eq. (4.22) to be computed efficiently. The main contribution of this subsection is the identification of appropriate values for \mathbf{M} and \mathbf{N} , and an investigation of the consequences of that choice.

In the following, we set

$$\mathbf{M} = \gamma \mathbf{H}^{-2} + \mathbf{I}_{NT}, \quad \text{and} \quad \mathbf{N} = \text{diag}(\text{vec}(\mathbf{S}')). \quad (4.23)$$

where \mathbf{S}' is the binary matrix representing the complement of the set of selected nodes, i.e.

$$\mathbf{S}'_{nt} = \begin{cases} 1 & \text{if } (n, t) \notin \mathcal{S} \\ 0 & \text{otherwise} \end{cases} \quad (4.24)$$

In this way, the update equation is given by

$$\text{vec}(\mathbf{F}_{k+1}) = (\mathbf{H}^{-2} + \mathbf{I})^{-1} \text{diag}(\text{vec}(\mathbf{S}')) \text{vec}(\mathbf{F}_k) + (\mathbf{H}^{-2} + \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \quad (4.25)$$

Note that this splitting is valid since $(\mathbf{H}^{-2} + \mathbf{I})^{-1}$ is guaranteed to exist. It can also be readily computed as we already have the eigendecomposition of \mathbf{H} . Noting the decomposed definition of \mathbf{H} given in eq. (4.8), this can be written as

$$\begin{aligned} \mathbf{M}^{-1} &= (\mathbf{H}^{-2} + \mathbf{I})^{-1} \\ &= \left((\mathbf{U}_T \otimes \mathbf{U}_N) \text{diag}(\text{vec}(\mathbf{G}))^{-1} (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) + \mathbf{I} \right)^{-1} \\ &= (\mathbf{U}_T \otimes \mathbf{U}_N) \left(\text{diag}(\text{vec}(\mathbf{G}))^{-1} + \mathbf{I} \right)^{-1} (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \\ &= (\mathbf{U}_T \otimes \mathbf{U}_N) \text{diag}(\text{vec}(\mathbf{J})) (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \end{aligned} \quad (4.26)$$

where $\mathbf{J} \in \mathbb{R}^{N \times T}$ has elements defined by

$$\mathbf{J}_{nt} = \frac{\mathbf{G}_{nt}^2}{\mathbf{G}_{nt}^2 + \gamma}. \quad (4.27)$$

Note that the update formula can be computed with $O(N^2T + NT^2)$ complexity at each step.

$$\mathbf{F}_{k+1} = \mathbf{U}_N (\mathbf{J} \circ (\mathbf{U}_N^\top (\mathbf{S}' \circ \mathbf{F}_k) \mathbf{U}_T)) \mathbf{U}_T^\top + \mathbf{F}_0 \quad (4.28)$$

$$\text{with} \quad \mathbf{F}_0 = \mathbf{U}_N (\mathbf{J} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T)) \mathbf{U}_T^\top \quad (4.29)$$

Furthermore, this is reduced to $O(N^2T + NT \log T)$ in the case of T-V problems, and to $O(NT \log NT)$ for data residing on a grid (see section 4.2).

It is well-known that a given splitting will be convergent if the largest eigenvalue λ_{\max} of the matrix $\mathbf{M}^{-1}\mathbf{N}$ has an absolute value of less than one. This attribute, $\rho = |\lambda_{\max}|$, is known as the spectral radius.

Whilst the spectral radius of $\mathbf{M}^{-1}\mathbf{N}$ cannot be computed directly, we can derive an upper bound based on the properties of \mathbf{M} and \mathbf{N} individually.

Consider the spectral radius of \mathbf{M}^{-1} . By directly inspecting eq. (4.26), it is clear that $\rho(\mathbf{M}^{-1})$ will be the maximum entry in the matrix \mathbf{J} since \mathbf{M}^{-1} is already diagonalised in the basis $\mathbf{U}_T \otimes \mathbf{U}_N$. Consider now the definition of \mathbf{J} given in eq. (4.27). By definition, $g(\cdot)$ has a maximum value of one on the non-negative reals, achieved when its argument is zero. Since the graph Laplacian is guaranteed to have at least one zero eigenvalue, the maximum entry in the matrix \mathbf{J} , and therefore the spectral radius of \mathbf{M}^{-1} , is surely given by

$$\rho(\mathbf{M}^{-1}) = \frac{1}{1 + \gamma} \quad (4.30)$$

Next, consider the spectral radius of \mathbf{N} . This can be extracted directly as 1, since it is a diagonal binary matrix. Since both \mathbf{M}^{-1} and \mathbf{N} are positive semi-definite, we can apply the theorem

$$\rho(\mathbf{AB}) \leq \rho(\mathbf{A}) \rho(\mathbf{B}) \quad (4.31)$$

[Bhatia, 1997]. Therefore, the spectral radius of $\mathbf{M}^{-1}\mathbf{N}$ is guaranteed to be less than or equal to $1/(1 + \gamma)$. Since γ is strictly positive, this is less than one and, as such, convergence is guaranteed. We return to the question of convergence more thoroughly in section 4.3.1.

Finally, the update formulas given in eqs. (4.28) and (4.29) can be written equivalently as

$$\begin{aligned} \Delta \mathbf{F}_0 &= \mathbf{U}_N (\mathbf{J} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T)) \mathbf{U}_T^\top \\ \Delta \mathbf{F}_{k+1} &= \mathbf{U}_N (\mathbf{J} \circ (\mathbf{U}_N^\top (\mathbf{S}' \circ \Delta \mathbf{F}_k) \mathbf{U}_T)) \mathbf{U}_T^\top \end{aligned} \quad (4.32)$$

In this form, the iterations can be easily terminated when $|\Delta \mathbf{F}_k|$ is sufficiently small. The complete procedure is given in [algorithm 1](#).

Algorithm 1 Stationary iterative method with matrix splitting

Input: Observation matrix $\mathbf{Y} \in \mathbb{R}^{N \times T}$
Input: Sensing matrix $\mathbf{S} \in [0, 1]^{N \times T}$
Input: Space-like graph Laplacian $\mathbf{L}_N \in \mathbb{R}^{N \times N}$
Input: Time-like graph Laplacian $\mathbf{L}_T \in \mathbb{R}^{T \times T}$
Input: Regularisation parameter $\gamma \in \mathbb{R}^+$
Input: Graph filter function $g(\cdot; \boldsymbol{\beta} \in \mathbb{R}^2)$
 Decompose \mathbf{L}_N into $\mathbf{U}_N \boldsymbol{\Lambda}_L \mathbf{U}_N^\top$ and \mathbf{L}_T into $\mathbf{U}_T \boldsymbol{\Lambda}_T \mathbf{U}_T^\top$
 Compute $\mathbf{G} \in \mathbb{R}^{N \times T}$ as $\mathbf{G}_{nt} = g\left(\begin{bmatrix} \lambda_t^{(T)} \\ \lambda_n^{(N)} \end{bmatrix}, \boldsymbol{\beta}\right)$
 Compute $\mathbf{J} \in \mathbb{R}^{N \times T}$ as $\mathbf{J}_{nt} = \mathbf{G}_{nt}^2 / (\mathbf{G}_{nt}^2 + \gamma)$
 $\mathbf{S}' \leftarrow \mathbf{1} \in \mathbb{R}^{N \times T} - \mathbf{S}$
 $\Delta \mathbf{F} \leftarrow \mathbf{U}_N (\mathbf{J} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T)) \mathbf{U}_T^\top$
 $\mathbf{F} \leftarrow \Delta \mathbf{F}$
while $|\Delta \mathbf{F}| > \text{tol}$ **do**
 $\Delta \mathbf{F} \leftarrow \mathbf{U}_N (\mathbf{J} \circ (\mathbf{U}_N^\top (\mathbf{S}' \circ \Delta \mathbf{F}) \mathbf{U}_T)) \mathbf{U}_T^\top$
 $\mathbf{F} \leftarrow \mathbf{F} + \Delta \mathbf{F}$
end while
Output: \mathbf{F}

4.2.3 A conjugate gradient method

The second approach we consider for computing the posterior mean is to use the Conjugate Gradient Method (CGM). First proposed in 1952, the CGM is part of the Krylov subspace family, and is perhaps the most prominent iterative algorithm for solving linear systems [Hestenes and Stiefel, 1952]. In computational terms, the method only requires repeated forward multiplication of vectors by the coefficient matrix which, in the standard CGM, must be PSD. It is therefore effective in applications where this process can be performed efficiently.

In brief, the CGM seeks to solve the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ by minimising, at the k -th iteration, some measure of error in the affine space $\mathbf{x}_0 + \mathcal{K}_k$ where \mathcal{K}_k is the k -th Krylov subspace given by

$$\mathcal{K}_k = \text{span}(\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0)$$

The residual \mathbf{r}_k is given by

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}$$

and the k -th iterate of the CGM minimises

$$\phi(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{x}^\top \mathbf{b}$$

over $\mathbf{x}_0 + \mathcal{K}_k$ [Kelley, 1995].

The CGM works best when the coefficient matrix \mathbf{A} has a low condition number κ (that is, the ratio between the largest and smallest eigenvalue is small) and, as such, a preconditioning step is often necessary. The purpose of a preconditioner is to reduce κ by solving an equivalent transformed problem. This can be achieved by right or left multiplying the linear system by a preconditioning matrix Φ . However, this likely means the coefficient matrix is no longer PSD, meaning the CGM cannot be used in its basic form. (Other approaches modified for non-PSD matrices exists, e.g. the CGNE or GIMRES [Elman, 1982, Saad and Schultz, 1986]). A preconditioner can also multiply the coefficient matrix on the right by a preconditioner Φ^\top and the left by Φ . This preserves the symmetry meaning we can continued to use the regular CGM.

In our case, where the coefficient matrix is given by $(\text{diag}(\text{vec}(\mathbf{S})) + \gamma \mathbf{H}^{-2})$, preconditioning will be essential for convergence. To see why, consider the definition of \mathbf{H} in equation (4.8). A low-pass filter function $g(\cdot)$ may be close to zero when applied to the high-frequency eigenvalues of the graph Laplacian, meaning elements of $\text{diag}(\text{vec}(\mathbf{G}))^{-2}$ may be very high. In the worst case, for example with a band-limited filter, the matrix \mathbf{H} will be singular, no matrix \mathbf{H}^{-2} will exist, and the condition number of the coefficient matrix will be, in effect, infinite. Therefore, the primary purpose of this subsection is to find a preconditioner that maintains efficient forward multiplication and is effective at reducing the condition number of the coefficient matrix.

References such as [Saad, 2003] give a broad overview of the known approaches to finding a preconditioner. Standard examples include the Jacobi preconditioner which is given by the inverse of the coefficient matrix diagonal and is effective for diagonally dominant matrices, and the Sparse Approximate Inverse preconditioner [Grote and Huckle, 1997]. However, such preconditioners generally require direct evaluation of parts of the coefficient matrix or are computationally intensive to calculate.

In the following, we derive an effective symmetric preconditioner that allows forward multiplication of the coefficient matrix to be performed efficiently. First consider the transformed variable \mathbf{Z} , related to \mathbf{F} in the following way.

$$\mathbf{F} = \mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top \quad (4.33)$$

Here, \mathbf{Z} can be interpreted as set of frequency coefficients, which are subsequently scaled according to the graph filter function, and then reverse Fourier transformed back into the node domain. Matrices \mathbf{Z} which are distributed according to a spherically symmetric distribution, result in signals \mathbf{F} which are smooth with respect to the graph topology. Since this transform filters out the problematic high-frequency Fourier components, the system defined by this transformed variable \mathbf{Z} is naturally far better conditioned.

By substituting this expression for \mathbf{F} back into the likelihood in equation (4.13), and the prior of equation (4.14), one can derive a new expression for the posterior mean of \mathbf{Z} . This is done explicitly in theorem A.2. The end result is that the new linear system for the transformed variable \mathbf{Z} is given by

$$\text{vec}(\mathbf{Z}) = (\mathbf{C} + \gamma \mathbf{I}_T \otimes \mathbf{I}_N)^{-1} \text{vec}(\mathbf{G} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T)) \quad (4.34)$$

where \mathbf{C} is the symmetric PSD matrix

$$\mathbf{C} = \mathbf{D}_\mathbf{G} (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \mathbf{D}_\mathbf{S} (\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_\mathbf{G} \quad (4.35)$$

where we have abbreviated $\text{diag}(\text{vec}(\mathbf{G}))$ and $\text{diag}(\text{vec}(\mathbf{S}))$ as $\mathbf{D}_\mathbf{G}$ and $\mathbf{D}_\mathbf{S}$ respectively. Note that the conditioning of the coefficient matrix $\mathbf{C} + \gamma \mathbf{I}$ is greatly improved from the untransformed problem, as we will discuss in greater detail in section 4.3.2. Note also that multiplication of a vector $\text{vec}(\mathbf{R})$ by the coefficient matrix can be computed efficiently as

$$\gamma \mathbf{R} + \mathbf{G} \circ \left(\mathbf{U}_N^\top \left(\mathbf{S} \circ (\mathbf{U}_N (\mathbf{G} \circ \mathbf{R}) \mathbf{U}_T^\top) \right) \mathbf{U}_T \right)$$

This has $O(N^2T + NT^2)$ complexity at each step which may be reduced to $O(N^2T + NT \log T)$ in the case of T-V problems, and to $O(NT \log NT)$ for data residing on a grid (see section 4.2).

The linear system defined eq. (4.34) can be understood as a two-sided symmetrically preconditioned version of the original linear system given in eq. (4.18). In particular, the new expression can be constructed by modifying the original system in the following way.

$$\left(\Phi^\top (\mathbf{D}_S + \gamma \mathbf{H}^{-2}) \Phi\right) \left(\Phi^{-1} \text{vec}(\mathbf{F})\right) = \Phi^\top \text{vec}(\mathbf{Y}), \quad (4.36)$$

where

$$\Phi = (\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_G. \quad (4.37)$$

Since preconditioning of the coefficient matrix on the left is achieved with Φ^\top and on the right with Φ , symmetry is preserved. This ensures that one can continue to utilise algorithms tailored to work with PSD matrices. In [algorithm 2](#), we outline a conjugate gradient method based on this new formulation. Just as in [section 4.2.2](#), this process takes advantage of several identities involving Kronecker products and vectorization, to give $O(N^2T + NT^2)$ complexity per iteration. Note also that, in the case of a time-vertex problem where \mathcal{G}_T is given by a path graph, this complexity can be further reduced to $O(N^2T + NT \log T)$.

4.3 Convergence properties

In this section, we describe theoretical bounds for the convergence rate of both the SIM and the CGM in terms of the hyperparameter γ .

4.3.1 Convergence of the SIM

First, denote the difference between the true solution $\text{vec}(F)$ and at the k -th iteration as \mathbf{E}_k , we can see that the error is updated according to the following equation.

$$\begin{aligned} \text{vec}(\mathbf{E}_k) &= \text{vec}(\mathbf{F}_k) - \text{vec}(\mathbf{F}_{k-1}) \\ &= \mathbf{M}^{-1} \mathbf{N} \text{vec}(\mathbf{F}_{k-1}) + \mathbf{M}^{-1} \text{vec}(\mathbf{Y}) - (\mathbf{M}^{-1} \mathbf{N} \text{vec}(\mathbf{F}_{k-2}) + \mathbf{M}^{-1} \text{vec}(\mathbf{Y})) \\ &= \mathbf{M}^{-1} \mathbf{N} (\text{vec}(\mathbf{F}_{k-1}) - \text{vec}(\mathbf{F}_{k-2})) \\ &= \mathbf{M}^{-1} \mathbf{N} \text{vec}(\mathbf{E}_{k-1}) \\ &= (\mathbf{M}^{-1} \mathbf{N})^{k-1} \text{vec}(\mathbf{E}_1) \end{aligned} \quad (4.38)$$

From this it is clear to see that convergence will be achieved so long as the maximum eigenvalue of $\mathbf{M}^{-1} \mathbf{N}$ is less than 1, since, if this condition holds then,

Algorithm 2 Conjugate gradient method with graph-spectral preconditioner

Input: Observation matrix $\mathbf{Y} \in \mathbb{R}^{N \times T}$
Input: Sensing matrix $\mathbf{S} \in [0, 1]^{N \times T}$
Input: Space-like graph Laplacian $\mathbf{L}_N \in \mathbb{R}^{N \times N}$
Input: Time-like graph Laplacian $\mathbf{L}_T \in \mathbb{R}^{T \times T}$
Input: Regularisation parameter $\gamma \in \mathbb{R}$
Input: Graph filter function $g(\cdot; \beta)$

Decompose \mathbf{L}_N into $\mathbf{U}_N \mathbf{\Lambda}_N \mathbf{U}_N^\top$ and \mathbf{L}_T into $\mathbf{U}_T \mathbf{\Lambda}_T \mathbf{U}_T^\top$

Compute $\mathbf{G} \in \mathbb{R}^{N \times T}$ as $\mathbf{G}_{nt} = g\left(\begin{bmatrix} \lambda_t^{(T)} \\ \lambda_n^{(N)} \end{bmatrix}, \beta\right)$

Initialise $\mathbf{Z} \in \mathbb{R}^{N \times T}$ randomly

$\mathbf{R} \leftarrow \mathbf{G} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T) - \gamma \mathbf{Z} - \mathbf{G} \circ \left(\mathbf{U}_N^\top (\mathbf{S} \circ (\mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top)) \mathbf{U}_T \right)$

$\mathbf{D} \leftarrow \mathbf{R}$

while $|\Delta \mathbf{R}| > \text{tol}$ **do**

$\mathbf{A}_D \leftarrow \gamma \mathbf{D} + \mathbf{G} \circ \left(\mathbf{U}_N^\top (\mathbf{S} \circ (\mathbf{U}_N (\mathbf{G} \circ \mathbf{D}) \mathbf{U}_T^\top)) \mathbf{U}_T \right)$

$\alpha \leftarrow \text{tr}(\mathbf{R}^\top \mathbf{R}) / \text{tr}(\mathbf{R}^\top \mathbf{A}_D \mathbf{R})$

$\mathbf{Z} \leftarrow \mathbf{Z} + \alpha \mathbf{D}$

$\mathbf{R} \leftarrow \mathbf{R} - \alpha \mathbf{A}_D$

$\delta \leftarrow \text{tr}(\mathbf{R}^\top \mathbf{R}) / \text{tr}((\mathbf{R} + \alpha \mathbf{A}_D)^\top (\mathbf{R} + \alpha \mathbf{A}_D))$

$\mathbf{D} \leftarrow \mathbf{R} + \delta \mathbf{D}$

end while

Output: $\mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top$

$$\lim_{k \rightarrow \infty} \text{vec}(\mathbf{E}_{k+1}) = \lim_{k \rightarrow \infty} (\mathbf{M}^{-1} \mathbf{N})^k \text{vec}(\mathbf{E}_1) = \mathbf{0}. \quad (4.39)$$

For the stationary iterative method, equation (4.38) describes how the magnitude of the error, $\text{vec}(\mathbf{E}_k)$, is reduced at each iteration. Convergence is guaranteed as long as the spectral radius of $\mathbf{M}^{-1} \mathbf{N}$, denoted ρ , is less than one which we have demonstrated must be true. In particular, the error at iteration $k + 1$ will take on the following value.

$$\begin{aligned} \text{vec}(\mathbf{E}_{k+1}) &= (\mathbf{M}^{-1} \mathbf{N})^k \text{vec}(\mathbf{E}_1) \\ &= (\mathbf{U} \mathbf{D}_J \mathbf{U}^\top \mathbf{D}_{S'})^k \text{vec}(\mathbf{E}_1) \end{aligned} \quad (4.40)$$

where we have used the shorthand $\mathbf{U} = \mathbf{U}_T \otimes \mathbf{U}_N$, $\mathbf{D}_J = \text{diag}(\text{vec}(\mathbf{J}))$ and $\mathbf{D}_{S'} = \text{diag}(\text{vec}(\mathbf{S}'))$. The slowest possible rate of convergence will be achieved when $\text{vec}(\mathbf{E}_1)$

is proportional to the eigenvector of $\mathbf{U}\mathbf{D}_\mathbf{J}\mathbf{U}^\top\mathbf{D}_{\mathbf{S}'}$ with the largest associated eigenvalue, namely the spectral radius $\rho(\mathbf{U}\mathbf{D}_\mathbf{J}\mathbf{U}^\top\mathbf{D}_{\mathbf{S}'})$. In this case, we have that

$$\text{vec}(\mathbf{E}_{k+1}) = \rho(\mathbf{U}\mathbf{D}_\mathbf{J}\mathbf{U}^\top\mathbf{D}_{\mathbf{S}'})^k \text{vec}(\mathbf{E}_1) \quad (4.41)$$

Consider the value of $\text{vec}(\mathbf{E}_1)$ that maximises the norm of the error for some value of k , subject to the constraint that $\|\text{vec}(\mathbf{E}_1)\|^2 = 1$.

$$\begin{aligned} \mathbf{e}^\star &= \underset{\text{vec}(\mathbf{E}_1)}{\text{argmax}} \quad \|\text{vec}(\mathbf{E}_{k+1})\|^2 \\ &= \underset{\text{vec}(\mathbf{E}_1)}{\text{argmax}} \quad \frac{(\mathbf{M}^{-1}\mathbf{N})^k \text{vec}(\mathbf{E}_1)}{\text{vec}(\mathbf{E}_1)^\top \text{vec}(\mathbf{E}_1)} \end{aligned} \quad (4.42)$$

In particular, the smallest possible rate of error reduction occurs when $\text{vec}(\mathbf{E}_0) = k\mathbf{1}$, i.e. it is proportional to the first eigenvector of the graph Laplacian, which is constant over all nodes. In this scenario, the worst-case error reduction follows

$$\text{vec}(\mathbf{E}_k) = \frac{1}{(1+\gamma)^k} \text{vec}(\mathbf{E}_0) \quad (4.43)$$

This implies that the number of iterations, n , required to reduce the error to a fraction ϵ follows

$$\epsilon = \frac{1}{(1+\gamma)^n} \quad \rightarrow \quad n = \frac{-\log \epsilon}{\log(1+\gamma)}$$

Therefore, for a given requirement in error reduction, the complexity of [algorithm 1](#) is bounded by

$$\frac{m}{\log(1+\gamma)} = m \left(\frac{1}{\gamma} + \frac{1}{2} - \frac{\gamma}{12} + \frac{\gamma^2}{24} - \dots \right)$$

where m is the number of multiplications required to complete each iteration. From the Taylor series expansion, we can see that the dominant behaviour for small γ is $O(\gamma^{-1})$. Therefore, for small γ , the overall run-time complexity of the SIM is given by

$$O\left(\frac{N^2T + NT^2}{\log(1+\gamma)}\right) \approx O\left(\frac{N^2T + NT^2}{\gamma}\right) \quad (4.44)$$

4.3.2 Convergence of the CGM

The conjugate gradient method, by contrast, is known to have a theoretical time complexity of $O(\sqrt{\kappa})$, where κ is the condition number of the coefficient matrix [Kelley \[1995\]](#). In our case, the coefficient matrix is given in equation (4.34) as $\mathbf{C} + \gamma \mathbf{I}$, and its associated condition number can be bounded.

Consider the definition for \mathbf{C} given in equation (4.35). Note that it can be described as the product of matrices, $\mathbf{C} = \mathbf{C}_1 \mathbf{C}_2 \mathbf{C}_1$ where

$$\mathbf{C}_1 = \mathbf{D}_{\mathbf{G}}, \quad \mathbf{C}_2 = (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \mathbf{D}_{\mathbf{S}} (\mathbf{U}_T \otimes \mathbf{U}_N)$$

The spectral radius $\rho(\mathbf{C}_1)$, is clearly 1, since the maximum output of $g(\cdot)$ on the domain \mathbb{R}^+ , and correspondingly the maximum element contained in the matrix \mathbf{G} , is 1. Similarly, $\rho(\mathbf{C}_2)$ is also 1. This is clear to see since it is already diagonalised in the basis $\mathbf{U}_T \otimes \mathbf{U}_N$, with the matrix \mathbf{S} being binary. Taking these facts together, and making use of equation (4.31), we see that the maximum eigenvalue of the matrix \mathbf{C} is 1. On the other hand, the minimum possible eigenvalue of \mathbf{C} is zero. This could occur if, say, \mathbf{G} represented a band-limited frequency response, i.e. \mathbf{G} contains zero elements. Clearly, this means \mathbf{C} transforms non-trivial vectors to zero.

We know that the maximum possible eigenvalue of the matrix \mathbf{C} is 1, and the minimum possible eigenvalue is 0. Accordingly, the maximum eigenvalue of the coefficient matrix $\mathbf{C} + \gamma \mathbf{I}$ is $1 + \gamma$ and the minimum is γ . This gives an upper bound for the condition number.

$$\kappa \leq \frac{1 + \gamma}{\gamma} \tag{4.45}$$

This limits the time complexity of the whole procedure stated in [algorithm 2](#) to an upper bound of

$$m \sqrt{\frac{1 + \gamma}{\gamma}} = m \left(\frac{1}{\sqrt{\gamma}} + \frac{\gamma}{2\sqrt{\gamma}} - \frac{\gamma^3}{8\sqrt{\gamma}} + \dots \right)$$

From this expansion, we can see that the dominant behaviour for small γ is $O(\gamma^{-1/2})$. Therefore, for small γ , the overall run-time complexity of the CGM is given by

$$O \left(\sqrt{\frac{1 + \gamma}{\gamma}} (N^2 T + N T^2) \right) \approx O \left(\frac{N^2 T + N T^2}{\sqrt{\gamma}} \right) \tag{4.46}$$

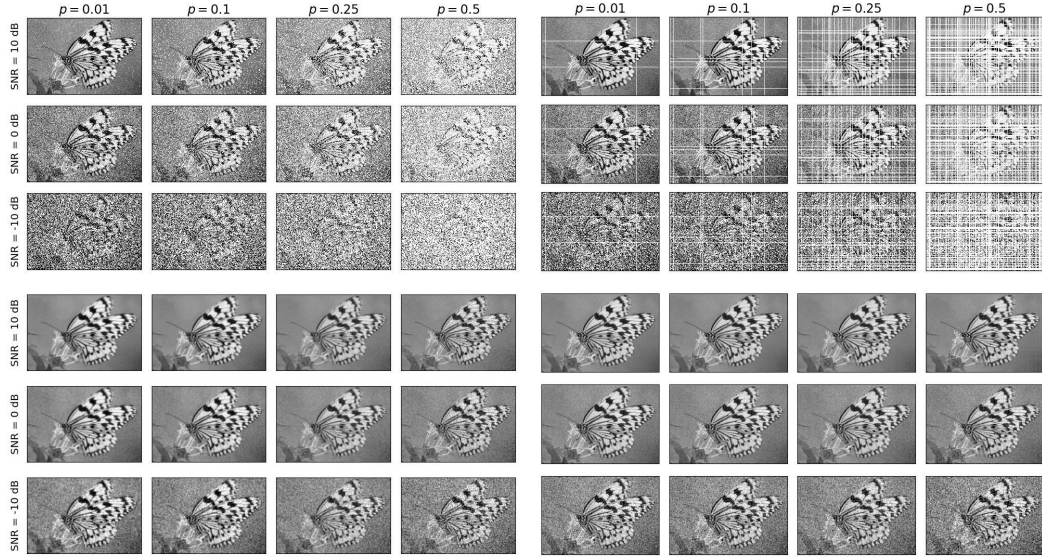


FIGURE 4.4: The output from the first experiment is depicted. In the top left quadrant, the input images are shown across a range of noise levels and missing pixel percentages, with missing pixels chosen uniformly at random. Below that, in the lower left quadrant, the corresponding reconstructed images are shown. The right half of the plot is the same, except here entire columns and rows of pixels are removed at random.

Furthermore, in the case of a time-vertex problem, or a

4.4 Image processing experiments

In this section, we run several small experiments to corroborate the properties of the CGM and SIM. In particular, we verify a) that both algorithms result in the same output for a given problem; b) that the runtime of the SIM and CGM increases at a significantly lower rate than naive Gaussian elimination as nodes are added; and c) that the number of iterations required for convergence as a function of γ is less than or equal to the upper bounds derived in section 4.3.

In order to perform these checks, we use a small image denoising/reconstruction task. In particular, we use a grey-scale image of width 571 and height of 856 pixels which has been zero-centred and normalised. Note that an image can be considered a graph signal, with each pixel site representing a node connected to the other pixels immediately adjacent. Furthermore, the data itself lies on a two-dimensional lattice which is a special case of a product graph, with each sub-graph in the product being a simple chain, in this case with $N = 856$ and $T = 571$ respectively. Therefore, images can serve as a simple test case for checking the behaviour of product graph algorithms. Note that the purpose of this section is not to compare graph signal reconstruction to existing algorithms specialised

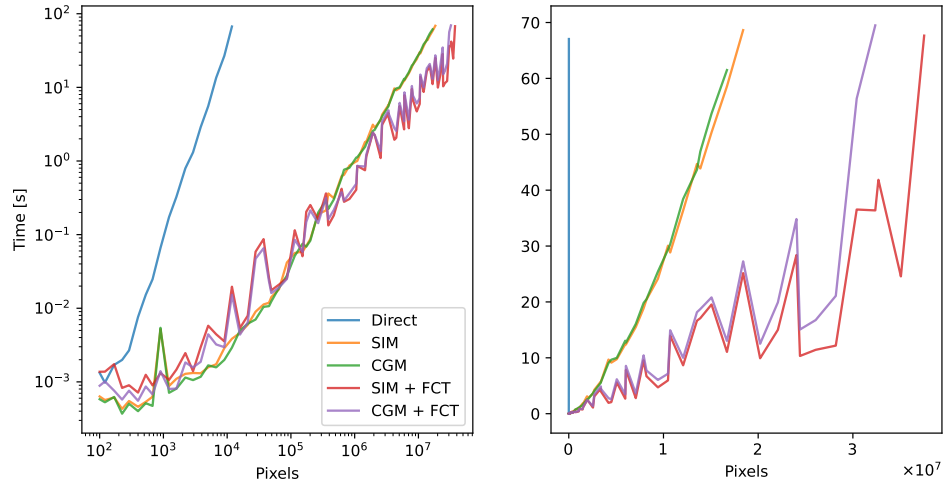


FIGURE 4.5: The total runtime in seconds for the SIM and CGM compared to a naive Gaussian elimination approach is shown as a function of the total number of nodes using a quad-core Intel i7-7700HQ CPU.

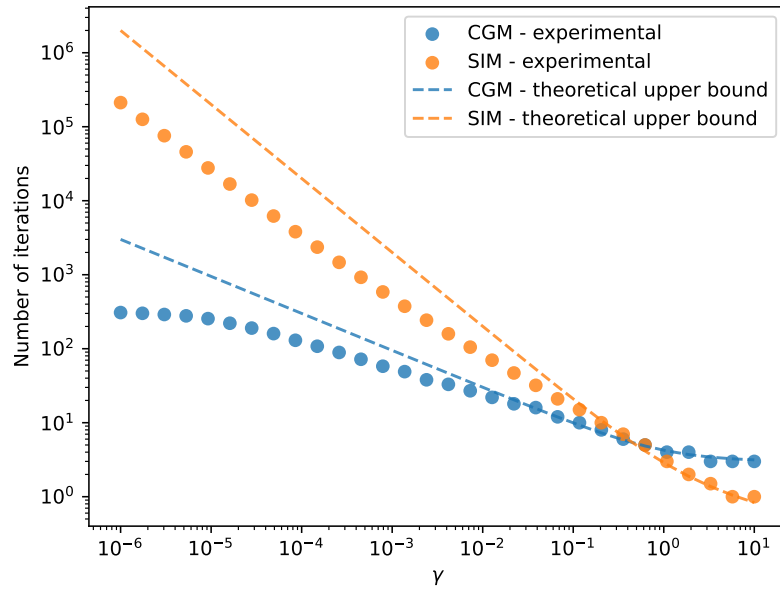


FIGURE 4.6: The number of iterations required to reach a certain level of precision is shown experimentally, along with the theoretical upper bound, for the SIM and CGM

for image denoising/reconstruction, but to examine the computational properties of GSR.

First, we simulate a graph signal reconstruction task across a range of noise levels and missing data distributions. In particular, we add white Gaussian noise to the raw image with a Signal to Noise Ratio (SNR) of -10, 0 and 10 dB. We also remove pixels according to two rules. First, pixels are removed uniformly at random, and second, we remove

entire rows and columns uniformly at random such that a fixed total percentage p of the pixels is missing. We do this for p equal to 0.01, 0.1, 0.25 and 0.5 giving 24 unique trials which are shown in [figure 2](#). We find that both the SIM and the CGM converge quickly and result in identical predicted outputs, as expected. Furthermore, the statistical model seems relatively robust to noise and missing data of both types, visually performing adequately at the level of $p = 0.5$ with SNR=-10dB, which is typically considered a challenging image reconstruction task.

Next, we measured the total runtime of the SIM and the CGM as the number of nodes was increased, and compared this to a more naive approach to solving equation (4.18) which is direct Gaussian elimination. In particular, we fixed the SNR at 0dB and the missingness to $p = 0.5$, with missing pixels chosen uniformly at random, and created square images of increasing size from 10^2 to 10^7 total pixels. The results are shown in [figure 3](#). As is visible, the runtime of the SIM and CGM scales at a significantly slower rate than a naive approach, remaining tractable at well above 10^7 nodes.

Finally, we test the effect that varying the hyperparameter γ has on the number of iterations required for convergence. To do this, we fixed the SNR at 0dB and the missingness to $p = 0.5$, with missing pixels chosen uniformly at random. Next, we varied γ in logarithmically spaced increments from 10^{-6} to 10^1 . For each unique value of γ , we ran both the SIM and the CGM and counted the number of steps required for each algorithm to reach a specific level of precision (10^{-8} across all elements). The results are shown in [figure 4](#). We also plot the theoretical upper bound derived in the previous subsection for each method.

As is visible, both methods converge within the bounds of their theoretical worst-case complexity. As expected, the CGM takes significantly fewer iterations than the SIM to converge at small γ . It is also visible that the CGM seems to outperform the theoretical worst-case scaling of $\gamma^{-0.5}$, converging with a rate closer to $\gamma^{-0.3}$ over the majority of the range, and even showing signs of reducing further at very small γ . The SIM on the other hand empirically converges slightly faster but close to the theoretical worst-case rate of γ^{-1} . Interestingly, the SIM converges faster than the CGM at relatively high γ , with the cross-over occurring at around $\gamma = 0.5$. Although the absolute number of iterations required for convergence in this domain is low (between 1 and 10), this could still be significant for very large problems, meaning it still has some value as a solution.

4.5 Kernel Graph Regression with Unrestricted Missing Data Patterns

Hello

4.5.1 Cartesian product graphs and KGR

Hello

4.5.2 Convergence properties

Hello

4.6 Regression with Network Cohesion

Hello

4.6.1 Regression with node-level covariates

Hello

4.6.2 Convergence properties

Hello

Chapter 5

Multi-Dimensional Cartesian Product Graphs

Hello

5.1 Fast computation with d -dimensional Kronecker products

Hello

5.2 Signal reconstruction

Hello

5.3 Kernel Graph Regression

Hello

5.4 Regression with Network Cohesion

Chapter 6

Signal Uncertainty: Estimation and Sampling

6.1 Introduction

6.2 Posterior Estimation

6.2.1 Log-variance prediction

6.2.2 Estimation models

6.2.3 Query strategies

6.2.4 Comparison and analysis

6.3 Posterior Sampling

6.3.1 Perturbation optimization

6.4 Estimation vs Sampling

6.4.1 Experiments

Chapter 7

Working with Binary-Valued Graph Signals

7.1 Logistic Graph Signal Reconstruction

7.2 Logistic Kernel Graph Regression

7.3 Logistic Regression with Network Cohesion

7.4 Approximate Sampling via the Laplace Approximation

Chapter 8

Conclusions

8.1 Main Section 1

Appendix A

Proofs

Theorem A.1. *The posterior distribution for \mathbf{F} is given by*

$$\text{vec}(\mathbf{F}) \mid \mathbf{Y} \sim \mathcal{N}(\boldsymbol{\Sigma} \text{vec}(\mathbf{Y}), \boldsymbol{\Sigma}) \quad (\text{A.1})$$

where

$$\boldsymbol{\Sigma} = \left(\text{diag}(\text{vec}(\mathbf{S})) + \gamma \mathbf{H}^{-2} \right)^{-1} \quad (\text{A.2})$$

Proof. Consider the matrix \mathbf{S}_ϵ defined in the following manner.

$$(\mathbf{S}_\epsilon)_{nt} = \begin{cases} 1 & \text{if } (n, t) \in \mathcal{S} \\ \epsilon & \text{otherwise} \end{cases} \quad (\text{A.3})$$

We can use this definition to rewrite equation 4.13 for the probability distribution of $\mathbf{Y} \mid \mathbf{F}$.

$$\text{vec}(\mathbf{Y}) \mid \mathbf{F} \sim \lim_{\epsilon \rightarrow 0} \left[\mathcal{N}(\text{vec}(\mathbf{S}_\epsilon \circ \mathbf{F}), \text{diag}(\text{vec}(\mathbf{S}_\epsilon))) \right] \quad (\text{A.4})$$

In this way, the negative log-likelihood of an observation $\mathbf{Y} \mid \mathbf{F}$ is given by

$$-\log \pi(\mathbf{Y} \mid \mathbf{F}) = \lim_{\epsilon \rightarrow 0} \left[\frac{1}{2} \text{vec}(\mathbf{S}_\epsilon \circ \mathbf{F} - \mathbf{Y})^\top \text{diag}(\text{vec}(\mathbf{S}_\epsilon))^{-1} \text{vec}(\mathbf{S}_\epsilon \circ \mathbf{F} - \mathbf{Y}) \right] \quad (\text{A.5})$$

up to an additive constant which does not depend on \mathbf{F} . Note that, since $\mathbf{Y} = \mathbf{S}_\epsilon \circ \mathbf{Y}$, we can rewrite $\text{vec}(\mathbf{S}_\epsilon \circ \mathbf{F} - \mathbf{Y})$ as

$$\begin{aligned} \text{vec}(\mathbf{S}_\epsilon \circ \mathbf{F} - \mathbf{Y}) &= \text{vec}(\mathbf{S}_\epsilon \circ (\mathbf{F} - \mathbf{Y})) \\ &= \text{diag}(\text{vec}(\mathbf{S}_\epsilon)) \text{vec}(\mathbf{F} - \mathbf{Y}) \end{aligned} \quad (\text{A.6})$$

Therefore, equation A.5 can be rewritten as

$$\begin{aligned} -\log \pi(\mathbf{Y}|\mathbf{F}) &= \lim_{\epsilon \rightarrow 0} \left[\frac{1}{2} \text{vec}(\mathbf{F} - \mathbf{Y})^\top \text{diag}(\text{vec}(\mathbf{S}_\epsilon)) \text{vec}(\mathbf{F} - \mathbf{Y}) \right] \\ &= \frac{1}{2} \text{vec}(\mathbf{F} - \mathbf{Y})^\top \text{diag}(\text{vec}(\mathbf{S})) \text{vec}(\mathbf{F} - \mathbf{Y}) \end{aligned} \quad (\text{A.7})$$

Now consider the full log-posterior. Using Bayes rule, this can be written as

$$\begin{aligned} -\log \pi(\text{vec}(\mathbf{F}) | \mathbf{Y}) &= \frac{1}{2} \text{vec}(\mathbf{F} - \mathbf{Y})^\top \text{diag}(\text{vec}(\mathbf{S})) \text{vec}(\mathbf{F} - \mathbf{Y}) + \\ &\quad \frac{\gamma}{2} \text{vec}(\mathbf{F})^\top \mathbf{H}^{-2} \text{vec}(\mathbf{F}) \end{aligned} \quad (\text{A.8})$$

Up to an additive constant not dependent \mathbf{F} , this can be written as

$$-\log \pi(\text{vec}(\mathbf{F}) | \mathbf{Y}) = \frac{1}{2} \left(\text{vec}(\mathbf{F})^\top (\text{diag}(\text{vec}(\mathbf{S})) + \gamma \mathbf{H}^{-2}) \text{vec}(\mathbf{F}) - 2 \text{vec}(\mathbf{Y})^\top \mathbf{F} \right) \quad (\text{A.9})$$

Using the conjugacy of the normal distribution, by direct inspection we can conclude that the posterior covariance is given by

$$\mathbf{\Sigma} = \left(\text{diag}(\text{vec}(\mathbf{S})) + \gamma \mathbf{H}^{-2} \right)^{-1} \quad (\text{A.10})$$

and that the posterior mean is given by $\mathbf{\Sigma} \text{vec}(\mathbf{Y})$.

□

Theorem A.2. Consider the random matrix \mathbf{Z} which is related to the random matrix \mathbf{F} as follows.

$$\mathbf{F} = \mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top$$

or, equivalently,

$$\text{vec}(\mathbf{F}) = (\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_\mathbf{G} \text{vec}(\mathbf{Z})$$

Then the posterior mean for $\mathbf{Z}|\mathbf{Y}$ is given by

$$\mathbb{E}[\mathbf{Z}|\mathbf{Y}] = (\mathbf{C} + \gamma \mathbf{I}_T \otimes \mathbf{I}_N)^{-1} \text{vec}(\mathbf{G} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T))$$

where

$$\mathbf{C} = \mathbf{D}_\mathbf{G} (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \mathbf{D}_\mathbf{S} (\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_\mathbf{G}$$

(Here we have abbreviated $\text{diag}(\text{vec}(\mathbf{G}))$ and $\text{diag}(\text{vec}(\mathbf{S}))$ as $\mathbf{D}_\mathbf{G}$ and $\mathbf{D}_\mathbf{S}$ respectively.)

Proof. The conditional distribution of $\mathbf{Y}|\mathbf{Z}$ is obtained by substituting in the definition of \mathbf{F} in terms of \mathbf{Z} into the original conditional likelihood expression.

$$\text{vec}(\mathbf{Y}) | \mathbf{Z} \sim \mathcal{N}\left(\text{vec}\left(\mathbf{S} \circ (\mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top)\right), \mathbf{D}_\mathbf{S}\right)$$

Similarly, since the prior specified for \mathbf{F} is $\mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{H}^2)$, this implies that the prior over \mathbf{Z} is simply

$$\text{vec}(\mathbf{Z}) \sim \mathcal{N}(\mathbf{0}, \gamma^{-1} \mathbf{I}_{NT})$$

To see this, consider the following

$$\begin{aligned} \text{Cov}[\text{vec}(\mathbf{F})] &= \text{Cov}[(\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_\mathbf{G} \text{vec}(\mathbf{Z})] \\ &= (\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_\mathbf{G} \text{Cov}[\text{vec}(\mathbf{Z})] \mathbf{D}_\mathbf{G} (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \end{aligned}$$

If $\text{vec}(\mathbf{Z})$ has covariance $\gamma^{-1} \mathbf{I}$, then $\text{vec}(\mathbf{F})$ has covariance given by

$$\begin{aligned}\text{Cov}[\text{vec}(\mathbf{F})] &= \gamma^{-1}(\mathbf{U}_T \otimes \mathbf{U}_N) \mathbf{D}_{\mathbf{G}}^2 (\mathbf{U}_T^\top \otimes \mathbf{U}_N^\top) \\ &= \gamma^{-1} \mathbf{H}^2\end{aligned}$$

by the definition of \mathbf{H} .

Now consider the transformed posterior

$$\begin{aligned}-\log p(\mathbf{Z}|\mathbf{Y}) &= -\log p(\mathbf{Y}|\mathbf{Z}) - \log p(\mathbf{Z}) \\ &= \frac{1}{2} \text{vec}(\mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top - \mathbf{Y})^\top \times \\ &\quad \mathbf{D}_{\mathbf{S}} \text{vec}(\mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top - \mathbf{Y}) \\ &\quad + \frac{\gamma}{2} \text{vec}(\mathbf{Z})^\top \text{vec}(\mathbf{Z})\end{aligned}$$

Up to an additive constant, this is equal to

$$\begin{aligned}-\log p(\mathbf{Z}|\mathbf{Y}) &= \frac{1}{2} \text{vec}(\mathbf{Z})^\top (\mathbf{C} + \gamma \mathbf{I}_{NT}) \text{vec}(\mathbf{Z}) \\ &\quad - \text{vec}(\mathbf{U}_N (\mathbf{G} \circ \mathbf{Z}) \mathbf{U}_T^\top)^\top \text{vec}(\mathbf{Y}) \\ &= \frac{1}{2} \text{vec}(\mathbf{Z})^\top (\mathbf{C} + \gamma \mathbf{I}_{NT}) \text{vec}(\mathbf{Z}) \\ &\quad - \text{vec}(\mathbf{Z})^\top \text{vec}(\mathbf{G} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T))\end{aligned}$$

By inspection, again, we can see that the posterior mean for \mathbf{Z} is

$$(\mathbf{C} + \gamma \mathbf{I}_T \otimes \mathbf{I}_N)^{-1} \text{vec}(\mathbf{G} \circ (\mathbf{U}_N^\top \mathbf{Y} \mathbf{U}_T))$$

□

Bibliography

- Ahmed, N., Natarajan, T., and Rao, K. (1974). Discrete cosine transform. *IEEE Transactions on Computers*, C-23(1):90–93.
- Barik, S., Bapat, R. B., and Pati, S. (2015). On the laplacian spectra of product graphs. *Applicable Analysis and Discrete Mathematics*, 9:39–58.
- Barik, S., Kalita, D., Pati, S., and Sahoo, G. (2018). Spectra of graphs resulting from various graph operations and products: a survey. *Special Matrices*, 6:323 – 342.
- Bhatia, R. (1997). *Matrix analysis*. Number 169 in Graduate texts in mathematics. Springer, New York.
- Brenner, S. C., Scott, L. R., and Scott, L. R. (2008). *The mathematical theory of finite element methods*, volume 3. Springer.
- Cooley, J. W. and Tukey, J. W. (1965). An algorithm for the machine calculation of complex fourier series. *Mathematics of Computation*, 19:297–301.
- Elman, H. C. (1982). *Iterative methods for large, sparse, nonsymmetric systems of linear equations*. PhD thesis.
- Fiedler, M. (1973). Algebraic connectivity of graphs. *Czechoslovak Mathematical Journal*, 23:298–305.
- Grassi, F., Loukas, A., Perraudin, N., and Ricaud, B. (2018). A time-vertex signal processing framework: Scalable processing and meaningful representations for time-series on graphs. *IEEE Transactions on Signal Processing*, 66(3):817–829.
- Grote, M. J. and Huckle, T. (1997). Parallel preconditioning with sparse approximate inverses. *SIAM Journal on Scientific Computing*, 18(3):838–853.
- Harzheim, E. (2005). Chapter 4 - products of orders. In *Ordered Sets*, volume 7 of *Advances in Mathematics*. Springer-Verlag, New York.
- Hestenes, M. R. and Stiefel, E. (1952). Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49:409–435.

- Imrich, W. and Klavžar, S. (2000). *Product Graphs: Structure and Recognition*. A Wiley-Interscience publication. Wiley.
- Isufi, E., Loukas, A., Simonetto, A., and Leus, G. (2017). Autoregressive moving average graph filtering. *IEEE Transactions on Signal Processing*, 65(2):274–288.
- Jiang, J. (2012). Introduction to spectral graph theory.
- Kaveh, A. and Alinejad, B. (2011). Laplacian matrices of product graphs: applications in structural mechanics. *Acta Mechanica*, 222:331–350.
- Kelley, C. T. (1995). *Iterative Methods for Linear and Nonlinear Equations*. Society for Industrial and Applied Mathematics.
- Loukas, A. and Foucard, D. (2016). Frequency analysis of time-varying graph signals. In *2016 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 346–350.
- Makhoul, J. (1980). A fast cosine transform in one and two dimensions. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(1):27–34.
- Mieghem, P. v. (2010). *Graph Spectra for Complex Networks*. Cambridge University Press.
- Newman, M. (2018). *Networks*. Oxford University Press.
- Ortega, A., Frossard, P., Kovačević, J., Moura, J. M. F., and Vandergheynst, P. (2018). Graph signal processing: Overview, challenges, and applications. *Proceedings of the IEEE*, 106(5):808–828.
- Rao, K. and Yip, P. (1990). *Discrete Cosine Transform: Algorithms, Advantages, Applications*. Elsevier Science & Technology Books.
- Saad, Y. (2003). *Iterative Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematics, second edition.
- Saad, Y. and Schultz, M. H. (1986). Gmres: a generalized minimal residual algorithm for solving nonsymmetric linear systems. *Siam Journal on Scientific and Statistical Computing*, 7:856–869.
- Sayama, H. (2016). Estimation of laplacian spectra of direct and strong product graphs. *Discrete Applied Mathematics*, 205:160–170.
- Shuman, D. I., Narang, S. K., Frossard, P., Ortega, A., and Vandergheynst, P. (2013). The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine*, 30(3):83–98.