

Tales Cesar de Oliveira Imbiriba

**NONLINEAR HYPERSPECTRAL UNMIXING:
STRATEGIES FOR NONLINEAR MIXTURE DETECTION,
ENDMEMBER ESTIMATION AND BAND-SELECTION**

Tese submetida ao Programa de Pós-graduação
em Engenharia Elétrica para a obtenção
do Grau de Doutor.

Orientador

Universidade Federal de Santa Catarina:
Prof. José Carlos Moreira Bermudez, Ph.D.

Florianópolis

2016

Ficha de identificação da obra elaborada pelo autor através do
Programa de Geração Automática da Biblioteca Universitária da UFSC.

Imbiriba, Tales Cesar de Oliveira

Nonlinear hyperspectral unmixing : strategies for
nonlinear mixture detection, endmember estimation and
band-selection / Tales Cesar de Oliveira Imbiriba ;
orientador, José Carlos Moreira Bermudez - Florianópolis,
SC, 2016.

171 p.

Tese (doutorado) - Universidade Federal de Santa
Catarina, Centro Tecnológico. Programa de Pós-Graduação em
Engenharia Elétrica.

Inclui referências

1. Engenharia Elétrica. 2. Imagem hiperespectral.
3. Detecção. 4. RKHS. 5. Seleção de Bandas. I.
Bermudez, José Carlos Moreira. II. Universidade Federal de
Santa Catarina. Programa de Pós-Graduação em Engenharia
Elétrica. III. Título.

Tales Cesar de Oliveira Imbiriba

**NONLINEAR HYPERSPECTRAL UNMIXING: STRATEGIES
FOR NONLINEAR MIXTURE DETECTION, ENDMEMBER
ESTIMATION AND BAND-SELECTION**

Esta Tese foi julgada aprovada para a obtenção do Título de “Doutor”, e aprovada em sua forma final pelo Programa de Pós-graduação em Engenharia Elétrica.

Florianópolis, 10 de novembro 2016.

Prof. Marcelo Lobo Heldwein, Dr.
Coordenador
Universidade Federal de Santa Catarina

Banca Examinadora:

Prof. José Carlos Moreira Bermudez, Ph.D.
Orientador
Universidade Federal de Santa Catarina

Prof. Márcio Holsbach Costa, Dr.
Universidade Federal de Santa Catarina

Prof. Fermín Sinforiano Viloche Bazán, Dr.
Universidade Federal de Santa Catarina

Prof. Leonardo Silva Resende, Dr.
Universidade Federal de Santa Catarina

Prof. Alejandro C. Frery Orgambide, Ph.D.
Universidade Federal de Alagoas

Prof. Aldebaro B. da Rocha Kautau Junior, Ph.D.
Universidade Federal do Pará

Dedico esta tese a minha amada esposa,
Fernanda.

AGRADECIMENTOS

Tenho muito a agradecer. Agradeço por todas as oportunidades que tive nos meus já 34 anos de idade, em parte atribuídas à minha maravilhosa família, amigos, professores e tantas outras pessoas importantes que passaram pela minha vida, e em parte ao acaso de ter tido tudo isso. Claro que nesse percurso algumas pessoas merecem especial reconhecimento.

Começo pelo meu núcleo familiar, e como não poderia deixar de ser, agradeço a minha mãe, mais conhecida como Chuca, a dona, ou professora, Nazaré. Essa mulher inquieta que com sua energia inacreditável ainda baila alegremente em festas e comemorações como a que marcou a minha defesa de doutorado; Que me impregnou não só de valores como honestidade, sinceridade, desapego às coisas materiais, mas também de amor à cultura, à música, às pessoas e à vida. Poderia me estender por mais de uma légua nesta prosa gostosa, mas ainda tenho muita gente para agradecer. Agradeço ao meu pai, Thomas, que foi meu parceiro, meu amigo, minha inspiração maior quando o tema é otimismo ou entrega àquilo que se faz ou se ama. Como apenas um pai não é suficiente na vida, agradeço ao meu pai Miguel, pela vida, pelos sonhos, pelo amor. Aos meus irmãos..., são muitos! De cima para baixo: Miguel, o mais velho, Breno, Lucas, Andrea, Milena, Paola, Tácio e Luisa, a caçula. Agradeço a todos pelos mil momentos compartilhados, cumplicidade, parceria, brigas e tudo mais que bons irmãos podem desfrutar. De todos poderia destacar momentos importantes, mas agradeço especialmente ao Breno, que com absoluta certeza sempre foi a minha maior referência e fonte de inspiração no que se relaciona à tudo que envolve ciência ou matemática. Agradeço aos meus tios Lício e Luca por todas as aventuras, histórias, acolhimento e carinho. E a todos as minhas tias, primos e primas membros dessa grande e amorosa família. Agradeço a minha esposa Fernanda pelo incentivo incondicional, pela parceria e todo o amor que me deu e dá.

Fora do meu núcleo familiar, percebe-se que só pela quantidade de irmãos poderia discursar por longa data, também tenho muita gente para agradecer. Ao meu orientador José Carlos Moreira Bermudez, pela grande influência científica que ficará para sempre marcada em mim e em tudo que eu produzir academicamente daqui para frente. Ao prof. Cédric Richard pela fantástica contribuição, claro, com o trabalho e artigos, mas também com amizade sempre dotada de extraordinário humor e elegância. Ao prof. Márcio Costa eu agradeço pelas conversas, na hora do almoço, e principalmente as de outras horas, sempre ajudando, não só a mim, mas todos com quem tem contato, a suportar toda a pressão e ansiedade que

esta pequena empreitada nos exige. Muito obrigado, Márcio! Aos meus queridos colegas de laboratório Johan Malaver, Rafael Chiea e Fábio Iturriet, muito obrigado. Especialmente agradeço Juan Liber, Ricardo Borsoi, Renata Borges, Daniel Montezano e Marcos Maruo (vulgo, Chen) pela amizade, parceria, histórias e colaborações que vão do tecnicismo da nossa área a coisas menos palpáveis, de cunho mais emocional, e seladas com incontáveis cervejas. Agradeço também a todos os professores que de alguma forma participaram da construção desse trabalho. Agradeço aos professores Jean Yves Tourneret, Leonardo Resende, Alejandro Frery e Fermín Bazán. A todos, muito obrigado. Especialmente agradeço ao prof. Aldebaro Klautau Junior por ter sido o primeiro condutor da minha descoberta científica, pela amizade e grande fonte de inspiração.

Finalmente, não poderia deixar de agradecer ao Wilson da Silva Costa e ao Marcelo Siqueira do PPGEEL que dedicadamente tocam toda a burocracia do programa, da forma mais humana possível e com grande consideração por todos que ali passam. Muito obrigado mesmo! Eles são uma amostra das extraordinárias pessoas que fazem parte da Universidade Federal de Santa Catarina, em nome de quem agradeço enormemente essa oportunidade que me foi concedida.

É melhor pedir desculpas do que autorização.
(Marcos Hideo Maruo, 2014)

RESUMO ESTENDIDO

Imagem hiperespectral (HI) é uma imagem em que cada pixel contém centenas (ou até milhares) de bandas estreitas e contíguas amostradas num amplo domínio do espectro eletromagnético. Sensores hiperespectrais normalmente trocam resolução espacial por resolução espectral devido principalmente a fatores como a distância entre o instrumento e a cena alvo, e limitada capacidade de processamento, transmissão e armazenamento históricas, mas que se tornam cada vez menos problemáticas. Este tipo de imagem encontra ampla utilização em uma gama de aplicações em astronomia, agricultura, imagens biomédicas, geociências, física, vigilância e sensoriamento remoto. A usual baixa resolução espacial de sensores espectrais implica que o que se observa em cada pixel é normalmente uma mistura das assinaturas espectrais dos materiais presentes na cena correspondente (normalmente denominados de *endmembers*). Assim um pixel em uma imagem hiperespectral não pode mais ser determinado por um tom ou cor mas sim por uma assinatura espectral do material, ou materiais, que se encontram na região analisada.

O modelo mais simples e amplamente utilizado em aplicações com imagens hiperespectrais é o modelo linear, no qual o pixel observado é modelado como uma combinação linear dos *endmembers*. No entanto, fortes evidências de múltiplas reflexões da radiação solar e/ou materiais intimamente misturados, i.e., misturados em nível microscópico, resultam em diversos modelos não-lineares dos quais destacam-se os modelos *bilineares*, modelos de pós não-linearidade, modelos de mistura íntima e modelos não-paramétricos.

Define-se então o problema de *desmistura espectral* (ou em inglês *spectral unmixing* – SU), que consiste em determinar as assinaturas espectrais dos *endmembers* puros presentes em uma cena e suas proporções (denominadas de abundâncias) para cada pixel da imagem. SU é um problema inverso e por natureza cego uma vez que raramente estão disponíveis informações confiáveis sobre o número de *endmembers*, suas assinaturas espectrais e suas distribuições em uma dada cena. Este problema possui forte conexão com o problema de separação cega de fontes mas difere no fato de que no problema de SU a independência de fontes não pode ser considerada já que as abundâncias são de fato proporções e por isso dependentes (abundâncias são positivas e devem somar 1). A determinação dos *endmembers* é conhecida como *extração de endmem-*

bers e a literatura apresenta uma gama de algoritmos com esse propósito. Esses algoritmos normalmente exploram a geometria convexa resultante do modelo linear e da restrições sobre as abundâncias. Quando os *endmembers* são considerados conhecidos, ou estimados em um passo anterior, o problema de SU torna-se um problema supervisionado, com pares de entrada (*endmembers*) e saída (pixels), reduzindo-se a uma etapa de inversão, ou regressão, para determinar as proporções dos *endmembers* em cada *pixel*. Quando modelos não-lineares são considerados, a literatura apresenta diversas técnicas que podem ser empregadas dependendo da disponibilidade de informações sobre os *endmembers* e sobre os modelos que regem a interação entre a luz e os materiais numa dada cena. No entanto, informações sobre o tipo de mistura presente em cenas reais são raramente disponíveis. Nesse contexto, métodos kernelizados, que assumem modelos não-paramétricos, têm sido especialmente bem sucedidos quando aplicados ao problema de SU. Dentre esses métodos destaca-se o SK-Hype, que emprega a teoria de *mínimos quadrados-máquinas de vetores de suporte* (LS-SVM), numa abordagem que considera um modelo linear com uma flutuação não-linear representada por uma função pertencente a um espaço de Hilbert de *kernel* reprodutivos (RKHS). Nesta tese de doutoramento diferentes problemas foram abordados dentro do processo de SU de imagens hiperespectrais não-lineares como um todo. Contribuições foram dadas para a detecção de misturas não-lineares, estimação de *endmembers* quando uma parte considerável da imagem possui misturas não-lineares, e seleção de bandas no espaço de Hilbert de kernels reprodutivos (RKHS). Todos os métodos foram testados através de simulações com dados sintéticos e reais, e considerando *unmixing* supervisionado e não-supervisionado.

No Capítulo 4, um método semi-paramétrico de detecção de misturas não-lineares é apresentado para imagens hiperespectrais. Esse detector compara a performance de dois modelos: um linear paramétrico, usando mínimos-quadrados (LS), e um não-linear não-paramétrico usando processos Gaussianos. A idéia da utilização de modelos não-paramétricos se conecta com o fato de que na prática pouco se sabe sobre a real natureza da não-linearidade presente na cena. Os erros de ajuste desses modelos são então comparados em uma estatística de teste para a qual é possível aproximar a distribuição na hipótese de misturas lineares e, assim, estimar um limiar de detecção para uma dada probabilidade de falso-alarme. A performance do detector proposto foi estudada considerando problemas supervisionados e não-supervisionados, sendo mostrado que a melhoria obtida no desempenho SU utilizando o detector proposto é estatística-

mente consistente. Além disso, um grau de não-linearidade baseado nas energias relativas das contribuições lineares e não-lineares do processo de mistura foi definido para quantificar a importância das parcelas linear e não-linear dos modelos. Tal definição é importante para uma correta avaliação dos desempenhos relativos de diferentes estratégias de detecção de misturas não-lineares.

No Capítulo 5 um algoritmo iterativo foi proposto para a estimação de *endmembers* como uma etapa de pré-processamento para problemas SU não supervisionados. Esse algoritmo intercala etapas de detecção de misturas não-lineares e estimação de *endmembers* de forma iterativa, na qual uma etapa de estimação de *endmembers* é seguida por uma etapa de detecção, na qual uma parcela dos pixels “mais não-lineares” é descartada. Esse processo é repetido por um número máximo de execuções ou até um critério de parada ser atingido. Demonstra-se que o uso combinado do detector proposto com um algoritmo de estimação de *endmembers* leva a melhores resultados de SU quando comparado com soluções do estado da arte. Simulações utilizando diferentes cenários corroboram as conclusões.

No Capítulo 6 dois métodos para SU não-linear de imagens hiperespectrais, que empregam seleção de bandas (BS) diretamente no espaço de Hilbert de *kernels* reprodutivos (RKHS), são apresentados. O primeiro método utiliza o algoritmo *Kernel K-Means* (KKM) para encontrar *clusters* diretamente no RKHS onde cada centroide é então associada ao vetor espectral mais próximo. O segundo método é centralizado e baseado no critério de coerência, que incorpora uma medida da qualidade do dicionário no RKHS para a SU não-linear. Essa abordagem centralizada é equivalente a resolver um problema de máximo clique (MCP). Contrariamente a outros métodos concorrentes que não incluem uma escolha eficiente dos parâmetros do modelo, o método proposto requer apenas uma estimativa inicial do número de bandas selecionadas. Os resultados das simulações empregando dados, tanto sintéticos como reais, ilustram a qualidade dos resultados de *unmixing* obtidos com os métodos de BS propostos. Ao utilizar o SK-Hype, para um número reduzido de bandas, são obtidas estimativas de abundância tão precisas quanto aquelas obtidas utilizando o método SK-Hype com todo o espectro disponível, mas com uma pequena fração do custo computacional.

Palavras-chave: Imagem hiperespectral. Otimização. Detecção. Processo Gaussiano. Kernel. RKHS. Seleção de bandas.

ABSTRACT

Mixing phenomena in hyperspectral images depend on a variety of factors such as the resolution of observation devices, the properties of materials, and how these materials interact with incident light in the scene. Different parametric and nonparametric models have been considered to address hyperspectral unmixing problems. The simplest one is the linear mixing model. Nevertheless, it has been recognized that mixing phenomena can also be nonlinear. Kernel-based nonlinear mixing models have been applied to unmix spectral information of hyperspectral images when the type of mixing occurring in the scene is too complex or unknown. However, the corresponding nonlinear analysis techniques are necessarily more challenging and complex than those employed for linear unmixing. Within this context, it makes sense to search for different strategies to produce simpler and/or more accurate results. In this thesis, we tackle three distinct parts of the complete spectral unmixing (SU) problem. First, we propose a technique for detecting nonlinearly mixed pixels. The detection approach is based on the comparison of the reconstruction errors using both a Gaussian process regression model and a linear regression model. The two errors are combined into a detection test statistics for which a probability density function can be reasonably approximated. Second, we propose an iterative endmember extraction algorithm to be employed in combination with the detection algorithm. The proposed detect-then-unmix strategy, which consists of extracting endmembers, detecting nonlinearly mixed pixels and unmixing, is tested with synthetic and real images. Finally, we propose two methods for band selection (BS) in the reproducing kernel Hilbert space (RKHS), which lead to a significant reduction of the processing time required by nonlinear unmixing techniques. The first method employs the kernel k-means (KKM) algorithm to find clusters in the RKHS. Each cluster centroid is then associated to the closest mapped spectral vector. The second method is centralized, and it is based upon the coherence criterion, which sets the largest value allowed for correlations between the basis kernel functions characterizing the unmixing model. We show that the proposed BS approach is equivalent to solving a maximum clique problem (MCP), that is, to searching for the largest complete subgraph in a graph. Furthermore, we devise a strategy for selecting the coherence threshold and the Gaussian kernel bandwidth using coherence bounds for linearly independent

bases. Simulation results illustrate the efficiency of the proposed method.

Keywords: Hyperspectral Images. Optimization. Detection. Gaussian Process. Kernel. RKHS. Band Selection.

LIST OF FIGURES

Figure 1	Remote sensing.	31
Figure 2	An observed pixel is in fact a mixture of spectral signatures.	31
Figure 3	Graphical representation of the hypercube collected by the AVIRIS instrument from the Cuprite mining district [1].	32
Figure 4	Different forms of Solar interaction.	39
Figure 5	Simplex.	41
Figure 6	ROC Curve. Different detectors presented in different colors.	65
Figure 7	DC embedded in noise.	67
Figure 8	Hypothesis PDFs.	67
Figure 9	Test statistic PDFs. The shaded area corresponds to the probabilities of detection (light gray) and false alarm (darker gray).	69
Figure 10	ROC Curve.	69
Figure 11	Empirical ROCs for: (a) the Robust LS detector [2], (b) the proposed GP detector, (c) the two detectors for $\eta_d = 0.5$. All curves were obtained for 8000 pixels (4000 linearly mixed and 4000 nonlinearly mixed) and SNR = 21dB. Nonlinear mixtures were generated using the simplified GBM described in Section 4.2.1.	80
Figure 12	Histograms for (a) the squared norm of the GP fitting error, (b) the least-squares fitting error, and (c) the test statistics (4.20).	81
Figure 13	Histogram of the test statistics under \mathcal{H}_0 and the adjusted Beta distribution.	82
Figure 14	ROCs for different proportions of nonlinearly mixed pixels and $\eta_d = 0.5$. Endmember extraction using VCA.	85
Figure 15	ROCs for different degrees of nonlinearity η_d and 50% of nonlinearly mixed pixels in the image. Endmember extraction using Algorithm 1.	89
Figure 16	Graphical illustration of the endmember estimation process using the proposed iterative algorithm. The data set consists of 2000 pixels, with a proportion of 50% nonlinearly mixed pixels obtained with the GMB model and $\eta_d = 0.5$. Green dots are the current estimated endmembers, and black dots are the data projected onto the subspace spanned by the columns of the current matrix \mathbf{M} . The true endmembers are shown as black circles at the vertices of the true simplex drawn with black lines. The data discarded at the corresponding iteration are shown within blue circles.	91

Figure 17 Cuprite mining site. The green box corresponds to the alunite hill scene.....	94
Figure 18 (a) Plot of the alunite hill with bands 30, 70 and 100. (b) Reconstruction of the scene using the LMM. (c) Adding 30 % of nonlinearly mixed pixels and WGN to give a 30dB SNR.....	94
Figure 19 Endmember estimations for the nonlinearly mixed image with different extraction techniques.....	94
Figure 20 The black circles are the real endmembers, the black dots are the data projected in the columns of \mathbf{M} . The blue circles are the estimated endmembers with the proposed algorithm after 10 iterations. The simplex for the “true” and estimated endmembers are also drawn.	95
Figure 21 Detection map and true nonlinear map. Linearly mixed pixels in gray, nonlinearly mixed pixels in white, and misclassified pixels in black.	95
Figure 22 Indian Pines test site representation selecting 3 different bands in (a), and 3 other bands in (b).	97
Figure 23 Detection of nonlinearly mixed pixels in Indian Pines hyperspectral image. Black pixels were detected as nonlinearly mixed ones by the proposed detector.	98
Figure 24 Estimated endmembers and USGS spectra.....	101
Figure 25 Abundance maps.....	102
Figure 26 Cuprite scene and reconstruction errors.....	102
Figure 27 The maximum clique problem (MCP).....	109
Figure 28 Pavia University. In the left, the Pavia University HI is represented using the bands 5, 30, and 50. In the right, the classified areas are labelled from 1 to 9, while 0 corresponds to unclassified areas.	119
Figure 29 Cuprite scene used in [3].....	120

LIST OF TABLES

Table 1	Characteristics of 8 hyperspectral instruments.	30
Table 2	Detection threshold.	77
Table 3	Abundance estimation RMSE for \mathbf{M} known and using the GBM mixing model (SNR = 21dB, $\eta_d = 0.5$).	83
Table 4	Abundance estimation RMSE for \mathbf{M} known and using the PNMM mixing model (SNR = 21dB, $\eta_d = 0.5$).	83
Table 5	One-tailed Wilcoxon signed rank test for Image I (Significance level 0.05).	84
Table 6	One-tailed Wilcoxon signed rank test for Image II (Significance level 0.05).	84
Table 7	Mean RMSE for endmember estimation.	92
Table 8	RMSE for the abundances in the alunite hill scene.	93
Table 9	Indian Pines classes by region.	98
Table 10	Subimages organization	99
Table 11	Indian Pines reconstruction error (RMSE) by subimage.	99
Table 12	Spectral angles (in rad) between estimated and USGS spectra.	101
Table 13	RMSE. 100 runs, 2000 pxl., 8 endmembers (Cuprite), SNR=21dB, GBM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).	117
Table 14	RMSE. 100 runs, 2000 pxl., 8 endmembers (Cuprite), SNR=21dB, PNMM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).	118
Table 15	RMSE. 100 runs, 2000 pxl., 8 endmembers (Pavia), SNR=21dB, GBM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).	121
Table 16	RMSE. 100 runs, 2000 pxl., 8 endmembers (Pavia), SNR=21dB, PNMM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).	122
Table 17	Cuprite image. RMSE between the abundances estimated with SK-Hype (all bands) and BS + SK-Hype.	123
Table 18	Pavia University image. RMSE between the abundances estimated with SK-Hype (all bands) and BS + SK-Hype.	124

LIST OF ABBREVIATIONS AND ACRONYMS

SU	Spectral unmixing	29
HI	Hyperspectral images	29
HI	Hyperspectral image	31
GP	Gaussian process	33
MCP	Maximal clique problem	34
LMM	Linear mixing model	38
WGN	White Gaussian noise	40
GBM	Generalized bilinear model	43
PNMM	Post-nonlinear mixing model	45
EEA	Endmember extraction algorithm	47
LS	Least squares	48
FCLS	Fully constrained least squares	49
PDF	Probability density function	50
BS	Band selection	53
RKHS	Reproducing kernel Hilbert space	61
PFA	Probability of false alarm	64
PD	Probability of detection	64
ROC	Receiver operating characteristic	65
RMSE	Root mean square error	82
KKM	Kernel k -means	105

LIST OF SYMBOLS

M	Mixing matrix	37
m_i	i th column of M	37
$m_{\lambda_\ell}^\top$	ℓ th line of M	37
r	Observerd pixel	37
L	Number of spectral bands	37
R	Number of endmembers	37
Ψ	Unknown mixing function	37
n	Additive noise vector	37
α	Abundance vector	38
\mathcal{N}	Gaussian PDF	40
I	Identity matrix	40
\mathcal{S}_α	Simplex formed by the constraints over the abundances	40
\mathcal{S}_r	Simplex formed by the endmembers	40
$\gamma_{i,j}$	Parameter governing the amount of nonlinear contribution for the GBM	43
R	Hyperspectral image	47
A	Abundance matrix	47
ψ	Mixing function in RKHS	55
$\kappa(\cdot, \cdot)$	Reproducing kernel	55
\mathcal{H}	RKHS	56
σ^2	Kernel bandwidth	57
σ_n^2	WGN power	58
\mathcal{L}	Lagrangian function	59
T	Detection test statistic	64
τ	Detection threshold	64
\mathcal{H}_0	Null hypothesis	66
\mathcal{H}_1	Alternative hypothesis	66
e_{lin}	Linear estimator error	72
χ	Chi-square distribution	73
e_{nonlin}	GP estimation error	74
\mathcal{D}	Dictionary of kernel functions	105
\mathcal{I}_D	Dictionary index set	105

M	Size of the dictionary	105
\mathcal{C}	Cluster	105
μ	Coherence	107
μ_0	Coherence parameter.....	107

CONTENTS

1	INTRODUCTION	29
1.1	CONTRIBUTION	33
1.2	ORGANIZATION OF THE DOCUMENT	35
2	STATE-OF-THE-ART	37
2.1	NOTATION	37
2.2	SPECTRAL MIXING MODELS	37
2.2.1	Linear Mixing Model	38
2.2.1.1	Geometry of the LMM	40
2.2.2	Bilinear Mixing Models	40
2.2.2.1	Fan's Model	42
2.2.2.2	Nascimento's Model	42
2.2.2.3	Generalized Bilinear Model	43
2.2.3	Post-Nonlinear Mixing Model - PNMM	45
2.2.4	Intimate Mixing Models	46
2.3	SPECTRAL UNMIXING: PROBLEM DEFINITION	47
2.3.1	Endmembers Estimation	47
2.3.2	Linear SU	48
2.3.3	Nonlinear SU	50
2.3.3.1	Nonlinear SU using Parametric Models	50
2.3.3.2	Model-free Nonlinear SU	51
2.4	DETECTION OF NONLINEARLY MIXED PIXELS	52
2.5	BAND SELECTION	53
3	PRELIMINARY THEORETICAL CONCEPTS	55
3.1	KERNEL REGRESSION	55
3.1.1	Mercer kernels and RKHS	55
3.1.2	Gaussian Process for regression	57
3.1.3	LS-SVR	59
3.1.4	SK-Hype	60
3.1.4.1	Solving with respect to ψ	61
3.1.4.2	Solving with respect to u	63
3.2	BASIC DETECTION CONCEPTS	63
3.2.1	The detection problem	64
3.2.2	Example: Detection of a DC level embedded in Gaussian noise	66
4	NONLINEAR MIXTURE DETECTOR	71
4.1	DETECTION OF NONLINEARLY MIXED PIXELS	72
4.1.1	The detection problem	72

4.1.2	Linear estimation error	72
4.1.3	Nonlinear estimation error with GP	74
4.1.4	The test statistics	75
4.1.5	Determining the detection threshold	76
4.2	SIMULATIONS	76
4.2.1	Degree of nonlinearity	77
4.2.1.1	Synthetic data generation with GBM	78
4.2.2	Synthetic data generation with PNMM	79
4.2.3	Simulations with known M	79
4.2.4	Simulations with an unknown endmember matrix M	84
4.3	PRELIMINARY CONCLUSIONS	85
5	EEA FOR NONLINEARLY MIXED HYPERSPECTRAL IMAGES	87
5.1	ENDMEMBER EXTRACTION IN NONLINEARLY MIXED HYPERSPECTRAL IMAGES	87
5.2	SIMULATIONS	89
5.2.1	Simulations with an unknown endmember matrix M	89
5.2.2	Choosing the parameters r_f, N_{\max}, and ε	90
5.2.3	Simulation with synthetic data extracted from a real scene	92
5.2.4	Real Data	96
5.2.4.1	Indian Pines	96
5.2.4.2	Cuprite	100
5.3	PRELIMINARY CONCLUSIONS	101
6	BAND SELECTION IN RKHS	103
6.1	REVISITING THE KERNEL FRAMEWORK	104
6.2	KERNEL K-MEANS FOR BAND SELECTION	105
6.3	COHERENCE-BASED BAND SELECTION	106
6.3.1	Coherence criterion for dictionary selection	106
6.3.2	Band selection as a maximum clique problem	108
6.3.3	The maximum clique problem	108
6.3.4	Coherence-based BS algorithms	110
6.3.4.1	Automatic parameter settings	110
6.3.4.2	Algorithms	111
6.4	SIMULATIONS	114
6.4.1	Simulation with synthetic data	114
6.4.2	Simulation with real data	116
6.5	PRELIMINARY CONCLUSIONS	120
7	CONCLUSIONS	125
7.1	FUTURE WORK	126
7.2	PUBLICATIONS	126
7.3	SOURCE CODE	127

APPENDIX A – Convex Optimization in RKHS	131
APPENDIX B – Gaussian Process Regression	151
REFERENCES.....	157

1 INTRODUCTION

Emerged in the 1960s with the first multispectral scanners, the spectral unmixing problem (SU) [4, 5] consists of identifying target materials and estimating their proportions in a given scene. Instruments capable of sampling contiguously a wide range of the solar radiation produce two-dimensional images called *hyperspectral images* (HIs). Due to historic downlink and computer processing limitations, specially in the early 1970s [4], hyperspectral images often trade spatial for spectral resolution [5]. Such trade-off is especially evident in remote sensing applications, and are caused by the large distance between sensors and target scenes. The observed reflectances then result from spectral mixtures of several pure material signatures. As a consequence, spectral unmixing has become an important issue for hyperspectral data processing [6], and is still a hot topic today, see [7, 8] and references there in.

Modern instruments produce HIs with tens to hundreds of narrow contiguous bands, in an increasingly wide portion of the spectrum, ranging from the visible light to the far-infrared [7]. Clearly, increasing the number of bands, i.e., using higher spectral resolution, results in a linear increase of the amount of data and, consequently, a proportional increase in the complexity of processing algorithms. This contrasts with the polynomial data growth that would result from increasing the spatial resolution of a given image. In HIs, a pixel cannot be identified by its color or tone, but by its spectral signature containing several (usually hundreds) of samples in different wavelengths. The low spatial resolution implies that the spectrum observed in a given pixel is often a mixture of the spectral signatures of the materials present in the scene. Figure 1 illustrates the acquisition process of HIs in which the solar radiation is reflected by the materials in the Earth's surface and is measured by the hyperspectral sensor in a satellite (spaceborne sensor). Figure 2 shows how an observed pixel can be seen as a mixture of spectral signatures.

Table 1, replicated from [7], presents characteristics of 8 hyperspectral instruments that are airborne (HYDICE and AVIRIS) and spaceborne (HYPERION EnMAP, PRISMA, CHRIS, HypsIRI, and IASI). These instruments differ in their operation altitude, spatial and spectral resolutions, number of pixels, number of bands, and spectral range. The spectral range for HYDICE, AVIRIS, HYPERION, EnMAP, PRISMA, and HypsIRI covers the visible, near-infrared, and the short-wave infrared spectra, while CHRIS covers only the visible and the IASI covers the mid- and far-infrared spectral regions.

Hyperspectral images represent a target scene with $n_1 \times n_2$ pixels, where each pixel has L contiguous narrow bands. Thus, a hyperspectral im-

Table 1: Characteristics of 8 hyperspectral instruments.

Parameters	HYDICE	AVIRIS	HYPERION	EnMAP	PRISMA	CHRIS	HypSIRI	IASI
Altitude (<i>km</i>)	1.6	20	705	653	614	556	626	817
Spatial resol. (<i>m</i>)	0.75	20	30	30	5-30	36	60	V: 1.2 <i>km</i> H: 25 <i>km</i>
Spectral resol. (<i>nm</i>)	7-14	10	10	6.5-10	10	1.3-12	4-12	0.5 <i>cm</i> ⁻¹
Spectral range (<i>μm</i>)	0.4-2.5	0.4-2.5	0.4-2.5	0.4-2.5	0.4-2.5	0.4-1.0	0.38-2.5 and 7.5-12	3.62-15.5 (645-2760 <i>cm</i> ⁻¹)
# of bands	210	224	220	228	238	63	217	8461
Spectral cube (lines× columns×bands)	200×320 ×210	512×614 ×224	660×236 ×220	1000×1000 ×228	400×880 ×238	748×748 ×63	620×512 ×210	765×120 ×8461

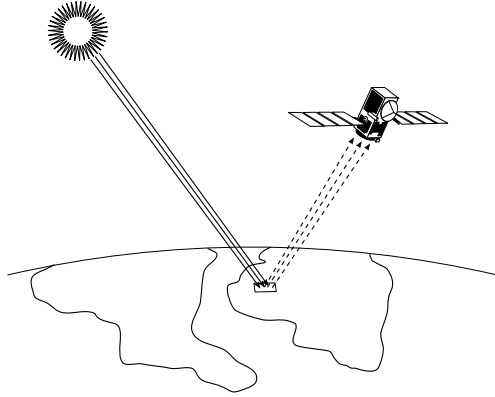


Figure 1: Remote sensing.

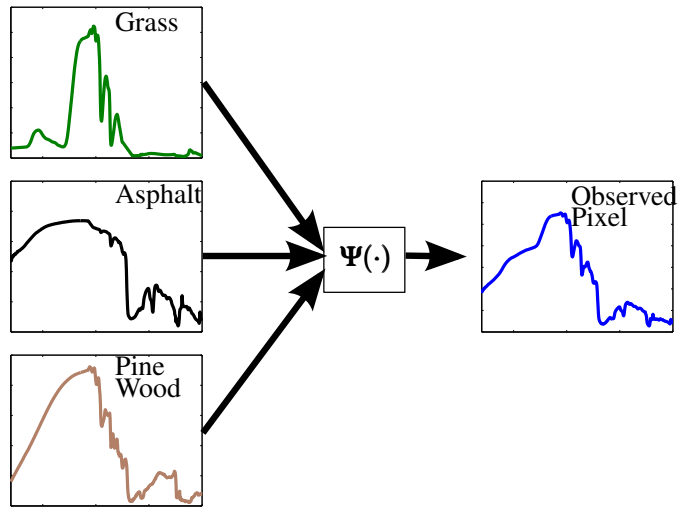


Figure 2: An observed pixel is in fact a mixture of spectral signatures.

age (HI) can be seen as L grayscale images of the same target, that is, one image for each of the L spectral bands. These images are usually stacked and presented in a three-dimensional hypercube $\mathcal{B} \in \mathbb{R}^{n_1 \times n_2 \times L}$. Figure 3 left shows a graphical representation of a hypercube captured by the AVIRIS hyperspectral sensor from the Cuprite mining district in Nevada-USA. Each of the 512×614 pixels collected in this scene has 224 spectral bands ranging

from 0.7 to 2.5 micrometers. Figure 3 right plots the spectral signature of the pixel indicated with a black square.

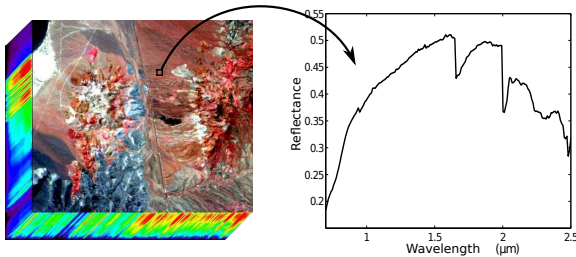


Figure 3: Graphical representation of the hypercube collected by the AVIRIS instrument from the Cuprite mining district [1].

In the specific parlance of the field, the materials in a given scene are called *endmembers*, and their proportions are named abundances. Now, we can divide the spectral unmixing problem in three basic problems:

- 1) define a mathematical model to govern the interactions between the light beams and the endmembers (mixture) present in a scene;
- 2) determine the endmembers (i.e., the spectral signature of the pure materials) in this scene;
- 3) estimate the abundances of each endmember for all pixels.

Several mathematical models were proposed in the literature in the past few years concerning different forms of endmember interaction, single or multiple reflections of solar light beams, size of endmembers, land relief, etc. Such models are usually classified into linear and nonlinear *mixing models* [6, 8] and the most relevant ones will be discussed in Section 2.2. The second problem is usually addressed by endmember extraction techniques, which are discussed in Section 2.3.1. The third problem is often called *inversion* since it falls in the category of inverse problems, and a more formal definition is given in Section 2.3.

The solution of the problems (1-3) presented above is greatly hampered by factors such as spectral variability [9, 10] (which can be seasonal or along the same image where the spectrum of an endmember slightly varies in different parts of the scene), illumination effects and solar incidence angle, atmospheric interference, and instrument calibration. Hyperspectral sensors usually measure the incident radiance of the scene being analyzed. Such information is then converted to surface apparent reflectance as a prior step to

unmixing [11], since SU algorithms are concerned with unveiling spectral and abundance information in the ground. The hampering factors mentioned above and the radiance-to-reflectance conversion are out of the scope of the present work. Nevertheless, we invite the interested reader to refer to [7] and references there in for more details.

This work focuses in nonlinear spectral unmixing, targeting three specific problems: *nonlinear mixture detection*, *endmember estimation* in nonlinearly mixed HIs, and *band selection* for nonparametric kernel based SU of nonlinearly mixed HIs.

1.1 CONTRIBUTION

In this thesis we deal with the problem of nonlinear spectral unmixing of hyperspectral images. We propose new approaches for detecting nonlinearly mixed pixels, endmember estimation for HIs with nonlinearly mixing, and band selection for kernelized unmixing methods.

The detector of nonlinearly mixed pixels discussed in this work was initially proposed in [12], and [3]. We use both least-squares and Gaussian processes (GP) to model the unknown mixing process occurring in Nature. The reconstruction errors for both methods are then combined into a novel test statistics for which a probability density function can be reasonably approximated. The proposed detector, namely GP detector for short, is nonparametric, and little is assumed regarding the type of nonlinearity occurring in the mixing process. Simulations show that the GP detector outperforms other parametric and nonparametric detectors found in the literature. However, it is verified that poorly estimated endmembers lead to degraded detection and unmixing performance.

The problems of extracting endmembers, detecting nonlinearly mixed pixels and unmixing are interlaced, and addressing them jointly is not a trivial task. For instance, most nonlinear unmixing techniques assume the endmembers to be known or to be estimated by an endmember extraction algorithm (EEA) [13, 14, 15, 16, 17, 18, 19, 20, 21]. However, most endmember extraction algorithms rely on the convex geometry associated with the linear mixing model [22, 23, 24, 25, 26], which obviously does not apply to nonlinearly mixed pixels. Endmember extraction techniques designed for situations where a significant part of the image is composed of nonlinear mixtures are rarely addressed in the literature. In fact, most of the techniques considering nonlinearly mixed pixels are part of a completely unsupervised unmixing strategy [27, 28]. Thus, we propose an (*Minimum Volume Enclosing Simplex*) MVES-based iterative endmember extraction algorithm to be employed

in combination with the GP detection algorithm to jointly detect nonlinearly mixed pixels and extract the image endmembers from the linearly mixed pixels. Simulations demonstrate the effectiveness of the method for improving detection and endmember estimation performances. This jointly iterative approach was also published in [3].

One of the problems in practical implementation of nonlinear unmixing algorithms is the profusion of spectral bands generated in the acquisition process, leading to high computational costs. Such inherent complexity, associated with the high redundancy within the complete set of bands, make the search of band selection techniques natural and relevant [29]. When considering kernel methods, the data is mapped to a high-dimensional reproducing kernel Hilbert space (RKHS) where the problem is solved linearly. Thus, selecting bands directly in the RKHS has shown to be quite effective in both reducing the complexity and preserving the accuracy. We propose two different approaches for selecting bands in the RKHS. The first [30], applies a kernel k -means algorithm to identify nonlinearly separable clusters of spectral bands in the corresponding RKHS. The second [31], formulates the band selection (BS) problem as a maximal clique problem (MCP) [32, 33], using the coherence criterion as a similarity measure among the mapped samples.

Briefly, the main contributions of this work are the following:

- a) a model-free detector of nonlinearly mixed pixels. The novel test statistics compares reconstruction errors of the observations modeled by a Gaussian Process and a linear regression;
- b) a novel recursive endmember estimation algorithm for scenes that are partly nonlinear;
- c) the definition of a degree of nonlinearity η_d which allows a meaningful comparison of detection results for images obtained using different mixing models;
- d) a kernel k -means based BS strategy;
- e) an MCP centralized strategy to perform BS using the coherence criterion;
- f) a meaningful methodology to select the coherence threshold and kernel parameter when considering the MCP for BS and unmixing.

For replicability, Matlab source codes (and datasets) which replicate the simulations presented in this work are available at https://github.com/talesim/NP_NL_Det_EE_HI/archive/master.zip (Chapters 4 and 5) and https://github.com/talesim/cliqeu_BS/archive/master.zip (Chapter 6).

1.2 ORGANIZATION OF THE DOCUMENT

Chapter 2 presents a more detailed description of the hyperspectral unmixing problem, and reviews the state of the art models and methods. In Chapter 3 preliminary theoretical concepts needed in Chapters 4–6 are presented. In Chapter 4 the proposed detection strategy is presented, while the joint endmember extraction and detection approach is discussed in Chapter 5. Both band selection strategies are presented in Chapter 6, and the work is concluded in Chapter 7. All chapters include detailed simulations to test and illustrate the application of the proposed methods to synthetic and real hyperspectral images. We also present two appendices to complement needed mathematical background. In Appendix A, the convex optimization in RKHS is discussed, and relevant theorems and definitions are presented. Finally, Appendix B is dedicated to a more in-depth Gaussian process application to nonlinear regression. Both appendices present simple examples to motivate the reader.

2 STATE-OF-THE-ART

This introductory chapter reviews the state of the art for the hyperspectral unmixing problem, and is organized as follows: in Section 2.1 the basic mathematical notation used is presented. In Section 2.2 the most relevant mixing models are discussed, while SU is defined in Section 2.3. In Section 2.4 the nonlinear mixture detector problem is presented and discussed.

2.1 NOTATION

In this work column vectors are represented as small bold letters such as \mathbf{x} , for which the i -th component is represented by x_i . Bold capital letters such as \mathbf{X} represent matrices with components $\mathbf{X}_{i,j}$, unless defined otherwise locally. Functions with scalar output are represented by letters in Latin or Greek alphabets, e.g., $f(\cdot)$ or $\psi(\cdot)$. Functions with vectorial output are represented by bold letters such as $\mathbf{f}(\cdot)$ or $\Psi(\cdot)$. The endmember matrix, i.e. the matrix containing the spectral signature of all endmembers, is represented by \mathbf{M} . The i -th column of \mathbf{M} is represented by a vector \mathbf{m}_i (one for each endmember) and is the spectral signature of the i -th endmember. An alternative notation considers the row vectors of \mathbf{M} , $\mathbf{m}_{\lambda_i}^\top$ (one vector for each wavelength λ_i).

2.2 SPECTRAL MIXING MODELS

As previously mentioned, in HIs each pixel can be represented as a mixture of the endmember spectra present in the scene. This section presents some of the most relevant mixing models considered in the literature of SU [6, 8, 7]. However, we first consider a general formulation in which a pixel \mathbf{r} can be represented as

$$\mathbf{r} = \Psi(\mathbf{M}) + \mathbf{n}, \quad (2.1)$$

where $\mathbf{r} = [r_1, \dots, r_L]^\top$ is a vector with L spectral components, $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_R]$ is a $L \times R$ matrix that contains the endmember spectral signatures \mathbf{m}_i , Ψ is an unknown function defining the interaction between the solar radiation and endmembers of \mathbf{M} , and \mathbf{n} is an independent additive noise assigned to unmodeled parts of the system. Different models considered in SU of HIs differ in the linearity (or not) of Ψ , and in the rôle of the abundances in the model.

Another more specifically parametrized modeling of (2.1) can be written as

$$\mathbf{r} = \Psi(\mathbf{M}, \boldsymbol{\alpha}) + \mathbf{n}, \quad (2.2)$$

where $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_R]^\top$ is the abundance vector, and α_i is the abundance of the i -th endmember. Abundances are frequently defined to represent the proportions of the contribution of each endmember to the total observed reflectance. Thus, abundances cannot be negative and must sum to one

$$\sum_{i=1}^R \alpha_i = 1, \quad \text{and} \quad \alpha_i \geq 0, \quad i = 1, \dots, R. \quad (2.3)$$

The physical phenomenon governing the interaction between sunlight and various surface materials is naturally nonlinear. The nature of this nonlinearity is mainly linked to multiple reflections and to the transmission mechanisms of light [34]. Figure 4 illustrates the three principal types of light-endmember interaction considered in modern HI processing [8]. In Figure 4a, each sunlight beam interacts with only one material resulting in the linear mixing model (LMM) [6]. In Figure 4b the incident beam interacts with multiple elements that happen to be intimately mixed, giving rise to the so-called intimate mixing model [35, 36]. Multiple reflection, also named multiple scattering, is illustrated in Figure 4c where each light beam can be reflected by more than one different material modifying the observed electromagnetic spectrum [37]. Multiple scattering has been considered in a variety of mixing models [37, 8].

Next, models that are pertinent to this work will be formally defined.

2.2.1 Linear Mixing Model

The simplest and most common model in SU is the linear mixing model (LMM) [6]. The LMM considers that the observed pixel \mathbf{r} is modeled as a linear combination of the endmember spectra plus an additive noise. Thus,

$$\begin{aligned} \mathbf{r} &= \mathbf{M}\boldsymbol{\alpha} + \mathbf{n}. \\ \text{subject to } & \sum_{i=1}^R \alpha_i = 1, \text{ and } \alpha_i \geq 0, \forall i \in \{1, \dots, R\}. \end{aligned} \quad (2.4)$$

The additive noise \mathbf{n} is usually modeled as a vector of zero-mean, uncorrelated jointly Gaussian random variables with variance σ_n^2 , and indepen-

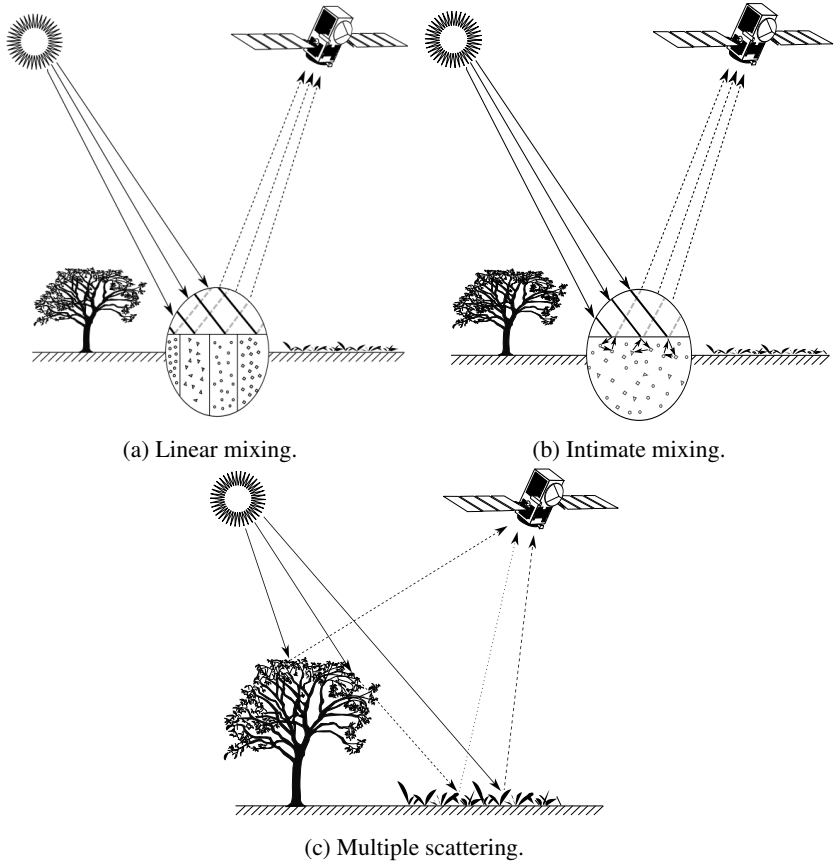


Figure 4: Different forms of Solar interaction.

dent of the endmembers \mathbf{m}_i . Thus, $\mathbf{n} \sim \mathcal{N}(0, \sigma_n^2 \mathbf{I})$, where σ_n^2 is the noise power and \mathbf{I} is the identity matrix. This type of noise is often referred as *white Gaussian noise* (WGN). The observation r_ℓ in the ℓ -th wavelength of (2.4) can be written as

$$r_\ell = \mathbf{m}_{\lambda_\ell}^\top \boldsymbol{\alpha} + n_\ell \quad (2.5)$$

where $\mathbf{m}_{\lambda_\ell}$ denotes the ℓ -th row of \mathbf{M} written as a column vector.

It is important to highlight that when working with HIs, the number of spectral bands L is much larger than the number of endmembers R , that is, $L \gg R$.

2.2.1.1 Geometry of the LMM

The constraints over the abundances considered in (2.4) define the simplex

$$\mathcal{S}_\alpha = \{\boldsymbol{\alpha} \in \mathbb{R}^R \mid \boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\alpha}^\top \mathbf{1} = 1\} \quad (2.6)$$

represented in Figure 5a. It is clear that the vectors in \mathcal{S}_α are confined to a $(R-1)$ -dimensional subspace due to the linear dependence among the abundance components, i.e., $\alpha_i = 1 - \sum_{j \neq i} \alpha_j$. Note, however, that considering (2.4) \mathbf{r} is a linear combination of the columns of \mathbf{M} , where the linear coefficients are the constrained abundances. Thus, all observations (neglecting the noise) are also confined to a simplex

$$\mathcal{S}_r = \{\mathbf{r} \in \mathbb{R}^L \mid \mathbf{r} = \mathbf{M}\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\alpha}^\top \mathbf{1} = 1\} \quad (2.7)$$

whose vertices are the endmembers. Then, it is clear that all the observed pixels also lie in an $(R-1)$ -dimensional space. This data simplex is illustrated in Figure 5b. It is important to emphasize that the convex geometry of the LMM is extensively exploited by a variety of algorithms proposed to solve the SU problem.

Clearly, the LMM neglects nonlinear interactions among endmembers as well as any other form of nonlinearity possibly present in the system. Other mixing models considers different types of nonlinearity, always trading between physical significance and mathematical tractability.

2.2.2 Bilinear Mixing Models

When considering multiple reflections of the solar light beam over the endmembers (see Figure 4c), we enter the realm of nonlinear mixing models.

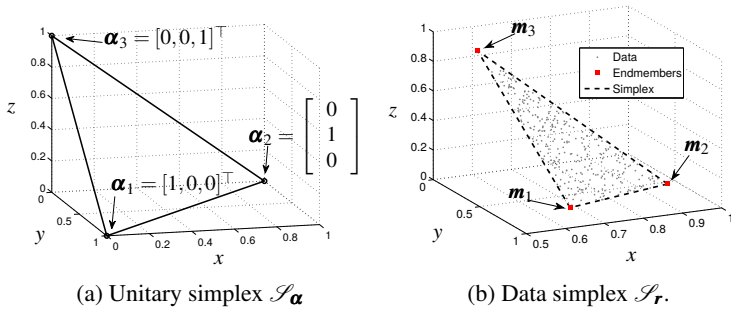


Figure 5: Simplex.

Bilinear models consider up to second order interactions, i.e., when the light beam interacts with up to two endmembers. In a general formulation, bilinear models represent a pixel \mathbf{r} as [8]

$$\mathbf{r} = \mathbf{f}(\mathbf{M}, \boldsymbol{\alpha}) + \mathbf{n}, \quad (2.8)$$

where

$$\mathbf{f}(\mathbf{M}, \boldsymbol{\alpha}) = \sum_{k=1}^R \alpha_k \mathbf{m}_k + \sum_{i=1}^{R-1} \sum_{j=i+1}^R \beta_{i,j} \mathbf{m}_i \odot \mathbf{m}_j, \quad (2.9)$$

and \odot represents the Hadamard (element-by-element) product

$$\mathbf{m}_i \odot \mathbf{m}_j = \begin{bmatrix} m_{1,i} \\ \vdots \\ m_{L,i} \end{bmatrix} \odot \begin{bmatrix} m_{1,j} \\ \vdots \\ m_{L,j} \end{bmatrix} = \begin{bmatrix} m_{1,i} m_{1,j} \\ \vdots \\ m_{L,i} m_{L,j} \end{bmatrix}.$$

In (2.9) the first term on the right side, also present in (2.4), corresponds to the linear parcel of the signal arriving at the sensor, while the second term is related to the nonlinear multiple scattering phenomenon. The coefficient $\beta_{i,j}$ governs the amount of nonlinear contribution of the interaction between \mathbf{m}_i and \mathbf{m}_j .

Different bilinear models assume different forms and constraints over the coefficients $\beta_{i,j}$, and abundances.

2.2.2.1 Fan's Model

The model proposed by W. Fan *at al.* [38] in 2009 is a bilinear model as presented in (2.9), in which the coefficients $\beta_{i,j}$ are the product of the abundances of the endmembers \mathbf{m}_i and \mathbf{m}_j , that is

$$\beta_{i,j} = \alpha_i \alpha_j.$$

Thus,

$$\mathbf{f}(\mathbf{M}, \boldsymbol{\alpha}) = \sum_{k=1}^R \alpha_k \mathbf{m}_k + \sum_{i=1}^{R-1} \sum_{j=i+1}^R \alpha_i \alpha_j \mathbf{m}_i \odot \mathbf{m}_j, \quad (2.10)$$

while keeping the constraints over the abundances as given in (2.3) and repeated here for convenience

$$\sum_{i=1}^R \alpha_i = 1, \quad \text{and} \quad \alpha_i \geq 0, \quad i = 1, \dots, R.$$

The mathematical expression of the Fan model is not physically motivated. It is derived from a polynomial approximation, where the first order terms in the Taylor expansion series were considered, leading to the coefficient products $\alpha_i \alpha_j$ in the second term of (2.10).

One problem of using this model for unmixing purposes is that the simplex, whose extremities are the endmembers, is defined only for the linear term. The nonlinear term can place the result of the mixture anywhere in the vector space.

2.2.2.2 Nascimento's Model

The nonlinear mixing model proposed in [39] generalizes Fan's model. The model has the same form presented in (2.9)

$$\mathbf{f}(\mathbf{M}, \boldsymbol{\alpha}) = \sum_{k=1}^R \alpha_k \mathbf{m}_k + \sum_{i=1}^{R-1} \sum_{j=i+1}^R \beta_{i,j} \mathbf{m}_i \odot \mathbf{m}_j, \quad (2.11)$$

but here $\beta_{i,j}$ is necessarily not given by the product of the abundances α_i, α_j . Applying the same concepts as in [37], Nascimento assumes the following coefficient constraints

- Positivity constraint: $\alpha_k \geq 0, \beta_{i,j} \geq 0$ for $\forall k$ and $\forall (i, j)$

- Sum-to-one constraint: $\sum_{k=1}^R \alpha_k + \sum_{i=1}^{R-1} \sum_{j=i+1}^R \beta_{i,j} = 1$.

Note that it makes sense to impose this sum-to-one constraint, so that each component of \mathbf{r} includes the total energy captured by the sensor in the corresponding wavelength.

This is a very tractable model because it preserves the simplex formed by the endmembers, which was an interesting characteristic of the linear mixing model. The nonlinear terms are incorporated to the model as “new endmembers” (the Hadamard product between original endmembers). Considering the pure endmembers and the new (mixed) endmembers, the unmixing problem can be treated as a linear unmixing problem. Thus, if we write $\tilde{\mathbf{M}} = [\mathbf{m}_1, \dots, \mathbf{m}_R, \mathbf{m}_1 \odot \mathbf{m}_2, \dots, \mathbf{m}_{R-1} \odot \mathbf{m}_R]$, with \tilde{R} the number of columns of $\tilde{\mathbf{M}}$, and $\tilde{\boldsymbol{\alpha}} = [\alpha_1, \dots, \alpha_R, \beta_{1,2}, \dots, \beta_{R-1,R}]^\top$, then Nascimento’s model can be written as

$$\begin{aligned} \mathbf{r} &= \tilde{\mathbf{M}} \tilde{\boldsymbol{\alpha}} + \mathbf{n} \\ \text{subject to } \sum_{k=1}^{\tilde{R}} \tilde{\alpha}_k &= 1, \quad \tilde{\alpha}_k \geq 0, \quad \forall k. \end{aligned} \quad (2.12)$$

Just as Fan’s model, this model was designed to handle multiple interaction between the solar radiation and the endmembers. Thus, in principle, this is not a suitable model for modelling intimate mixtures.

An issue about this model is that it assumes the previous knowledge of the pure endmember signatures prior to the SU in order to “build” the mixed endmembers.

2.2.2.3 Generalized Bilinear Model

The generalized bilinear model (GBM) [18] was proposed as a generalization of Nascimento’s model. It has the same bilinear form as presented in (2.9), and coefficients $\beta_{i,j}$ given by

$$\beta_{i,j} = \gamma_{i,j} \alpha_i \alpha_j,$$

where, the parameters $\gamma_{i,j} \in [0, 1]$ govern the amount of nonlinear contribution. The constraints over abundances are kept as presented in (2.9). Thus, we can rewrite (2.9) as

$$\mathbf{f}(\mathbf{M}, \boldsymbol{\alpha}) = \sum_{k=1}^R \alpha_k \mathbf{m}_k + \sum_{i=1}^{R-1} \sum_{j=i+1}^R \gamma_{i,j} \alpha_i \alpha_j \mathbf{m}_i \odot \mathbf{m}_j \quad (2.13)$$

$$\alpha_k \geq 0, \quad \forall k \in \{1, \dots, R\} \quad (2.14)$$

$$\sum_{k=1}^R \alpha_k = 1 \quad (2.15)$$

and

$$0 \leq \gamma_{i,j} \leq 1, \quad \forall i \in \{1, \dots, R-1\}, \quad \forall j \in \{i+1, \dots, R\}. \quad (2.16)$$

Comparing (2.13) with (2.11) the following differences can be observed:

- Data generated by both models (noiseless case) are in limited space regions. For model (2.11), there are $R(R-1)/2$ mixed endmembers. Hence, the data are in a simplex in $\mathbb{R}^{[R(R-1)/2]-1}$ whose extremities are the pure endmembers and the mixed endmembers. For model (2.13), the generated data will be in a nonlinearly distorted simplex in \mathbb{R}^R whose extremities correspond to the pure endmembers (like in the linear mixing model). The nature of this nonlinear distortion, however, is difficult to predict. An example has been shown in [40] for small coefficients $\gamma_{i,j}$ in which the distortion looks like a space curvature. However, there is no guarantee that this will always be the case.
- The nonlinear terms $\mathbf{m}_i \odot \mathbf{m}_j$ in (2.13) under constraints (2.15) and (2.16) can no longer be considered as new endmembers in a linear mixing model. This was possible in (2.11) because the coefficients $\beta_{i,j}$ were included in the sum-to-one constraint. Thus, some solutions used for the linear mixing model (such as those based on geometrical approaches) cannot be directly applied to the model (2.13).
- To apply model (2.13) and use a linear unmixing strategy the constraints would have to be modified to

$$\sum_{k=1}^R \alpha_k + \sum_{i=1}^R \sum_{j=1}^R \gamma_{i,j} \alpha_i \alpha_j = 1 \quad (2.17)$$

$$\alpha_k \geq 0, \quad k = 1, \dots, R \quad (2.18)$$

and

$$\gamma_{i,j} \alpha_i \alpha_j \geq 0, \quad i = 1, \dots, R, j = 1, \dots, R. \quad (2.19)$$

Note that in this case, $\gamma_{i,j}$ can be greater than one, which is not allowed in (2.16).

- The physical interpretation of $f(\mathbf{M}, \boldsymbol{\alpha})$ being the sum of the spectral

energies due to the pure endmembers \mathbf{m}_i and the nonlinear endmembers $\mathbf{m}_i \odot \mathbf{m}_j$ to compose the total spectral energy received is lost in (2.13).

- If an endmember \mathbf{m}_k is not present in a target pixel, then its interactions with other endmembers are automatically eliminated from (2.13) ($\alpha_k = 0$). This model also assumes that pure endmembers have been estimated before the unmixing.

2.2.3 Post-Nonlinear Mixing Model - PNMM

This is a large class of nonlinear unmixing models, for which the non-linearity is obtained by applying a nonlinear function to a linear combination of the endmembers. The general expression for the observations is

$$\mathbf{r} = \mathbf{g}(\mathbf{M}\boldsymbol{\alpha}) + \mathbf{n}. \quad (2.20)$$

This model has been initially proposed for source separation problems [41]. It is an interesting idea for mathematical modeling purposes. Although the physical motivation is not very clear, this model can be seen as a generalization of other bilinear models, such as GBM, falling in the particular cases if the function $\mathbf{g}(\cdot)$ is conveniently chosen. For instance, the PNMM considered in [13] is given by

$$\mathbf{r} = (\mathbf{M}\boldsymbol{\alpha})^\xi + \mathbf{n} \quad (2.21)$$

where $(\mathbf{v})^\xi$ denotes the exponentiation applied to each entry of the vector \mathbf{v} . For $\xi = 2$, (2.21) becomes a bilinear model closely related to the GBM but without a linear term. The PNMM was explored in other works considering different forms for \mathbf{g} applied to hyperspectral data unmixing [42, 43].

When the function $\mathbf{g}(\cdot)$ is modeled as a polynomial, the model is often called *post polynomial nonlinear mixing model* (PPNMM) [42]. A very simple form of this model was considered in [44, 17]. The nonlinearity was modeled by a polynomial of degree 2, and given by

$$g(s_i) = s_i + bs_i^2, \quad i = 1, \dots, L \quad (2.22)$$

where s_i is the i -th component of the vector $\mathbf{M}\boldsymbol{\alpha}$.

An interesting characteristic of this model is that it reduces to a bilinear model

$$\mathbf{r} = \mathbf{M}\boldsymbol{\alpha} + b(\mathbf{M}\boldsymbol{\alpha}) \odot (\mathbf{M}\boldsymbol{\alpha}) + \mathbf{n}. \quad (2.23)$$

On the other hand, a limiting feature is the fact that a single parameter b is

used for all the wavelengths corresponding to a given pixel. In [44] Altmann *et al.* states that the PPNMM should present SU results at least as good as those obtained using the linear model considered in [45]. This should be the case since the PPNMM becomes the LMM for $b = 0$. However, the model can add some distortion to the solution once b is the same for all wavelengths.

The second order PPNMM (2.23) was considered in recent works [44, 17]. In [44] a Bayesian approach was applied to solve the SU problem, while an optimization routine was considered in [17]. For both the PPNMM presented good results. A comparative analysis was also performed in [17] where the authors present results considering different mixing models. The PPNMM presented the smallest average error. This illustrates the potential of such models even when restricted simple cases, as (2.23), are considered.

2.2.4 Intimate Mixing Models

While the linear and bilinear models are normally applied to macroscopic level of spectral mixture, sometimes the mixture can occur at microscopic levels, as illustrated in Figure 4b. In this case the endmembers are said to be intimately mixed [36]. Classical examples where intimate mixing occurs are desert sands or mining fields where the endmembers are considered to be minerals composing the “sand” or the “field” and mixed at the microscopic level [46, 47]. Due to its physical meaning, the most popular approaches to model intimately mixed endmembers can be found in [36]. However, more tractable models based on the same concepts found in [36] have been reported in the literature [35, 48, 49]. Nevertheless, intimate models depend on parameters that are inherent to the experiment such as the geometrical positioning of the sensor in relation to the target sample, land relief, etc. This kind of dependence makes the unmixing problem even more complex and difficult.

The use of intimate mixture models makes sense if the endmembers are mixed at the microscopic scale. To illustrate this, consider a target scene containing sand and trees. For a specific application, sand can be considered as one endmember and trees as another. So, we would have just 2 endmembers. However, these endmembers are composed by other materials (minerals, leaves, wood, etc). Thus, the definition of endmember depends on the application, and what one considers to be a “pure” element. To circumvent this issue, the authors in [8] state that it makes sense to associate “pure components” to individual instances with size of the same order of magnitude as the resolution of the sensor used. For this reason, and because of the complexity of the models involved in this kind of mixture, we will not consider models of intimate mixtures in the remaining of this work.

2.3 SPECTRAL UNMIXING: PROBLEM DEFINITION

The spectral unmixing of a hyperspectral image $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_N]$, with N pixels, consists in determining the endmember matrix \mathbf{M} and the abundance matrix $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_N]$. This problem is non-supervised and equivalent to the *blind source separation* (BSS) problem in a typical signal processing view. Some authors refer to this problem as blind hyperspectral unmixing (BHU) [50]. However, given the abundance constraints (2.3), the statistical independence of the sources often assumed in BSS does not hold for BHU [51].

The analysis of HIs has demanded great attention from the scientific community in the past few years [52]. In this period, a variety of methods were proposed to solve the SU problem considering linear and nonlinear mixing [50, 8, 53]. Among the possible approaches we highlight the methods that exploit the convex geometry of the problem, specially if the model is linear. Such methods are usually based on constrained least squares, Bayesian approaches, projection techniques, and convex and non-convex optimization procedures.

When the endmember matrix \mathbf{M} is known, the SU problem is said to be supervised and reduces to inversion (or regression) step. Note, however, that this is a strong assumption since knowledge about the endmembers is rarely available. On the other hand, assuming \mathbf{M} known provides a controlled environment where the true potential of proposed techniques can be unveiled. In Section 2.3.1 the most common endmember extraction techniques are presented, while linear and nonlinear unmixing strategies are discussed in 2.3.2 and 2.3.3 respectively.

2.3.1 Endmembers Estimation

Endmember extraction algorithms (EEA) is the denomination given to a range of endmember estimation techniques that rely on little or no prior information regarding the observed pixels, i.e., $\{\mathbf{r}_n\}_{n=1}^N$. EEAs exploit the convex geometry of the LMM to identify the endmembers. Most EEA algorithms assume the existence of pure pixels within the image. Pure pixels are those that have only one endmember (in contrast to non-pure pixels that are composed of a mixture of endmembers), that is, if a given pixel \mathbf{r}_n is a pure pixel containing only the endmember \mathbf{m}_k then the corresponding abundance

vector $\boldsymbol{\alpha}_n$ has its entries $\alpha_{n,i}$ defined as

$$\alpha_{n,i} = \begin{cases} 1, & i = k \\ 0 & \text{otherwise.} \end{cases}$$

Pure-pixel based approaches assume the existence of at least one pure pixel per endmember. The majority of such algorithms use one of the following two properties:

- a) The extremes of the projection of the spectral vectors \mathbf{r}_n into any subspace correspond to the endmembers.
- b) The hypervolume defined by a set of p spectral vectors is maximum when these vectors are endmembers.

EEAs representing group a) are the *pixel purity index* (PPI) [22], *vertex component analysis* (VCA) [23], *simplex growing algorithm* (SGA) [24], *successive volume maximization* (SVMAX) [26], and *recursive algorithm for separable nonnegative matrix factorization* (RSSNMF) [54]; Algorithms representing the group b) are the N-FINDR [55], *iterative error analysis* (IEA) [56], *sequential maximum angle convex cone* (SMACC), and *alternating volume maximization* (AVMAX) [26].

2.3.2 Linear SU

Assuming the endmember matrix \mathbf{M} to be known¹, the problem boils down to the solution of an inverse problem. For the LMM, the inverse step consists in solving a linear system as in (2.4) for each of the N pixels. This type of linear system are overdetermined, $L \gg R$, and, therefore, has no exact solution [58]. However, an optimal solution can be achieved minimizing the squared error, i.e., the least-squares (LS) [59]. The LS solution for a given pixel \mathbf{r}_n is its orthogonal projection onto the space spanned by the columns of \mathbf{M} , and the projection coefficients are the abundances. Thus, the abundances can be found by solving the following convex quadratic problem

$$\boldsymbol{\alpha}^* = \underset{\boldsymbol{\alpha}}{\operatorname{arg\,min}} \|\mathbf{r} - \mathbf{M}\boldsymbol{\alpha}\|_2^2, \quad (2.24)$$

¹ \mathbf{M} can be previously estimated using an EEA technique, or using other procedures such as local measurements or using digital libraries [57].

for which the linear combination coefficients $\boldsymbol{\alpha}$ has a closed analytical expression given by

$$\boldsymbol{\alpha}^* = \left(\mathbf{M}^\top \mathbf{M}\right)^{-1} \mathbf{M}^\top \mathbf{r}. \quad (2.25)$$

Numerically more stable, the regularized LS [60] can be also applied to find the abundances solving the following convex problem

$$\boldsymbol{\alpha}^* = \arg \min_{\boldsymbol{\alpha}} \|\mathbf{r} - \mathbf{M}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2. \quad (2.26)$$

Solving (2.26) the expression for the $\boldsymbol{\alpha}$'s becomes

$$\boldsymbol{\alpha}^* = \left(\mathbf{M}^\top \mathbf{M} + \lambda \mathbf{I}\right)^{-1} \mathbf{M}^\top \mathbf{r}, \quad (2.27)$$

where λ is the regularization parameter. Although widely used, the presented LS approaches have no constraints over the abundances. In this context, constrained LS approaches can be found in [61, 62], where only the some-to-one constraint is considered, and then, solving the dual unconstrained problem using Lagrange multipliers. A version known as Fully Constrained Least Squares (FCLS) is presented in [63].

When pure pixels are not available in a target scene, different approaches have been proposed in the literature. Among them, we highlight those based on the minimum volume (MV) simplex, where the abundances and endmembers are jointly estimated. From an optimization perspective, MV algorithms can be formulated as

$$\begin{aligned} \mathbf{M}^*, \mathbf{A}^* &= \arg \min_{\mathbf{M}, \mathbf{A}} \|\mathbf{R} - \mathbf{M}\mathbf{A}\|_F^2 + \lambda V(\mathbf{M}) \\ \text{subject to } & \mathbf{A} \succeq 0, \mathbf{1}_R^\top \mathbf{A} = \mathbf{1}_N^\top, \end{aligned} \quad (2.28)$$

where $\|\mathbf{X}\|_F = \sqrt{\text{trace}\{\mathbf{X}\mathbf{X}^\top\}}$ is the Frobenius norm, \succeq denotes the entry-wise \geq operator, i.e., $\alpha_{n,i} \geq 0, \forall n \in \forall i$, $\mathbf{1}_z$ represents a column vector with all z components equal to 1, $V(\mathbf{M})$ is a simplex volume penalty term which promotes a minimum volume estimation, and λ is the regularization parameter. Initially proposed in [64], this approach underlines several geometrical based unmixing algorithms, minimizing successively with respect to \mathbf{M} and \mathbf{A} . This is the case for the *iterative constrained endmembers (ICE) algorithm* [65], and the *minimum volume transform-nonnegative matrix factorization (MVC-NMF)* [66], whose main differences are related with the way they define the regularizer $V(\mathbf{M})$. For variations of these ideas recently introduced, see [67]. It is important to highlight that problem (2.28) is not convex and its solutions are highly dependent on the initialization. A convex formulation, named *min-*

imum volume enclosing simplex (MVES) can be found in [25], by reformulating (2.28) with respect to \mathbf{M}^{-1} instead of \mathbf{M} .

Bayesian methods were also widely employed to solve SU problems. This kind of approach brings great flexibility by incorporating constraints in priors, i.e.; probability density functions (PDFs) *a priori*. It also estimates \mathbf{M} and \mathbf{A} jointly, in a hierarchical Bayesian model using Monte Carlo Markov Chain (MCMC) methods. This type of approach allows posterior PDFs to be estimated even when considering very complex density functions [68]. MCMC Bayesian methods have been applied to linear SU problems in [45, 69, 70, 71].

2.3.3 Nonlinear SU

In recent years, promising methods were employed to nonlinear SU of HIs. A large portion of these methods are based on specific nonlinear parametric models as presented in Sections 2.2.2 and 2.2.4. However, some approaches assume only general characteristics about the type of nonlinearity or nonlinear function. This second group has been called in the literature *model-free* (or *model independent*) nonlinear spectral unmixing [8].

2.3.3.1 Nonlinear SU using Parametric Models

For a given parametric model, the problem of SU can be formulated as a constrained nonlinear regression or as a nonlinear source separation problem, depending on whether the endmembers are known or not.

Assuming that the mixing matrix \mathbf{M} is known, various approaches have been proposed for supervised bilinear models. In this context, SU of a given pixel \mathbf{r} can be formulated in a general way as the following minimization problem

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} \|\mathbf{r} - \boldsymbol{\varphi}(\mathbf{M}, \boldsymbol{\theta})\|_2^2 \\ &\text{subject to } \boldsymbol{\theta} \in \Omega, \end{aligned} \quad (2.29)$$

where $\boldsymbol{\theta}$ is a vector containing the abundances and any other model parameters, $\boldsymbol{\varphi}(\cdot)$ is a parametric bilinear function, and Ω defines the feasible region for the vector $\boldsymbol{\theta}$, i.e., $\boldsymbol{\alpha} \in \mathcal{S}_{\boldsymbol{\alpha}}$, and further constraints (possibly) imposed over other model parameters. Since the Nascimento model can be interpreted as a linear model with a new endmember basis (2.12), linear approaches can be used to solve the problem, as in [39]. In [38, 40] the authors considered the

FCLS proposed in [72] to estimate θ for a linearized version (via Taylor series expansion) of $\varphi(\cdot)$. A fully Bayesian approach based on Monte Carlo approximations was conducted in [18] for the GBM. This same strategy was applied in [17] but considering the PPNMM.

Methods for unsupervised nonlinear SU are also reported in the literature. In [73] the authors extended the Bayesian methodology considered in [17] to solve an unsupervised problem. Another work worth mentioning is the method proposed by Heylen and Scheunders [74]. In this work the problem was approached from a geometric point of view considering manifold learning techniques, using an integral formulation to calculate the “true” geodesic distances in the manifold induced by the GBM. Within this concept, strategies for nonlinear unmixing, dimensionality reduction, clustering, or classification can be reformulated.

Artificial neural networks have also been considered when adopting intimate mixture models [75]. More details about these and other methods applied to supervised nonlinear SU can be found in [8] and in their references.

2.3.3.2 Model-free Nonlinear SU

When the type of nonlinearity is unknown, more flexible approaches should be sought. In this context, methods based on *reproducing kernels* [76, 77] are of particular interest due to their capacity to approximate functions without the need for a strict parametric model. Several kernel-based methods were presented in the nonlinear SU literature [8]. However, not all of them take into account that the type of nonlinearity involved in the mixing is mainly due to multiple interactions between light and different endmembers. Chen *et al.* [13] presents a supervised kernel-based formulation considering multiple interactions between light and endmembers. For this, a partially linear model was proposed leading to algorithms called K-Hype and SK-Hype (the latter is used in this work and discussed in Section 3.1.4).

Within an unsupervised setting, two approaches are particularly relevant for model-free unmixing. In [28] a Bayesian approach was proposed employing a Gaussian process latent variable model (GPLVM) as a tool for probabilistic nonlinear dimensionality reduction. Although the authors make no prior strong assumption about the spectral signatures in \mathbf{M} , they consider the number of endmembers (R) to be known. The second approach is related to manifold learning [78, 79], where geodesic distances are approximated based on graphs obtained directly from the data [27].

2.4 DETECTION OF NONLINEARLY MIXED PIXELS

It is now acknowledged that nonlinear unmixing algorithms can lead to a better understanding of the individual spectral contributions. On the other hand, nonlinear analysis techniques are necessarily more challenging and complex than those employed for linear unmixing. As hyperspectral images tend to include both linearly and nonlinearly mixed pixels, there are two important reasons to match the unmixing method to the nature of each pixel in the image. First, nonlinear unmixing algorithms are always more complex to implement than linear unmixing algorithms. Second, unmixing linearly mixed pixels with nonlinear unmixing algorithms leads to poorer results than doing it with linear unmixing algorithms. Hence, it makes sense to detect the nonlinearly mixed pixels in an image prior to its analysis, and then employ the simplest and more accurate available unmixing technique to analyze each pixel. However, detecting nonlinearly mixed pixels in a hyperspectral image is also a complex task. Physically motivated models [34, 80] usually tend to be too complex for application in practical detection strategies.

In a pioneer work in the analysis of hyperspectral data Han and Goodenough [81] used surrogate data, borrowed from analysis of nonlinear dynamical time series [82], to test nonlinear hypothesis in HIs. In [81] a pixel was seen as a realization of a dynamical nonlinear system along the wavelengths, and considered the same approach as in [83]. However, nonlinearity in HIs are modeled in the amplitude relation within each band (see Section 2.2), not in its dynamics, i.e., along the bands. One possible approach is to consider a simplified parametric model for the nonlinearity. The parameters of this nonlinear model are then estimated from the image, and hypothesis tests are derived based on these estimates. For instance, a single-parameter polynomial post-nonlinear model is assumed in [42]. The main question regarding parametric modeling of nonlinear mixing mechanisms is whether the chosen model can capture the actual nonlinear effects present in a scene. When nothing or little is known about the nonlinear mixing mechanism, a direct strategy is to exploit the property of linear mixing models to confine the noiseless data to a simplex. The hypothesis test proposed in [2] is based on the distance between the observed pixel and this simplex. Though this test is robust to nonlinear mixing mechanisms, it conveys too little information about the nonlinearity as a trade-off to guarantee simplicity. An alternative strategy is to use nonparametric techniques to extract information about the nonlinearity directly from the observations. A nonparametric unmixing technique based on kernel expansions is presented in [13], but this work does not address nonlinearity detection. A nonlinear mixing model for joint unmixing and nonlinearity detection is proposed in [84]. It assumes that the observed reflectances

result from linear spectral mixtures corrupted by a residual nonlinear component. This model is rather similar to the model initially introduced in [13], but the estimation method relies on a computationally intensive Bayesian procedure.

2.5 BAND SELECTION

Nonlinear methods have been successfully applied to unmix nonlinearly mixed HIs [8], where the size of the input data equals the number of spectral bands in a space with dimension equal to the number of endmembers [8]. It means that when dealing with HIs, and the profusion of spectral bands generated in the acquisition process, these methods must deal with matrices composed of hundreds or even thousands of vectors for each pixel. Such inherent complexity, associated with the high redundancy within the complete set of bands, make the search for *band selection* (BS) techniques natural and relevant [29]. Several BS algorithms have been proposed for linearly mixed HIs, which generally require solving an optimization problem [85]. However, BS for nonlinear unmixing presents an even more challenging problem.

Band selection has been an active topic for classification of spectral patterns, see [86, 87, 88, 89, 90] and references there in. When concerning unmixing of HIs, band selection approaches [91, 85] are usually deprecated in relation to subspace projection techniques [92, 23, 93]. This happens because the more consolidated literature assumes linear mixing models which confine the data into a low dimensional simplex [6]. Such assumption is lost when nonlinear mixing models are considered. It is the case when considering kernelized methods such as SK-Hype [13]. We highlight, however, that mutual information based BS strategies [85] might be considered under nonlinear modeling of the mixing occurring in hyperspectral images. This, however, must be done with the proper care and is out of the scope of this work.

3 PRELIMINARY THEORETICAL CONCEPTS

3.1 KERNEL REGRESSION

This section describes the two kernel frameworks for supervised non-linear regression considered in this Thesis. The representation is rigorous but, at the same time, lets the data speak for themselves. This characteristic is desirable when little is known about the functions to be estimated. Using some knowledge obtained from the observations about the endmember matrix, we propose a supervised learning strategy to make inference on $\boldsymbol{\psi}$. Consider the training set $\{\mathbf{M}, \mathbf{r}\}$ with inputs $\mathbf{M} = [\mathbf{m}_{\lambda_1}, \dots, \mathbf{m}_{\lambda_L}]^\top$, and outputs or observations $\mathbf{r} = [r_1, \dots, r_L]^\top$, where $\mathbf{m}_{\lambda_\ell}$ is a column vector of the R endmember signatures at the ℓ -th wavelength, that is, $\mathbf{m}_{\lambda_\ell}^\top$ is the ℓ -th line of the \mathbf{M} matrix. By analogy with the linear mixing model (2.5), we write the ℓ -th row of (2.1) as

$$r_\ell = \boldsymbol{\psi}(\mathbf{m}_{\lambda_\ell}) + n_\ell, \quad (3.1)$$

with r_ℓ the ℓ -th entry of the observation \mathbf{r} , $\boldsymbol{\psi}$ a real-valued function in a (reproducing kernel) Hilbert space \mathcal{H} , and n_ℓ an additive WGN in the ℓ -th band.

Next, a discussion about Mercer kernels, and regression using Gaussian process and least-squares support vector regression (LS-SVR) is presented. The theory related to these methods involves concepts from functional analysis, convex optimization and functional derivatives. These concepts are presented in Appendix A – Convex optimization in RKHS. In Appendix B a more didactic discussion about Gaussian process is presented where the kernel framework is shown to be a generalization of the standard Bayesian linear regression under a few assumptions.

3.1.1 Mercer kernels and RKHS

The theory of positive definite kernels emerged from the study of positive definite integral operators [94], and was further generalized in the study of positive definite matrices [95]. It was established that to every positive definite function $\kappa(\cdot, \cdot) : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$, defined over a non-empty compact $\mathcal{M} \subset \mathbb{R}^d$, there corresponds one and only one class of real-valued functions on \mathcal{M} forming a Hilbert space \mathcal{H} with a uniquely defined inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ in it, and admitting κ as a *reproducing kernel* (r.k.) [96]. By definition, see [97] or Appendix A.3.2, κ is a r.k. of \mathcal{H} if $\kappa(\cdot, \mathbf{m}_\lambda) \in \mathcal{H}$ for all

$\mathbf{m}_\lambda \in \mathcal{M}$, and the *reproducing property*

$$\psi(\mathbf{m}_\lambda) = \langle \psi, \kappa(\cdot, \mathbf{m}_\lambda) \rangle_{\mathcal{H}} \quad (3.2)$$

holds, for all $\psi \in \mathcal{H}$ and all $\mathbf{m}_\lambda \in \mathcal{M}$. For the existence of a r.k. κ it is necessary and sufficient that for every $\mathbf{m}_\lambda \in \mathcal{M}$ the evaluation (Dirac) functional $\delta_{\mathbf{m}_\lambda}$ to be a continuous (or equivalently bounded) functional for any $\psi \in \mathcal{H}$. On the other hand, supposing $\delta_{\mathbf{m}_\lambda}$ to be also linear there exists a function $\varphi_{\delta_{\mathbf{m}_\lambda}} \in \mathcal{H}$ such that $\delta_{\mathbf{m}_\lambda}[\psi] = \psi(\mathbf{m}_\lambda) = \langle \psi, \varphi_{\delta_{\mathbf{m}_\lambda}} \rangle_{\mathcal{H}}$ (Riesz representation theorem [98, pg.188]). Thus, if \mathcal{H} is a Hilbert space with continuous linear evaluation functional, then \mathcal{H} is called a *reproduced kernel Hilbert space* (RKHS) admitting κ as its unique reproducing kernel, and $\kappa(\cdot, \mathbf{m}_\lambda) = \varphi_{\delta_{\mathbf{m}_\lambda}}$ is called representer of the evaluation at \mathbf{m}_λ . Furthermore, as a direct consequence of the Riesz theorem, we have that $\kappa(\cdot, \mathbf{m}_\lambda)$ depends on $\delta_{\mathbf{m}_\lambda}$, is uniquely defined by $\delta_{\mathbf{m}_\lambda}$, and has norm $\|\kappa(\cdot, \mathbf{m}_\lambda)\|_{\mathcal{H}} = \|\delta_{\mathbf{m}_\lambda}\|_{\mathcal{H}'}^1$. The RKHS \mathcal{H} is then formed by a class of functions generated by all functions of the form $\psi(\cdot) = \sum_j \alpha_j \kappa(\cdot, \mathbf{m}_{\lambda_j})$, with norm defined by the quadratic form $\|\psi\|_{\mathcal{H}}^2 = \sum_i \sum_j \alpha_i \alpha_j \kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})$.

In the context of machine learning, kernel methods are often related with the concept of building a high dimensional feature space \mathcal{H} , and a mapping

$$\Phi: \mathcal{M} \longrightarrow \mathcal{H} \quad (3.3)$$

$$\mathbf{m}_\lambda \longmapsto \Phi(\mathbf{m}_\lambda) \quad (3.4)$$

with inner product defined as $\kappa(\mathbf{m}_\lambda, \mathbf{m}_{\lambda'}) = \langle \Phi(\mathbf{m}_\lambda), \Phi(\mathbf{m}_{\lambda'}) \rangle_{\mathcal{H}}$. If κ is a r.k. of \mathcal{H} , then \mathcal{H} is a RKHS and also a feature space of κ with $\Phi(\mathbf{m}_\lambda) = \kappa(\cdot, \mathbf{m}_\lambda)$. In this case Φ is called the *canonical feature map* [99, pg. 120]. This leads to the so called “*kernel trick*” allowing one to compute inner products of data mapped into higher, or even infinite, dimensional feature spaces by evaluating a real function $\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})$ in the input space.

Several kernel functions have been considered in a variety of applications and algorithms during the past two decades of intense research activity [76, 77, 100]. Among the most frequently used kernels we highlight the Gaussian kernel. When defined over a compact metric space \mathcal{M} , Gaussian kernels, among other continuous kernels, are known to produce RKHSs \mathcal{H} that are dense in the space of continuous functions $f: \mathcal{M} \rightarrow \mathbb{R}$, namely $C(\mathcal{M})$. This means that for every function $f \in C(\mathcal{M})$ and all $\varepsilon > 0$ there exists an $\psi \in \mathcal{H}$ such that $\|f - \psi\|_\infty \leq \varepsilon$. Kernels having such property

¹ \mathcal{H}' denotes the set of all bounded linear functionals $\zeta: \mathcal{H} \rightarrow \mathbb{R}$. \mathcal{H}' is also a Hilbert space, and is called the *dual* of \mathcal{H} [98]. See Appendix A.3.1 Definition 7.

are often referred in the literature as *universal kernels* [99]. It is important to punctuate, however, that universal kernels can lead to overfitting learning curves into the data. This can be a problem specially when the learning process is embedded in high noise levels. Finally, for entries $\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j} \in \mathcal{M}$, the Gaussian kernel is given by

$$\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j}) = \exp\left(-\frac{\|\mathbf{m}_{\lambda_i} - \mathbf{m}_{\lambda_j}\|^2}{2\sigma^2}\right) \quad (3.5)$$

where the parameter $\sigma^2 > 0$ controls the the kernel bandwidth. Other examples of common used kernel functions are the linear kernel

$$\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j}) = \mathbf{m}_{\lambda_i}^\top \mathbf{m}_{\lambda_j} \quad (3.6)$$

and the polynomial kernel

$$\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j}) = (\mathbf{m}_{\lambda_i}^\top \mathbf{m}_{\lambda_j} + c)^d \quad (3.7)$$

where d is the polynomial degree and $c \geq 0$ is a real number.

3.1.2 Gaussian Process for regression

Gaussian process (GP) regression methods consist of defining stochastic models for functions and performing inference in functional spaces [100]. A more detailed presentation can be found in Appendix B. A Gaussian process is a collection of random variables, any finite number of which has a joint Gaussian distribution [100]. Considering the model presented in (3.1), replicated here for convenience,

$$r_\ell = \psi(\mathbf{m}_{\lambda_\ell}) + n_\ell,$$

we define a Gaussian prior distribution for ψ with mean and covariance functions given by

$$\begin{aligned} \mathbb{E}\{\psi(\mathbf{m}_{\lambda_\ell})\} &= 0 \\ \mathbb{E}\{\psi(\mathbf{m}_{\lambda_\ell})\psi(\mathbf{m}_{\lambda_{\ell'}})\} &= \kappa(\mathbf{m}_{\lambda_\ell}, \mathbf{m}_{\lambda_{\ell'}}) \end{aligned} \quad (3.8)$$

where κ is a positive definite kernel. For notational simplicity, it is common but not necessary to consider GPs with zero mean. This assumption is not overly restricting as the mean of the posterior distribution is not confined to

be zero (as shown by (3.11)). The prior on the noisy observation \mathbf{r} becomes:

$$\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma_n^2 \mathbf{I}), \quad (3.9)$$

with \mathbf{K} the Gram matrix whose entries $\mathbf{K}_{ij} = \kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})$ are given by the kernel covariance function evaluated at \mathbf{m}_{λ_i} and \mathbf{m}_{λ_j} , σ_n^2 the noise power, and \mathbf{I} the $L \times L$ identity matrix.

To obtain the predictive distribution for $\psi_* \triangleq \psi(\mathbf{m}_{\lambda_*})$ at any test point \mathbf{m}_{λ_*} , we can write the joint distribution of the observation \mathbf{r} and $\psi(\mathbf{m}_{\lambda_*})$ as [100]

$$\begin{bmatrix} \mathbf{r} \\ \psi_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \boldsymbol{\kappa}_* \\ \boldsymbol{\kappa}_*^\top & \kappa_{**} \end{bmatrix} \right) \quad (3.10)$$

with $\boldsymbol{\kappa}_* = [\kappa(\mathbf{m}_{\lambda_*}, \mathbf{m}_{\lambda_1}), \dots, \kappa(\mathbf{m}_{\lambda_*}, \mathbf{m}_{\lambda_L})]^\top$ and $\kappa_{**} = \kappa(\mathbf{m}_{\lambda_*}, \mathbf{m}_{\lambda_*})$. The predictive distribution of ψ_* , or posterior of ψ_* , is then obtained by conditioning (3.10) on the observation as follows:

$$\begin{aligned} \psi_* | \mathbf{r}, \mathbf{M}, \mathbf{m}_{\lambda_*} \sim \mathcal{N} \left(\boldsymbol{\kappa}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{r}, \right. \\ \left. \kappa_{**} - \boldsymbol{\kappa}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \boldsymbol{\kappa}_* \right). \end{aligned} \quad (3.11)$$

The extension to a multivariate predictive distribution with test data $\mathbf{M}_* = [\mathbf{m}_{\lambda_{*1}}, \dots, \mathbf{m}_{\lambda_{*L}}]^\top$ yields:

$$\begin{aligned} \boldsymbol{\psi}_* | \mathbf{r}, \mathbf{M}, \mathbf{M}_* \sim \mathcal{N} \left(\mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{r}, \right. \\ \left. \mathbf{K}_{**} - \mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{K}_* \right) \end{aligned} \quad (3.12)$$

with $[\mathbf{K}_*]_{ij} = \kappa(\mathbf{m}_{\lambda_{*i}}, \mathbf{m}_{\lambda_{*j}})$ and $[\mathbf{K}_{**}]_{ij} = \kappa(\mathbf{m}_{\lambda_{*i}}, \mathbf{m}_{\lambda_{*j}})$. Finally, we arrive at the minimum mean square error (MMSE) estimator for GP regression:

$$\begin{aligned} \hat{\boldsymbol{\psi}}_* &= \mathbb{E}\{\boldsymbol{\psi}_* | \mathbf{r}, \mathbf{M}, \mathbf{M}_*\} \\ &= \mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{r}. \end{aligned} \quad (3.13)$$

In order to turn GP into a practical tool for processing hyperspectral data, it is essential to derive a method for estimating free parameters such as the noise variance σ_n^2 and possible kernel parameters defining the unknown parameter vector $\boldsymbol{\theta}$. We proceed as in [100] by maximizing the marginal

likelihood $p(\mathbf{r}|\mathbf{M}, \sigma_n^2, \boldsymbol{\theta})$ with respect to $(\sigma_n^2, \boldsymbol{\theta})$, which yields

$$(\hat{\sigma}_n^2, \hat{\boldsymbol{\theta}}) = \arg \max_{\sigma_n^2, \boldsymbol{\theta}} \left(-\frac{1}{2} \mathbf{r}^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{r} - \frac{1}{2} \log |\mathbf{K} + \sigma_n^2 \mathbf{I}| \right). \quad (3.14)$$

This problem has to be addressed with numerical optimization methods. There is no guarantee that the cost function does not suffer from multiple local optima. However, our practical experience with hyperspectral data indicates that local optima are not a critical problem in this context. The solutions to the optimization problem (3.14) for all examples reported in this work were determined using the GPML (*Gaussian Processes for Machine Learning*) toolbox [101].

In the sequel, we shall use the Gaussian kernel for its smoothness and non-informativeness, as we lack any knowledge about the unknown function ψ . Then, $\boldsymbol{\theta} = \sigma$ (scalar). Note that this kernel has been used successfully in many signal and image processing applications, in particular for hyperspectral data unmixing [13, 16].

3.1.3 LS-SVR

This section describes the use of a state-of-the-art kernel method for nonlinear unmixing of hyperspectral data. Consider an observation r_ℓ at the ℓ -th wavelength, modeled as in (3.1), with ψ a real-valued function in a RKHS \mathcal{H} that characterizes the nonlinear interactions between the endmembers, and n_ℓ an additive noise at the ℓ -th band. In order to estimate ψ in the least squares sense, we can formulate the following convex optimization problem, also called LS-SVR [102]²:

$$\begin{aligned} \min_{\psi \in \mathcal{H}} \quad & \frac{1}{2} \|\psi\|_{\mathcal{H}}^2 + \frac{1}{2\mu} \sum_{\ell=1}^L e_\ell^2 \\ \text{such that} \quad & e_\ell = r_\ell - \psi(\mathbf{m}_{\lambda_\ell}), \quad \ell = 1, \dots, L. \end{aligned} \quad (3.15)$$

Consider the Lagrangian function

$$\mathcal{L}(\boldsymbol{\psi}, \mathbf{e}, \boldsymbol{\beta}) = \frac{1}{2} \|\boldsymbol{\psi}\|_{\mathcal{H}}^2 + \frac{1}{2\mu} \sum_{\ell=1}^L e_\ell^2 - \sum_{\ell=1}^L \beta_\ell (e_\ell - r_\ell + \psi(\mathbf{m}_{\lambda_\ell})). \quad (3.16)$$

²This is a direct application of the LS-SVR presented in Appendix A.6.

where $\boldsymbol{\beta} = [\beta_1, \dots, \beta_L]^\top$ is the vector of Lagrange multipliers. Using the directional derivative with respect to $\boldsymbol{\psi}$ [103], the conditions for optimality with respect to the primal variables $\boldsymbol{\psi}$ and e_ℓ are given by

$$\boldsymbol{\psi}^* = \sum_{\ell=1}^L \beta_\ell \boldsymbol{\kappa}(\cdot, \mathbf{m}_{\lambda_\ell}) \quad (3.17)$$

$$e_\ell^* = \mu \beta_\ell \quad (3.18)$$

Substituting (3.17) and (3.18) in (3.16), we obtain the following function to be maximized with respect to $\boldsymbol{\beta}$:

$$\mathcal{L}(\boldsymbol{\psi}^*, \mathbf{e}^*, \boldsymbol{\beta}) = -\frac{1}{2} \boldsymbol{\beta}^\top (\mathbf{K} + \mu \mathbf{I}) \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{r}, \quad (3.19)$$

where \mathbf{K} is the Gram matrix whose (i, j) -th entry is defined by $\boldsymbol{\kappa}(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})$. Now we can state the following dual problem:

$$\boldsymbol{\beta}^* = \arg \max_{\boldsymbol{\beta}} -\frac{1}{2} \boldsymbol{\beta}^\top (\mathbf{K} + \mu \mathbf{I}) \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{r}. \quad (3.20)$$

Its solution is obtained by solving the linear system:

$$(\mathbf{K} + \mu \mathbf{I}) \boldsymbol{\beta} = \mathbf{r}. \quad (3.21)$$

Although the formulation (3.15)–(3.20) allows one to address an estimation problem in \mathcal{H} by solving the linear system (3.21), this approach is computationally demanding since it involves the inversion of $L \times L$ matrices. This issue is critical, as modern hyperspectral image sensors employ hundreds of contiguous bands with an ever increasing spatial resolution.

3.1.4 SK-Hype

This section reviews the SK-Hype algorithm³ for nonlinear unmixing of HIs [13]. It considers the mixing model consisting of a linear trend parametrized by the abundance vector $\boldsymbol{\alpha}$ and a nonlinear residual component $\boldsymbol{\psi}_{\text{nl}}$. Thus, it considers an extension of the model $\mathbf{r} = \boldsymbol{\psi}(\mathbf{M}) + \mathbf{n}$ (see Equation 3.1), where the underlying function $\boldsymbol{\psi}$ is

$$\boldsymbol{\psi}(\mathbf{m}_{\lambda_\ell}) = \theta \boldsymbol{\alpha}^\top \mathbf{m}_{\lambda_\ell} + \boldsymbol{\psi}_{\text{nl}}(\mathbf{m}_{\lambda_\ell}), \quad (3.22)$$

³Matlab code available at www.cedric-richard.fr

where $\mathbf{m}_{\lambda_\ell}^\top$ is the ℓ -th line of \mathbf{M} , θ is a new parameter affecting the linear mixing, which in turn is parameterized by the abundance vector $\boldsymbol{\alpha}$. In [13] and [104], ψ_{lin} is defined as a vector from a *reproducing kernel Hilbert space* \mathcal{H}_{lin} .

In this section we briefly review the LS-SVR problem solved in SK-Hype. This problem couples with the model presented in (3.22) and it is a simple extension of the LS-SVR presented in Section 3.1.3 to a multi-kernel setting. Here the endmember matrix \mathbf{M} is considered known, and the SK-Hype solves a supervised nonlinear kernelized regression problem. Thus, the LS-SVR problem presented in [13] is given by

$$\min_u J(u) \quad \text{subject to} \quad 0 \leq u \leq 1 \quad (3.23)$$

with

$$J(u) = \begin{cases} \min_{\psi} F(u, \psi) = \frac{1}{2} \left(\frac{1}{u} \|\psi_{\text{lin}}\|_{\mathcal{H}_{\text{lin}}}^2 + \frac{1}{1-u} \|\psi_{\text{nonlin}}\|_{\mathcal{H}_{\text{nonlin}}}^2 \right) + \frac{1}{2\mu} \sum_{\ell=1}^L e_\ell^2 \\ \text{subject to} \quad e_\ell = r_\ell - \psi(\mathbf{m}_{\lambda_\ell}) \quad \text{with} \quad \psi = \psi_{\text{lin}} + \psi_{\text{nonlin}} \\ \text{and} \quad \psi_{\text{lin}}(\mathbf{m}_{\lambda_\ell}) = \mathbf{h}^\top \mathbf{m}_{\lambda_\ell} \quad \text{with} \quad \mathbf{h} \succeq \mathbf{0} \end{cases} \quad (3.24)$$

where \mathcal{H}_{lin} and $\mathcal{H}_{\text{nonlin}}$ are RKHSs, $\mathbf{h} = \theta \boldsymbol{\alpha}$, and $u \in [0, 1]$ controls the linear and nonlinear contributions. Note that the vector \mathbf{h} does not have the sum-to-one constraint, this was done intentionally ensure the convexity of the problem. However, since $\mathbf{h} = \theta \boldsymbol{\alpha}$ and $\mathbf{1}^\top \boldsymbol{\alpha} = 1$, the optimal fully constrained abundances can be computed as $\boldsymbol{\alpha}^* = \mathbf{h}^* / \theta^*$ with $\theta^* = \mathbf{1}^\top \mathbf{h}^*$.

As stated, it can be shown that the problem (3.23)-(3.24) is convex and more details can be found in [13]. Chen *et al* [104] solve the problem (3.23)-(3.24) using a iterative optimization procedure that alternates the solutions of (3.23)-(3.24) with respect to ψ and with respect to u successively.

3.1.4.1 Solving with respect to ψ

By the strong duality property [105], the solutions to the primal problem $J(u) = F(u, \psi^*)$ and its dual are the same (see, Appendix A.2). The Lagrangean function for problem (3.24) can be written using the Lagrange

multipliers β_ℓ and γ_r as

$$G = \frac{1}{2} \left(\frac{1}{u} \|\mathbf{h}\|^2 + \frac{1}{1-u} \|\Psi_{\text{nlín}}\|_{\mathcal{H}_{\text{nlín}}}^2 \right) + \frac{1}{2\mu} \sum_{\ell=1}^L e_\ell^2 - \sum_{\ell=1}^L \beta_\ell (e_\ell - r_\ell + \psi(\mathbf{m}_{\lambda_\ell})) - \sum_{r=1}^R \gamma_r h_r \quad (3.25)$$

with $\gamma_r \geq 0$, and using $\|\Psi_{\text{nlín}}\|_{\mathcal{H}_{\text{nlín}}}^2 = \|\mathbf{h}\|^2$.⁴ The optimality conditions (see Appendix A.2) of G with respect to the primal variables are given by

$$\begin{cases} \mathbf{h}^* = u (\sum_{\ell=1}^L \beta_\ell^* \mathbf{m}_{\lambda_\ell} + \boldsymbol{\gamma}^*) \\ \Psi_{\text{nlín}}^* = (1-u) \sum_{\ell=1}^L \beta_\ell^* \kappa_{\text{nlín}}(\cdot, \mathbf{m}_{\lambda_\ell}) \\ e_\ell^* = \mu \beta_\ell^* \end{cases} \quad (3.26)$$

Replacing (3.26) in (3.25), we obtain the following dual problem

$$J(u) = \begin{cases} \max_{\boldsymbol{\beta}, \boldsymbol{\gamma}} G'(u, \boldsymbol{\beta}, \boldsymbol{\gamma}) = \\ -\frac{1}{2} \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}^\top \left(\begin{array}{c|c} \mathbf{K}_u + \mu \mathbf{I} & u \mathbf{M} \\ \hline u \mathbf{M}^\top & u \mathbf{I} \end{array} \right) \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} + \begin{pmatrix} \mathbf{r} \\ \mathbf{0} \end{pmatrix}^\top \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} \\ \text{subject to } \boldsymbol{\gamma} \succeq \mathbf{0} \end{cases} \quad (3.27)$$

where $\mathbf{K}_u = u \mathbf{M} \mathbf{M}^\top + (1-u) \mathbf{K}_{\text{nlín}}$. Solving (3.27) with respect to the Lagrange multipliers is equivalent to solving the linear system

$$\begin{pmatrix} \mathbf{K}_u + \mu \mathbf{I} & u \mathbf{M} \\ \hline u \mathbf{M}^\top & u \mathbf{I} \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} = \begin{pmatrix} \mathbf{r} \\ \mathbf{0} \end{pmatrix} \quad (3.28)$$

which requires a $(L+R) \times (L+R)$, $L \gg R$, matrix inverse.

The estimative of a pixel can be obtained using $\psi^*(\mathbf{m}_{\lambda_\ell}) = \mathbf{m}_{\lambda_\ell}^\top \mathbf{h}^* + \Psi_{\text{nlín}}^*(\mathbf{m}_{\lambda_\ell})$ for all wavelengths, i.e., $\mathbf{r}^* = [\psi^*(\mathbf{m}_{\lambda_1}), \dots, \psi^*(\mathbf{m}_{\lambda_L})]^\top$, where $\Psi_{\text{nlín}}^*$ is given by (3.26). Finally, the abundance vector is estimated as

$$\boldsymbol{\alpha}^* = \frac{\mathbf{M}^\top \boldsymbol{\beta}^* + \boldsymbol{\gamma}^*}{\mathbf{1}^\top (\mathbf{M}^\top \boldsymbol{\beta}^* + \boldsymbol{\gamma}^*)}. \quad (3.29)$$

⁴Note that $\Psi_{\text{nlín}}(\cdot) = (\cdot)^\top \mathbf{h}$. Thus, $\langle \Psi_{\text{nlín}}, \Psi_{\text{nlín}} \rangle_{\mathcal{H}_{\text{nlín}}} = \langle \mathbf{h}, \mathbf{h} \rangle_{\mathcal{H}_{\text{nlín}}} = \|\mathbf{h}\|^2$.

3.1.4.2 Solving with respect to u

In [104] a closed analytical form is presented for u^* . Note that

$$f_{p,q}(u) = \frac{p}{u} + \frac{q}{1-u}, \quad \text{com } p, q \geq 0 \quad (3.30)$$

is convex in the interval $]0, 1[$. Thus simple derivatives lead to the optimum solution given by

$$u^* = (1 + \sqrt{q/p})^{-1}. \quad (3.31)$$

Considering the problem (3.23), and using the stationarity conditions in (3.26) the optimum solution becomes

$$u^* = \left(1 + (1 - u_{-1}^*) \sqrt{(\boldsymbol{\beta}^{*\top} \mathbf{K}_{\text{nl in}} \boldsymbol{\beta}^*) / \|\mathbf{h}^*\|^2} \right) \quad (3.32)$$

where u_{-1}^* is the optimum u^* for the previous iteration.

3.2 BASIC DETECTION CONCEPTS

Modern detection theory is fundamental to the design of electronic signal processing system for decision making and information extraction. Such systems share the common goal of being able to decide when an event of interest occurs and then determine more information about that event. Basic signal detection theory often assumes cases where the event of interest is in fact a signal that may change its behavior in noisy measures. Such signal is often referred as *target* signal. A vast theory exists on the subject concerning different scenarios and considerations about the behavior of the target signal. The simplest case is when the target signal may or may not be present in noisy measures. When one is concerned about detecting changes in the underlying model of a signal that is always present, then such operation is often referred to as model change detection. Applications are found in numerous fields including radar, communication, sonar, image processing, etc. We emphasize, however, that for some physically motivated applications usual detection strategies may result in models that are too complex, or inefficient, and alternatives must be sought.

In this section a brief discussion of basic detection concepts is presented. More detailed information about detectors and detection systems can be found in [106].

3.2.1 The detection problem

A detection system aims at deciding when an event of interest occurs or not. Thus, two hypotheses are directly present and stated as

$$\mathcal{H}_0 : \text{The event did not occur,} \quad (3.33a)$$

$$\mathcal{H}_1 : \text{The event occurred.} \quad (3.33b)$$

Thus the detection problem is divided in two mutually exclusive hypotheses. The first hypothesis \mathcal{H}_0 is that the event of interest has not occurred, while the second considers that the event has occurred. This problem is known as binary hypothesis testing, since the outcome must be selected between two hypotheses. Thus, for a given signal sample x , the detector must assign x to one of the two hypotheses. Such problem is widely addressed in the literature and sometimes referred as binary classification [106, 107, 108]. For any binary detection problem two types of errors can be defined when assigning a signal sample x to one of the two possible hypotheses. The *type I* error occurs if x is assigned to \mathcal{H}_1 but \mathcal{H}_0 is true. This type of error is also referred as *false alarm*. The *type II* error is when x is assigned to \mathcal{H}_0 but \mathcal{H}_1 is true. The probability of false alarm (PFA) is the probability that the detector decides for \mathcal{H}_1 but \mathcal{H}_0 is true, that is $P(\mathcal{H}_1|\mathcal{H}_0)$. The probability of detection (PD) is the probability that the detector decides for \mathcal{H}_1 and \mathcal{H}_1 is true, that is $P(\mathcal{H}_1|\mathcal{H}_1)$. When designing a detector a common objective is to have a small PFA while maximizing the PD.

To perform the detection a test statistic T is needed. This test is then compared to a decision threshold τ to determine if a signal sample x follows \mathcal{H}_0 or \mathcal{H}_1 .

$$T \underset{\mathcal{H}_1}{\overset{\mathcal{H}_0}{\leq}} \tau \quad (3.34)$$

where, the notation $\underset{\mathcal{H}_1}{\overset{\mathcal{H}_0}{\leq}}$ presents the detection decision depending on the comparison result, that is if T is smaller than τ , then x is assigned to \mathcal{H}_0 , and to \mathcal{H}_1 otherwise. The problems of finding a meaningful test statistic and an optimal detection threshold τ are interlaced. In the literature many strategies have been used to solve such problems. One of the most widely used is the Neyman-Pearson (NP) strategy which address the two problems simultaneously by maximizing the probability of detection while fixing the probability of false alarm, that is $PFA = \xi$. Such strategy leads to the known likelihood

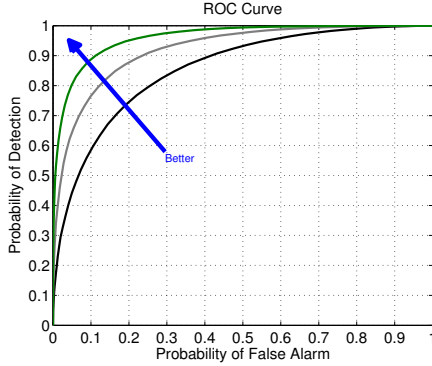


Figure 6: ROC Curve. Different detectors presented in different colors.

ratio test (LRT) given by

$$L(x) = \frac{P(x|\mathcal{H}_1)}{P(x|\mathcal{H}_0)} \underset{\mathcal{H}_1}{\overset{\mathcal{H}_0}{\gtrless}} \tau \quad (3.35)$$

where τ is found from

$$\text{PFA} = \int_{L(x) > \tau} P(x|\mathcal{H}_0) dx = \xi. \quad (3.36)$$

In many cases, however, considering a test statistic based on LRT or its variants can be intractable depending on the underlying signal models and their resulting PDFs. In such cases alternative test statistics must be sought. Such alternatives, however, must produce reasonable PDF distributions under \mathcal{H}_0 allowing the connection between the PFA and τ .

The performance of a given detector can be summarized by plotting the PD versus the PFA both computed for τ varying in the range $(-\infty, +\infty)$. This approach produces a curve that gives values for the PD and PFA for all possible values of τ allowing one to quickly access the operating characteristic of the detector in many situations. Such curve is named *Receiver Operating Characteristic* (ROC). This curve must always be above the identity function (the “45° line”). This is because the 45° ROC indicates that PD=PFA for all possible τ and then the test could be replaced by a simple flip of a coin. The ROC curve also allows to easily compare different detectors as illustrated in Figure 6.

In some applications, the distribution under \mathcal{H}_1 , i.e., $P(T|\mathcal{H}_1)$, may not be readily accessible not allowing one to produce analytical ROC curves.

In such cases, estimations of the detector performance can be accessed using empirical ROCs where the PD and PFA are computed empirically using synthetic or labeled data. Note, however, that the PD and the PFA depend on τ and thus, the empirical PD and PFA (*i.e.* PD_{Emp} and PFA_{Emp}), also need to be computed for $\tau \in (-\infty, +\infty)$. In practice, however, this range can be reduced to the range of the test statistic T . Considering N signal samples that must be assigned to \mathcal{H}_0 or \mathcal{H}_1 . Consider N_0 to be the number of samples following \mathcal{H}_0 and N_1 the number of samples following \mathcal{H}_1 . Consider that for a given threshold τ we have C_0 occurrences of false alarm, and C_1 correctly detected samples. Then, the empirical PD and PFA can be approximated as

$$\text{PD}_{\text{Emp}}(\tau) = \frac{C_1}{N_1} \quad (3.37)$$

and

$$\text{PFA}_{\text{Emp}}(\tau) = \frac{C_0}{N_0}. \quad (3.38)$$

3.2.2 Example: Detection of a DC level embedded in Gaussian noise

In this section we present a simple example to illustrate some of the basic concepts discussed above. For this, consider the following signal detection problem

$$\mathcal{H}_0 : x_i = n_i \quad (3.39a)$$

$$\mathcal{H}_1 : x_i = A + n_i \quad (3.39b)$$

where $i = 0, 1, \dots, N-1$ is the time index, x_i is the i -th sample of the observed signal which may or may not have a DC component $A = 0.25$ embedded in zero-mean WGN with power $\sigma_n^2 = 1$. Figure 7 shows the signal discussed above for $N = 200$ where the DC level is present from sample 65 to 135.

Since the noise is Gaussian distributed the PDFs for the each sample x_i under both hypotheses are clearly Gaussian and given by $P(x_i|\mathcal{H}_0) = \mathcal{N}(0, \sigma_n^2)$ and $P(x_i|\mathcal{H}_1) = \mathcal{N}(A, \sigma_n^2)$, both depicted in Figure 8. Considering the vector notation with $\mathbf{x} = [x_0, \dots, x_{N-1}]^T$ we have

$$P(\mathbf{x}|\mathcal{H}_0) = \mathcal{N}(0, \sigma_n^2 \mathbf{I}) \quad (3.40)$$

and

$$P(\mathbf{x}|\mathcal{H}_1) = \mathcal{N}(A, \sigma_n^2 \mathbf{I}). \quad (3.41)$$

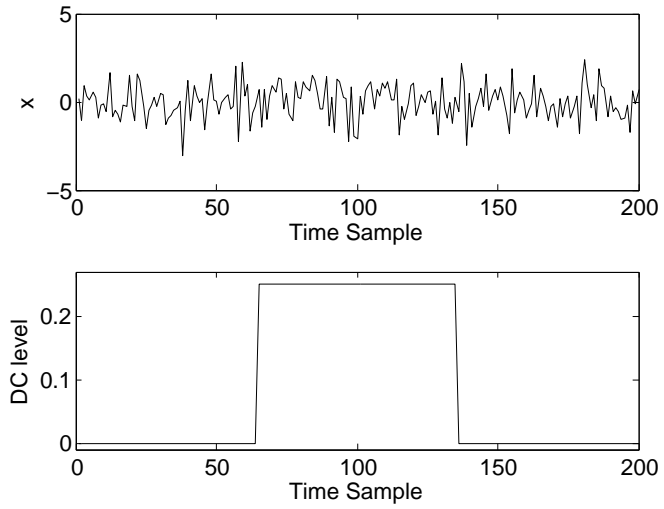


Figure 7: DC embedded in noise.

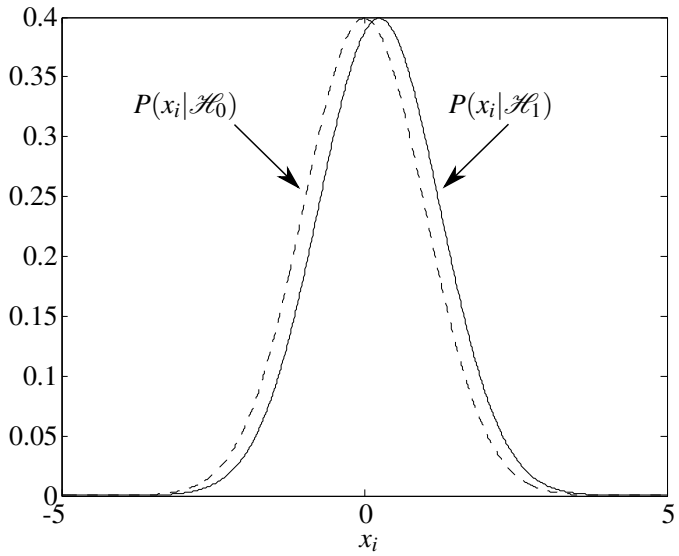


Figure 8: Hypothesis PDFs.

The NP strategy leads to the LRT that can be written as

$$L(\mathbf{x}) = \frac{\frac{1}{(2\pi\sigma_n^2)} \exp\left(-\frac{1}{2\sigma_n^2} \sum_{i=0}^{N-1} (x_i - A)^2\right)}{\frac{1}{(2\pi\sigma_n^2)} \exp\left(-\frac{1}{2\sigma_n^2} \sum_{i=0}^{N-1} (x_i)^2\right)} \underset{\mathcal{H}_1}{\overset{\mathcal{H}_0}{\leq}} \tau. \quad (3.42)$$

Taking the logarithm of both sides and making some simple algebraic simplifications we obtain

$$\frac{1}{N} \sum_{i=0}^{N-1} x_i \underset{\mathcal{H}_1}{\overset{\mathcal{H}_0}{\leq}} \frac{\sigma_n^2}{NA} \ln \tau + \frac{A}{2} = \tau' \quad (3.43)$$

where $\tau' = \frac{\sigma_n^2}{NA} \ln \tau + \frac{A}{2}$ is the new threshold which is compared to the sample mean \bar{x} as test statistic:

$$T = \frac{1}{N} \sum_{i=0}^{N-1} x_i \underset{\mathcal{H}_1}{\overset{\mathcal{H}_0}{\leq}} \tau'. \quad (3.44)$$

Note that T is Gaussian distributed under both hypotheses, that is

$$T \sim \begin{cases} \mathcal{N}(0, \sigma^2/N) & \text{under } \mathcal{H}_0 \\ \mathcal{N}(A, \sigma^2/N) & \text{under } \mathcal{H}_1. \end{cases} \quad (3.45)$$

Then the PFA and PD can be determined as

$$\text{PFA} = P(T(\mathbf{x}) > \tau' | \mathcal{H}_0) = Q\left(\frac{\tau'}{\sqrt{\sigma_n^2/N}}\right) \quad (3.46)$$

and

$$\text{PD} = P(T(\mathbf{x}) > \tau' | \mathcal{H}_1) = Q\left(\frac{\tau' - A}{\sqrt{\sigma_n^2/N}}\right), \quad (3.47)$$

where $Q(x) = 1 - \Phi(x)$, where $\Phi(\cdot)$ is the standard Gaussian cumulative distribution.

Figure 9 presents the distribution of the test statistic under both hypotheses, where the test threshold τ' was computed assuming a PFA = 0.05, and 200 samples were used to make the test, that is $N = 200$. The PD corresponds to the light gray shaded area while the PFA corresponds to the darker gray shaded area. Note that by considering several samples and LRT the resulting distribution for the test is much more separated than in Figure 8.

The ROC curve for the above detector is then presented in Figure 10, where the PD becomes 1 when PFA larger than 0.2 can be tolerated, but good performance can be achieved for low probability of false alarm.

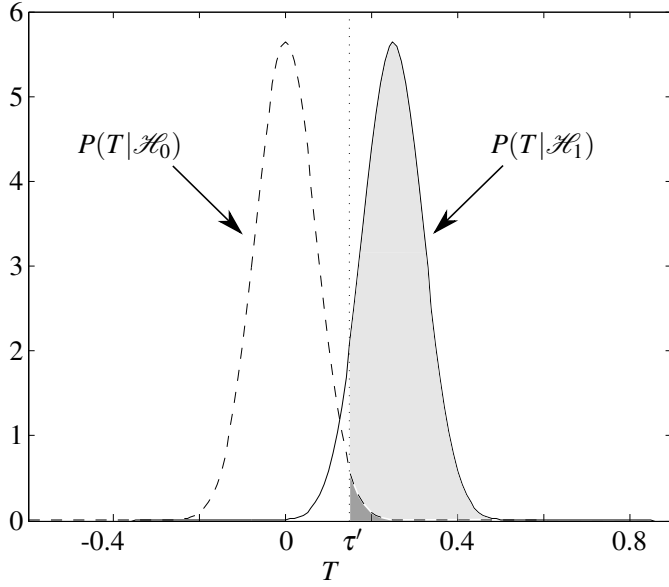


Figure 9: Test statistic PDFs. The shaded area corresponds to the probabilities of detection (light gray) and false alarm (darker gray).

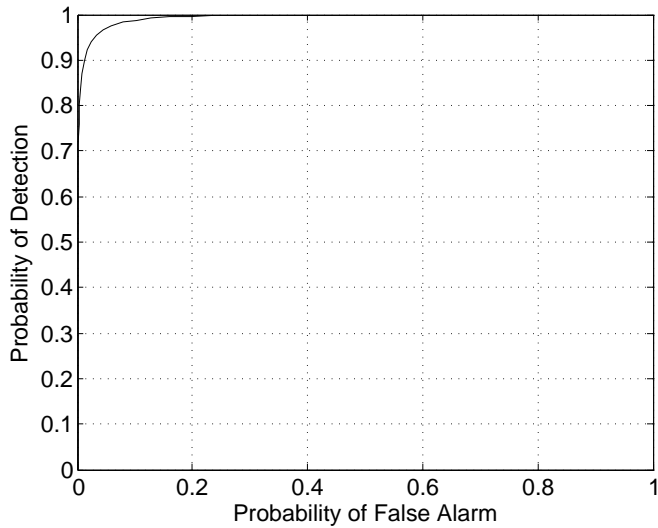


Figure 10: ROC Curve.

4 NONLINEAR MIXTURE DETECTOR

Strategies for detection of nonlinearly mixed pixels have been considered in the literature assuming strategies borrowed from analysis of nonlinear dynamical time series [81] (in this case assuming that the nonlinear behavior is along the bands), or using specific parametric nonlinear models for which the detection is performed based on the estimate of models parameters [42]. When nothing or little is known about the nonlinear mixing mechanism, a direct strategy is to exploit the property of linear mixing models to confine the noiseless data to a simplex. The hypothesis test proposed in [2] is based on the distance between the observed pixel and this simplex. Though this test is robust to nonlinear mixing mechanisms, it conveys too little information about the nonlinearity as a trade-off to guarantee simplicity. An alternative strategy is to use nonparametric techniques to extract information about the nonlinearity directly from the observations. A nonlinear mixing model for joint unmixing and nonlinearity detection is proposed in [84]. It assumes that the observed reflectances result from linear spectral mixtures corrupted by a residual nonlinear component. This model is rather similar to the model initially introduced in [13] and presented in (3.22), but the estimation method relies on a computationally intensive Bayesian procedure.

In this chapter, we present a model-free detector of nonlinenarly mixed pixels in hyperspectral images. To detect nonlinearly mixed pixels in an hyperspectral image, assuming Ψ in (2.1) is unknown, we propose to compare the reconstruction errors resulting from estimating Ψ with nonlinear and linear regression methods. As benchmark, we consider the LS-robust detector presented in [2] since it makes no strong assumption about the nonlinear mixing that actually occurs in the scene. The performance of the proposed technique is evaluated through simulations and followed by preliminary conclusions.

4.1 DETECTION OF NONLINEARLY MIXED PIXELS

4.1.1 The detection problem

Given an observation \mathbf{r} , we formulate the nonlinear mixture detector as the following binary hypothesis test problem

$$\mathcal{H}_0 : \mathbf{r} = \mathbf{M}\boldsymbol{\alpha} + \mathbf{n} \quad (4.1a)$$

$$\mathcal{H}_1 : \mathbf{r} = \boldsymbol{\psi}(\mathbf{M}) + \mathbf{n} \quad (4.1b)$$

where \mathbf{n} is a zero-mean WGN with variance σ_n^2 . We assume that the endmember matrix \mathbf{M} is available, or has been estimated from data using an endmember extraction technique [7]. We shall relax this hypothesis in Section 5.1, and use the nonlinear mixture detector to jointly perform this task.

We propose to compare the fitting errors resulting from estimating \mathbf{r} with a linear or a nonlinear estimator (3.13). Under \mathcal{H}_0 , both estimators should provide good estimates. Under \mathcal{H}_1 , the estimation error resulting from the linear estimator should be significantly larger than that obtained with the nonlinear estimator. We shall now evaluate these fitting errors.

4.1.2 Linear estimation error

The MMSE estimator (3.13) may be used with the linear kernel (3.6) to estimate $\boldsymbol{\alpha}$ in (4.1a). Nevertheless, this would require to solve (3.14) in order to estimate σ_n^2 . To avoid unnecessary computational effort, we shall limit the use of GP to nonlinear model estimation. The MMSE estimator for (4.1a) is given by:

$$\hat{\boldsymbol{\alpha}} = (\mathbf{M}^\top \mathbf{M})^{-1} \mathbf{M}^\top \mathbf{r} \quad (4.2)$$

resulting in the following estimation error:

$$\mathbf{e}_{\text{lin}} = \mathbf{r} - \hat{\mathbf{r}}_{\text{lin}} = \mathbf{P}\mathbf{r} \quad (4.3)$$

where $\mathbf{P} = \mathbf{I}_L - \mathbf{M}(\mathbf{M}^\top \mathbf{M})^{-1} \mathbf{M}^\top$ is an $L \times L$ projection matrix of rank $\rho = L - R$. Note that no constraint is imposed on the abundance vector $\boldsymbol{\alpha}$. The objective is to obtain the best linear estimator, since the purpose at this point is not to perform unmixing, but to decide on the linearity (or not) of the considered mixture.

Consider first the distribution for $\|\mathbf{e}_{\text{lin}}\|^2$. Under \mathcal{H}_1 , we have:

$$\mathbf{e}_{\text{lin}}|\mathcal{H}_1 = \mathbf{P}[\boldsymbol{\psi} + \mathbf{n}]. \quad (4.4)$$

This implies that

$$\mathbf{e}_{\text{lin}}|\mathcal{H}_1 \sim \mathcal{N}(\mathbf{P}\boldsymbol{\psi}, \sigma_n^2 \mathbf{P}) \quad (4.5)$$

where we use that the projection matrix \mathbf{P} is idempotent, that is, $\sigma_n^2 \mathbf{P}\mathbf{P}^\top = \sigma_n^2 \mathbf{P}$. Thus, the distribution for the each entry $e_{\text{lin},i}$ of \mathbf{e}_{lin} is given as

$$e_{\text{lin},i}|\mathcal{H}_1 \sim \mathcal{N}(\mathbf{p}_i^\top \boldsymbol{\psi}, \sigma_n^2 \mathbf{p}_i^\top \mathbf{p}_i) \quad (4.6)$$

where \mathbf{p}_i^\top denotes the i -th row of matrix \mathbf{P} . Under \mathcal{H}_0 , we have:

$$\mathbf{e}_{\text{lin}}|\mathcal{H}_0 \sim \mathcal{N}(0, \sigma_n^2 \mathbf{P}) \quad (4.7)$$

and

$$e_{\text{lin},i}|\mathcal{H}_0 \sim \mathcal{N}(0, \sigma_n^2 \mathbf{p}_i^\top \mathbf{p}_i). \quad (4.8)$$

Proper normalization of each squared entry $e_{\text{lin},i}$ of \mathbf{e}_{lin} yields the conditional distributions under the two hypotheses¹:

$$\begin{aligned} \frac{e_{\text{lin},i}^2}{\sigma_n^2 \mathbf{p}_i^\top \mathbf{p}_i} \Big| \mathcal{H}_1 &\sim \chi_1^2 \left(\frac{[\mathbf{p}_i^\top \boldsymbol{\psi}]^2}{\sigma_n^2 \mathbf{p}_i^\top \mathbf{p}_i} \right) \\ \frac{e_{\text{lin},i}^2}{\sigma_n^2 \mathbf{p}_i^\top \mathbf{p}_i} \Big| \mathcal{H}_0 &\sim \chi_1^2(0) \end{aligned} \quad (4.9)$$

where $\chi_n^2(\lambda)$ is the noncentral χ -square distribution with n degrees of freedom and centrality parameter λ [109].² As \mathbf{P} is idempotent and of rank $\rho = L - R$, which leads to $\|\mathbf{e}_{\text{lin}}\|^2 = \mathbf{r}^\top \mathbf{P}\mathbf{r}$, and assuming that the vector \mathbf{e}_{lin} has independent entries, we conclude that [106, p. 33]:

$$\frac{\|\mathbf{e}_{\text{lin}}\|^2}{\sigma_n^2} \Big| \mathcal{H}_0 \sim \chi_\rho^2(0). \quad (4.10)$$

¹Note that the normalization used in (4.9) produces unit-variance normal random variables.

²Given N independent zero-mean unit normal random variables U_1, \dots, U_N , and μ_1, \dots, μ_N constants, then the distribution of $\sum_{j=1}^N (U_j + \mu_j)^2$ is a non-central χ^2 distribution with N degrees of freedom and noncentrality parameter $\lambda = \sum_{j=1}^N \mu_j^2$ [109, p. 433].

4.1.3 Nonlinear estimation error with GP

Since our interest at this point is not to make predictions for new data, but to evaluate the fitting error between the model output and the available data, we define the GP estimation error as:

$$\mathbf{e}_{\text{nonlin}} = \mathbf{r} - \hat{\mathbf{r}}_{\text{nonlin}} \quad (4.11)$$

where $\hat{\mathbf{r}}_{\text{nonlin}}$ is given by (3.13) with $\mathbf{M}_* = \mathbf{M}$, and thus, $\mathbf{K}_* = \mathbf{K}$. Hence, using (3.13) in (4.11) yields

$$\mathbf{e}_{\text{nonlin}} = \mathbf{r} - \hat{\boldsymbol{\Psi}}_* \Big|_{\mathbf{M}_* = \mathbf{M}} = \mathbf{H}\mathbf{r} \quad (4.12)$$

where $\mathbf{H} = \mathbf{I}_L - \mathbf{K}^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}_L]^{-1}$ is a real-valued matrix. We shall now analyze the distribution of $\|\mathbf{e}_{\text{nonlin}}\|^2$ under hypotheses \mathcal{H}_0 and \mathcal{H}_1 . Under hypothesis \mathcal{H}_1 , we have:

$$\mathbf{e}_{\text{nonlin}} | \mathcal{H}_1 = \mathbf{H}(\boldsymbol{\psi} + \mathbf{n}). \quad (4.13)$$

This leads to the following conditional distribution

$$\mathbf{e}_{\text{nonlin}} | \mathcal{H}_1 \sim \mathcal{N}(\mathbf{H}\boldsymbol{\psi}, \sigma_n^2 \mathbf{H}\mathbf{H}^\top). \quad (4.14)$$

Thus, the distribution for each entry $e_{\text{nonlin},i}$ of $\mathbf{e}_{\text{nonlin}}$ can be written as

$$e_{\text{nonlin},i} | \mathcal{H}_1 \sim \mathcal{N}(\mathbf{h}_i^\top \boldsymbol{\psi}, \sigma_n^2 \mathbf{h}_i^\top \mathbf{h}_i) \quad (4.15)$$

where \mathbf{h}_i^\top denotes the i -th row of \mathbf{H} .

Under hypothesis \mathcal{H}_0 , the distribution for the error becomes

$$\mathbf{e}_{\text{nonlin}} | \mathcal{H}_0 \sim \mathcal{N}(\mathbf{H}\mathbf{M}\boldsymbol{\alpha}, \sigma_n^2 \mathbf{H}\mathbf{H}^\top). \quad (4.16)$$

The distribution of the i -th entry of $\mathbf{e}_{\text{nonlin}}$ is thus given by

$$e_{\text{nonlin},i} | \mathcal{H}_0 \sim \mathcal{N}(\mathbf{h}_i^\top \mathbf{M}\boldsymbol{\alpha}, \sigma_n^2 \mathbf{h}_i^\top \mathbf{h}_i). \quad (4.17)$$

Analogously to the procedure applied in (4.9), proper normalization of each

squared entry $e_{\text{nl},i}$ of \mathbf{e}_{nl} yields the following conditional distributions:

$$\begin{aligned} \frac{e_{\text{nl},i}^2}{\sigma_n^2 \mathbf{h}_i^\top \mathbf{h}_i} \Big|_{\mathcal{H}_1} &\sim \chi_1^2 \left(\frac{[\mathbf{h}_i^\top \boldsymbol{\Psi}]^2}{\sigma_n^2 \mathbf{h}_i^\top \mathbf{h}_i} \right) \\ \frac{e_{\text{nl},i}^2}{\sigma_n^2 \mathbf{h}_i^\top \mathbf{h}_i} \Big|_{\mathcal{H}_0} &\sim \chi_1^2 \left(\frac{[\mathbf{h}_i^\top \mathbf{M} \boldsymbol{\alpha}]^2}{\sigma_n^2 \mathbf{h}_i^\top \mathbf{h}_i} \right). \end{aligned} \quad (4.18)$$

Non-central χ -square distributions in (4.9) and (4.18) make the analysis of the test statistics in the next section intractable, even under \mathcal{H}_0 . In order to proceed, we argue that it is reasonable to assume that, under \mathcal{H}_0 , both the nonlinear GP regression method and the linear one should achieve the same level of accuracy. Considering (4.10), this approximation leads to

$$\frac{\|\mathbf{e}_{\text{nl}}\|^2}{\sigma_n^2} \Big|_{\mathcal{H}_0} = \chi_p^2(0). \quad (4.19)$$

We validated this approximation using extensive Monte Carlo simulations. Figures 12a and 12b illustrate this assumption for a representative example.

4.1.4 The test statistics

We propose to compare the squared norms of the two fitting error vectors \mathbf{e}_{nl} and \mathbf{e}_{lin} to decide between \mathcal{H}_0 and \mathcal{H}_1 . Also, the test statistic should allow for the adjustment of the detection threshold to a given probability of false alarm (PFA) for design purposes. Considering these two objectives, we propose the following statistical test

$$T = \frac{2\|\mathbf{e}_{\text{nl}}\|^2}{\|\mathbf{e}_{\text{nl}}\|^2 + \|\mathbf{e}_{\text{lin}}\|^2} \stackrel{\mathcal{H}_1}{\underset{\mathcal{H}_0}{\leq}} \tau \quad (4.20)$$

where τ is the detection threshold.

The reasoning behind the choice of T defined in (4.20) is as follows. Under \mathcal{H}_0 , both $\|\mathbf{e}_{\text{nl}}\|^2$ and $\|\mathbf{e}_{\text{lin}}\|^2$ are χ -square dependent random variables. Now, we write \mathbf{e}_{lin} as $\mathbf{e}_{\text{nl}} + \sqrt{2}\boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ is assumed to be also a zero-mean i.i.d. Gaussian vector³, and neglect the cross-term $\mathbf{e}_{\text{nl}}^\top \boldsymbol{\epsilon}$ when compared to $\|\boldsymbol{\epsilon}\|^2$ in evaluating $\|\mathbf{e}_{\text{lin}}\|^2$ under \mathcal{H}_0 . The latter approximation is due to the lack of correlation between \mathbf{e}_{nl} and $\boldsymbol{\epsilon}$, as the latter can be largely attributed to mismatches resulting from the numerical optimization

³The constant factor $\sqrt{2}$ is for notation purpose only.

required to solve (3.14). Under these considerations, (4.20) can be written as $T = \|\mathbf{e}_{\text{nonlin}}\|^2 / (\|\mathbf{e}_{\text{nonlin}}\|^2 + \|\boldsymbol{\epsilon}\|^2)$ with both $\|\mathbf{e}_{\text{nonlin}}\|^2$ and $\|\boldsymbol{\epsilon}\|^2$ independent and χ -square distributed. Such a statistic is known to follow a beta distribution [110].

As the GP estimator tends to fit better nonlinearly mixed data, T should be less than 1 under hypothesis \mathcal{H}_1 . Conversely, T should be close to one for linearly mixed pixels, as $\|\boldsymbol{\epsilon}\|^2$ tends to be much less than $2\|\mathbf{e}_{\text{nonlin}}\|^2$. Hence, as per (4.20), we accept hypothesis \mathcal{H}_0 if $T > \tau$ and we conclude for the nonlinear mixing hypothesis \mathcal{H}_1 if $T < \tau$.

4.1.5 Determining the detection threshold

Considering the assumption that the test statistic T has a beta distribution under \mathcal{H}_0 , a decision threshold τ can be determined for a given PFA as

$$\tau = \mathcal{B}_{\alpha,\beta}^{-1}(\text{PFA}) \quad (4.21)$$

where $\mathcal{B}_{\alpha,\beta}$ is the cumulative distribution function of the beta distribution with parameters (α, β) . The parameters of this function must be estimated from the data. To this end, we initially determine an estimate $\hat{\mathbf{A}}$ of the abundance matrix assuming the linear mixing model with the real observations $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_N]$ and the known endmember matrix \mathbf{M} . Then, using \mathbf{M} and $\hat{\mathbf{A}}$ we construct the synthetic image $\mathbf{R}_s = \mathbf{M}\hat{\mathbf{A}}$, which satisfies \mathcal{H}_0 . For this linearly mixed hyperspectral image, we then compute, say, N samples of the test statistics $T|\mathcal{H}_0$ defined in (4.20) and fit a beta distribution to these samples. The threshold τ for each PFA is then determined using (4.21). These steps are summarized in Table 2.

This procedure requires the knowledge of the endmember matrix \mathbf{M} . The next chapter proposes an iterative technique to estimate \mathbf{M} from an hyperspectral image, which we assume to contain linearly and nonlinearly mixed pixels.

4.2 SIMULATIONS

This section presents simulation results to validate the proposed approach for detecting nonlinearly mixed pixels, with synthetic images. The use of synthetic images is important as they provide a ground truth against which the performance of the detector can be verified. First, we propose a definition for a degree of nonlinearity of an hyperspectral image so that the

Table 2: Detection threshold.

For a given endmember matrix \mathbf{M} and a HI $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_N]$ the detection threshold τ can be computed for a given PFA as:

1. a synthetic image is created using the linear mixing model $\mathbf{R}_s = \mathbf{M}\hat{\mathbf{A}}$.
With the estimated abundance matrix $\hat{\mathbf{A}}$ obtained using the LS;
2. then, the test statistics $T|\mathcal{H}_0$ are computed as in (4.20) for \mathbf{R}_s ;
3. finally, a beta distribution is adjusted to the statistics $T|\mathcal{H}_0$, and τ is computed for a given PFA using (4.21).

relative performances of different detectors can be compared. This is necessary to quantify the relative energies associated with the linear and nonlinear mixing components in hyperspectral images generated with different nonlinear mixing models.

4.2.1 Degree of nonlinearity

Consider that a pixel vector can be written as the sum of a linear and a nonlinear mixing component⁴ as is the case for most existing nonlinear mixing models [38, 39, 37, 43, 18]:

$$\mathbf{r} = \mathbf{r}_{\text{lin}} + \mathbf{r}_{\text{nonlin}} \quad (4.22)$$

where \mathbf{r}_{lin} and $\mathbf{r}_{\text{nonlin}}$ are, respectively, the linear and nonlinear mixing contributions to \mathbf{r} . The energy of \mathbf{r} is given by

$$E = \|\mathbf{r}\|^2 = \|\mathbf{r}_{\text{lin}}\|^2 + 2\mathbf{r}_{\text{lin}}^\top \mathbf{r}_{\text{nonlin}} + \|\mathbf{r}_{\text{nonlin}}\|^2, \quad (4.23)$$

where $E_{\text{lin}} = \|\mathbf{r}_{\text{lin}}\|^2$ is the energy of the linear contribution and $E_{\text{nonlin}} = 2\mathbf{r}_{\text{lin}}^\top \mathbf{r}_{\text{nonlin}} + \|\mathbf{r}_{\text{nonlin}}\|^2$ is the part of the pixel energy affected by the nonlinear mixing. Given a mixing model, we define the degree of nonlinearity η_d as the ratio of the energy of the nonlinear contribution E_{nonlin} to the total energy E . Thus,

$$\eta_d = \frac{E_{\text{nonlin}}}{E} = \frac{1}{1+A} \quad (4.24)$$

⁴We do not account for noise contribution as it can be set by the user independently of the mixing model.

where $A = \|\mathbf{r}_{\text{lin}}\|^2 / (2\mathbf{r}_{\text{lin}}^\top \mathbf{r}_{\text{nonlin}} + \|\mathbf{r}_{\text{nonlin}}\|^2)$. Next, we show how to apply this definition for generating synthetic samples with two different mixing models.

4.2.1.1 Synthetic data generation with GBM

To be able to control the relative contributions of the linear and nonlinear mixing parts of the GBM model, we introduce a new scaling factor k into the generalized bilinear model (GBM) used in [2]. For an endmember matrix \mathbf{M} and an abundance vector $\boldsymbol{\alpha}$, we write the modified noiseless GBM model as

$$\mathbf{r} = k\mathbf{M}\boldsymbol{\alpha} + \gamma\mathbf{v} \quad (4.25)$$

where $0 \leq k \leq 1$, $\mathbf{v} = \sum_{i=1}^{R-1} \sum_{j=i+1}^R \alpha_i \alpha_j \mathbf{m}_i \odot \mathbf{m}_j$ is the nonlinear mixing term, γ is the scaling parameter for the nonlinear contribution, and \odot is the Hadamard product. The degree of nonlinearity is then

$$\eta_d = \frac{2k\gamma(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha}) + \gamma^2 \|\mathbf{v}\|^2}{k^2 \|\mathbf{M}\boldsymbol{\alpha}\|^2 + 2k\gamma(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha}) + \gamma^2 \|\mathbf{v}\|^2} = \frac{1}{1+A} \quad (4.26)$$

with $A = k^2 \|\mathbf{M}\boldsymbol{\alpha}\|^2 / (2k\gamma(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha}) + \gamma^2 \|\mathbf{v}\|^2)$. We have to determine the scaling factors k and γ so that the energy E is independent of $\eta_d \geq 0$. This condition can be expressed as $\|\mathbf{M}\boldsymbol{\alpha}\|^2 = k^2 \|\mathbf{M}\boldsymbol{\alpha}\|^2 + 2k\gamma(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha}) + \gamma^2 \|\mathbf{v}\|^2$, leading to

$$A = \frac{k^2}{1-k^2} \quad (4.27)$$

or

$$k = \sqrt{\frac{A}{1+A}} = \sqrt{1 - \eta_d}. \quad (4.28)$$

To obtain γ , note that the denominator of A can be written as $\gamma^2 \|\mathbf{v}\|^2 + 2k\gamma(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha}) = (1 - k^2) \|\mathbf{M}\boldsymbol{\alpha}\|^2$. Since γ must be positive, we have

$$\begin{aligned} \gamma &= \frac{1}{2\|\mathbf{v}\|^2} \left(-2k(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha}) \right. \\ &\quad \left. + \sqrt{4k^2(\mathbf{v}^\top \mathbf{M}\boldsymbol{\alpha})^2 + 4\|\mathbf{v}\|^2(1-k^2)\|\mathbf{M}\boldsymbol{\alpha}\|^2} \right). \end{aligned} \quad (4.29)$$

Once k and γ have been determined from η_d , we can generate data following the model in (4.25).

4.2.2 Synthetic data generation with PNMM

To match the noiseless PNMM model (2.21) with the proposed formulation (4.25), we complement it with a weighted linear mixture as follows:

$$\mathbf{r} = k\mathbf{M}\boldsymbol{\alpha} + \gamma\mathbf{v}, \quad (4.30)$$

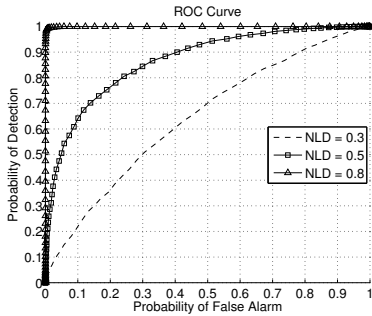
where $\mathbf{v} = (\mathbf{M}\boldsymbol{\alpha})^\xi$ denotes the exponential value ξ applied to each entry of $\mathbf{M}\boldsymbol{\alpha}$. Model (4.30) reduces to (2.21) for $k = 0$ and $\gamma = 1$. Again, parameters k and γ are scaling factors that control the relative amounts of linear and nonlinear contributions given η_d . As for the GBM, both can be set using (4.28) and (4.29).

4.2.3 Simulations with known \mathbf{M}

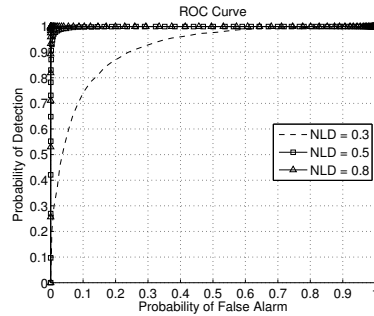
We now present simulations with synthetic data and a known endmember matrix \mathbf{M} . These simulations allow us to assess the detector performance disregarding estimation errors for the endmembers. Hence, they illustrate the potential of the proposed detector. To construct synthetic data, we used three materials ($R = 3$) extracted from the spectral library of the software ENVITM [111]: green grass, olive green paint and galvanized steel metal. Each endmember \mathbf{m}_r has $L = 826$ bands that were uniformly decimated by 3 to $L = 275$ bands.

To evaluate the performance of the proposed detector, we generated 8000 synthetic samples by mixing the three collected spectra. Among the 8000 pixels, 4000 were generated using the linear model in (2.4), and 4000 using the modified generalized bilinear model in (4.25). A fixed abundance vector $\boldsymbol{\alpha} = [0.3, 0.6, 0.1]^\top$ was used for all samples. Nonlinearly mixed samples were generated using different degrees of nonlinearity $\eta_d \in \{0.3, 0.5, 0.8\}$ to test the detector under different conditions. The power of the additive Gaussian noise was set to $\sigma_n^2 = 0.001$, which corresponds to $\text{SNR} = 21\text{dB}$.

Figure 11 shows the receiver operating characteristics (ROCs) of the proposed GP detector and the LS robust detector presented in [2] for the three values of η_d . The proposed detector performs better, especially for moderate to high degree of nonlinearity. For instance, Fig. 11c shows that the GP detector achieves a probability of detection of 1 for $\text{PFA} = 0.1$, while the LS robust detector yields a probability of detection of approximately 0.65 for the same PFA. Figure 12 shows the histograms of $\|\mathbf{e}_{\text{nlm}}\|^2$, $\|\mathbf{e}_{\text{lin}}\|^2$ and T for both linearly (\mathcal{H}_0) and nonlinearly (\mathcal{H}_1) mixed data. The proposed test statistics clearly leads to histograms that differ significantly under both hypothe-



(a) Robust LS detector.



(b) Proposed GP detector.

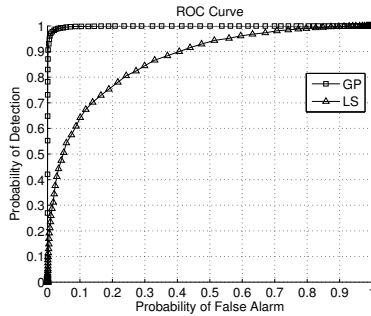
(c) Comparison of LS and GP ($\eta_d = 0.5$).

Figure 11: Empirical ROCs for: (a) the Robust LS detector [2], (b) the proposed GP detector, (c) the two detectors for $\eta_d = 0.5$. All curves were obtained for 8000 pixels (4000 linearly mixed and 4000 nonlinearly mixed) and SNR = 21dB. Nonlinear mixtures were generated using the simplified GBM described in Section 4.2.1.

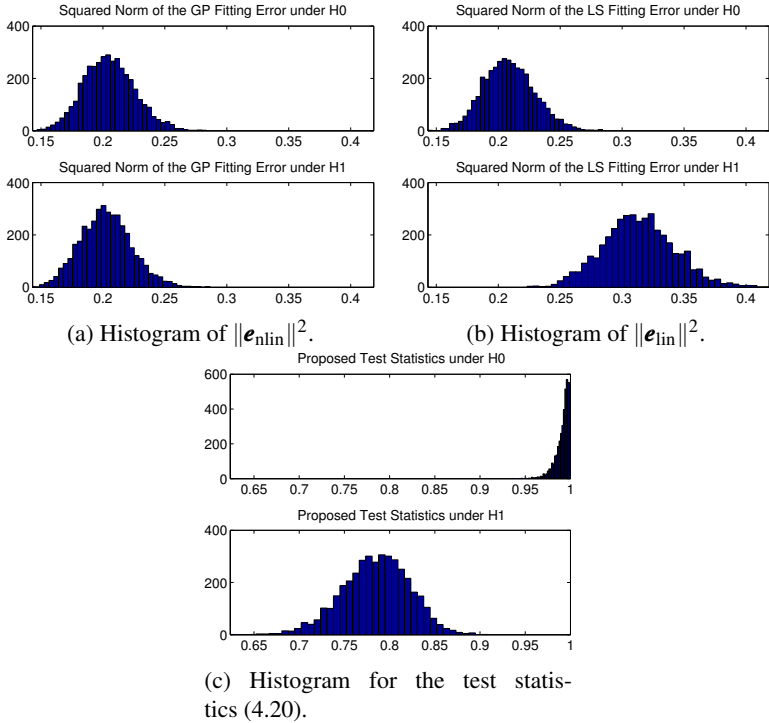


Figure 12: Histograms for (a) the squared norm of the GP fitting error, (b) the least-squares fitting error, and (c) the test statistics (4.20).

ses \mathcal{H}_0 and \mathcal{H}_1 , which explains the improvement in detection performance. Figure 13 compares the histogram of T under \mathcal{H}_0 with the fitted beta distribution, confirming that the distribution of T can be reasonably approximated by a beta distribution.

We considered two unmixing algorithms to assess the impact of the proposed detector on unmixing performance, one linear and one nonlinear. Linear unmixing was performed using the fully-constrained least-squares (FCLS) algorithm [72]. For nonlinear unmixing, we used the SK-Hype algorithm [13]. The two algorithms were employed in two unmixing strategies. First, each algorithm was used to unmix the complete hyperspectral image. In the second strategy called detect-then-unmix (D.+U.), the proposed detector (GP), and the detector of [2] (LS) were used as a pre-processing step. Then, FCLS was used to unmix pixels detected as linearly mixed and SK-Hype was used to unmix pixels detected as nonlinearly mixed. The detection threshold τ was

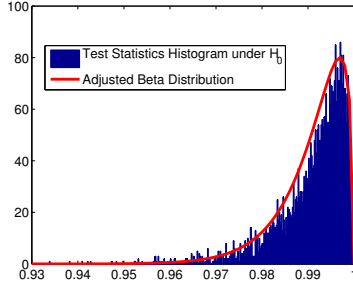


Figure 13: Histogram of the test statistics under \mathcal{H}_0 and the adjusted Beta distribution.

determined for $\text{PFA} = 0.01$. Two synthetic images were considered with 1000 pixels each, 500 being linearly mixed and 500 being nonlinearly mixed. Each image was constructed with a particular nonlinear mixing model, with a fixed degree of nonlinearity $\eta_d = 0.5$ in both cases, with abundance vectors uniformly sampled in the simplex. The GBM (4.25) was used for the first image (Image I), while the PNMM (4.30) with $\xi = 3$ was considered for the second image (Image II). The SNR was 21dB in both cases. Parameters k and γ were determined for each pixel to maintain the desired value of nonlinearity degree η_d for all simulations. To compare the results, we used the root mean square error (RMSE) of abundance estimation, defined as

$$\text{RMSE} = \sqrt{\frac{1}{NR} \sum_{n=1}^N \|\alpha_n - \hat{\alpha}_n\|^2} \quad (4.31)$$

where N is the number of pixels in each image.

The results ($\text{RMSE} \pm$ standard deviation) are presented in Tables 3 and 4. For each image, these tables indicate the RMSEs for the linearly mixed part (LMM), for the nonlinearly mixed part (NLM), and for the full image (F. Img.) using the three unmixing strategies. The results shown in bold blue are those with the lowest RMSE in each row of the tables. As expected, FCLS has the best results when unmixing linearly mixed pixels. The same observation can be made for SK-Hype with nonlinearly mixed pixels. Nevertheless, we verify that the results using the detect-then-unmix strategy and the proposed detector (D.+U. GP) are very close to the best results for both types of pixels, LMM and NLM. When processing the whole image without prior information on the mixing nature of each pixel, the best results were those obtained with the detect-then-unmix GP strategy. Additionally, we present the classification error (C. E.) in percentage for the detect-then-unmix strategy in both tables.

Table 3: Abundance estimation RMSE for \mathbf{M} known and using the GBM mixing model (SNR = 21dB, $\eta_d = 0.5$).

Image I: LMM + GBM				
Model	FCLS	SK-Hype	D.+U. GP (C.E.%)	D.+U. LS (C.E.%)
LMM	0.0095 ± 0.00010	0.0205 ± 0.00057	0.0097 ± 0.00012 (0.6)	0.0096 ± 0.00010 (0.2)
NLM	0.0624 ± 0.00384	0.0312 ± 0.00110	0.0324 ± 0.00119 (5.6)	0.0509 ± 0.00314 (51.4)
F.Img	0.0446 ± 0.00332	0.0264 ± 0.00092	0.0239 ± 0.00097 (3.1)	0.0366 ± 0.00255 (25.8)

Table 4: Abundance estimation RMSE for \mathbf{M} known and using the PNMM mixing model (SNR = 21dB, $\eta_d = 0.5$).

Image II: LMM + PNMM				
Model	FCLS	SK-Hype	D.+U. GP (C.E.%)	D.+U. LS (C.E.%)
LMM	0.0095 ± 0.00010	0.0205 ± 0.00057	0.0099 ± 0.00013 (1.2)	0.0095 ± 0.00010 (0)
NLM	0.0958 ± 0.00882	0.0440 ± 0.00204	0.0443 ± 0.00210 (0.8)	0.0483 ± 0.00276 (17)
F.Img	0.0681 ± 0.00772	0.0344 ± 0.00168	0.0321 ± 0.00176 (1)	0.0348 ± 0.00225 (8.5)

The last two columns in both tables clearly illustrate the better performance obtained using the proposed (GP) detector, as opposed to the detector of [2].

To verify the statistical significance of the results shown in Tables 3 and 4, we performed the one-tailed left nonparametric Wilcoxon signed rank test [112]. The test was performed to compare the abundance estimation RMSEs obtained with the proposed methodology (D.+U. GP) and with each of the alternative methods listed in Tables 3 and 4. The Wilcoxon signed rank test considers the samples to be paired, which corresponds to our case, and tests the following null hypothesis

$$\text{median}(\text{RMSE}_{\text{prop}}) = \text{median}(\text{RMSE}_{\text{alt}}) \quad (4.32)$$

where $\text{RMSE}_{\text{prop}}$ and RMSE_{alt} stand for the RMSEs obtained using the proposed and the alternative methods, respectively. Tables 5 and 6 show the results obtained for the simulations corresponding to Tables 3 and 4. We assigned the symbol \mathcal{A} if the null hypothesis was rejected with negative Z statistic, i.e., if there was enough evidence that $\text{RMSE}_{\text{prop}} < \text{RMSE}_{\text{alt}}$ at the 0.05 significance level. We assigned the symbol “-” if the null hypothesis could not be rejected. These results provide evidence that the improvement in abundance estimation obtained using the proposed technique is statistically consistent.

Table 5: One-tailed Wilcoxon signed rank test for Image I (Significance level 0.05).

	FCLS	SK-Hype	D.+U. LS
LMM	-	\mathcal{A}	-
NLM	\mathcal{A}	-	\mathcal{A}
F.Img.	\mathcal{A}	\mathcal{A}	\mathcal{A}

Table 6: One-tailed Wilcoxon signed rank test for Image II (Significance level 0.05).

	FCLS	SK-Hype	D.+U. LS
LMM	-	\mathcal{A}	-
NLM	\mathcal{A}	-	\mathcal{A}
F.Img.	\mathcal{A}	\mathcal{A}	\mathcal{A}

4.2.4 Simulations with an unknown endmember matrix \mathbf{M}

The simulations conducted in Section 4.2.3 assumed the endmember matrix \mathbf{M} to be known. Although this study is important to quantify the potential of the proposed detector, the endmembers are rarely known in practice. Hence, in this section, we study the sensitivity of the detection performance as a function of the endmember estimation accuracy and of the degree of nonlinearity. Endmember extraction was performed with the iterative method proposed in Section 5.1, and with VCA [23] for comparison.

Figure 14 presents the results of 4 experiments using synthetic images with 5000 samples, SNR = 21dB, abundances uniformly sampled in the simplex, a proportion of nonlinearly mixed pixels in the image varying from 10% to 50%, and $\eta_d = 0.5$. For every experiment, the endmember matrix was extracted using VCA. These results show how the detection performance can degrade as the number of nonlinear pixels increases and as VCA loses accuracy in extracting the endmembers from the image. These results confirm the importance of devising alternatives to VCA (or to other endmember extraction algorithms specifically designed for linearly-mixed images) for images containing nonlinearly-mixed pixels.

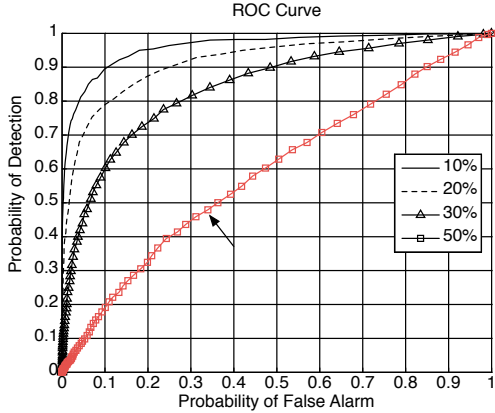


Figure 14: ROCs for different proportions of nonlinearly mixed pixels and $\eta_d = 0.5$. Endmember extraction using VCA.

4.3 PRELIMINARY CONCLUSIONS

In this chapter we presented a nonparametric method for detecting nonlinear mixtures in hyperspectral images. The performance of the detector was studied for supervised and unsupervised unmixing problems. When the endmember matrix is known, we showed that the improvement in the unmixing performance obtained when using the proposed detector is statistically consistent. Additionally, a degree of mixture nonlinearity based on the relative energies of the linear and nonlinear contributions to the mixing process was defined to quantify the importance of the linear and nonlinear model counterparts. Such a definition is important for a proper evaluation of the relative performances of different nonlinear mixture detection strategies. Finally, when considering unsupervised mixing problems we showed that the detection performance can degrade as the number of nonlinear pixels increases and as VCA loses accuracy in extracting the endmembers from the image. These results confirm the importance of devising alternatives to endmember extraction algorithms specifically designed for nonlinearly-mixed images.

5 EEA FOR NONLINEARLY MIXED HYPERSPECTRAL IMAGES

The presence of nonlinearly mixed pixels in a hyperspectral image tends to degrade the estimation accuracy of endmember extraction methods based on a linear mixing model. As a consequence, nonlinearly mixed pixels also affect the performance of algorithms using the endmember matrix such as the detection method presented in Chapter 4. There has been few papers addressing endmember estimation from nonlinearly mixed images. A nonlinear unmixing algorithm is derived in [73]. The pixel reflectances are supposed to be post-nonlinear functions of unknown pure spectral components. A Bayesian strategy is proposed to both unmix the data and estimate the endmembers. Both tasks are however mutually dependent and the unmixing model is very specific. A nonlinear endmember estimation algorithm based on the approximation of geodesic distances is introduced in [27, 78]. This algorithm can however suffer from the absence of pure pixels in the image, and the effectiveness of using manifold learning methods on real data still needs to be analyzed and confirmed. In this chapter, we propose an iterative technique for estimating the endmember matrix \mathbf{M} under the reasonable assumptions that the number R of endmembers is known [113, 114, 115], and that these endmembers are linearly mixed within at least a small part of the image. Nonlinear mixtures may however compose a significant part of the image. The proposed technique combines the detector of nonlinearly mixed pixels presented in Section 4.1 and the endmember estimation algorithm known as *Minimum Volume Enclosing Simplex* (MVES) [25].

5.1 ENDMEMBER EXTRACTION IN NONLINEARLY MIXED HYPERSPECTRAL IMAGES

The proposed procedure is described in Algorithm 1. It is a two-step iterative algorithm, and called *Iterative Endmember Estimation* (IEE) Algorithm. The first step consists of using MVES to estimate the endmembers (line 2 and 14 in Algorithm 1). The second step uses (4.20) to compute the detection statistics for all the L pixels in the image \mathbf{R}_{tmp} (line 7 in Algorithm 1). Then, all pixels detected as nonlinearly mixed, that is, whose detection statistic satisfies $T(i) \leq \tau_r$ are removed (line 9), where $\tau_r = r_f \times \tau$ (line 4 and 11) is the relaxed detection threshold. The use of a relaxed threshold is suggested to avoid discarding linear pixels during the first iterations, when the estimates of \mathbf{M} are still not sufficiently accurate. The relaxing fac-

tor is initialized for $r_f = 0.9$ and is increased by a factor $r_{inc} = 0.1/N_{max}$ at each iteration to improve pixel selection as the estimation of the matrix \mathbf{M} improves (line 10). The procedure is repeated until the linear and the non-linear GP models in (4.20) present similar fitting errors within the limit of ε . Using this procedure, τ_r tends to the desired threshold τ as the estimation of \mathbf{M} improves, leaving mostly linear pixels for which both models have similar performance. A maximum number of iterations N_{max} is also set to avoid discarding too much data.

Note that we have opted for the MVES algorithm for endmember extraction because it inscribes the data into a minimum-volume simplex. Thus, MVES is suitable to estimate \mathbf{M} in the absence of pure pixels. This feature is specially interesting for our purpose since the procedure described above discards data, which may even be pure or near-pure pixels during the first iterations. Nevertheless, any other endmember estimation algorithm valid in absence of pure pixel [116, 117, 118] could be potentially used with Algorithm 1.

Algorithm 1: Iterative endmember estimation (IEE)

Input : The hyperspectral image \mathbf{R} , and the number of endmembers R

Output: Estimated endmember matrix $\widehat{\mathbf{M}}$

- 1 Initialization: $T_{max} = 1, T_{min} = 0, \varepsilon = 0.05, \mathbf{R}_{tmp} = \mathbf{R}, N_{max} = 10, cc = 0, r_f = 0.9, r_{inc} = (1 - r_f)/N_{max}, PFA = 0.05;$
- 2 $\widehat{\mathbf{M}} = \text{MVES}(\mathbf{R}_{tmp}, R);$
- 3 Compute τ using (4.21);
- 4 $\tau_r = r_f \times \tau;$ %% (relaxed threshold)
- 5 **while** $T_{max} - T_{min} > \varepsilon$ & $cc < N_{max}$ **do**
- 6 **for** $i = 1$ **to** N_{pixels} **do**
- 7 | Compute $\mathbf{T}(i)$ using (4.20);
- 8 **end**
- 9 Remove all pixels with $\mathbf{T}(i) \leq \tau_r$ from $\mathbf{R}_{tmp};$
- 10 $r_f = r_f + r_{inc};$ %% (relaxing factor)
- 11 $\tau_r = r_f \times \tau;$
- 12 $T_{max} = \max(\mathbf{T}); T_{min} = \min(\mathbf{T});$
- 13 $cc = cc + 1;$
- 14 $\widehat{\mathbf{M}} = \text{MVES}(\mathbf{R}_{tmp}, R);$
- 15 **end**

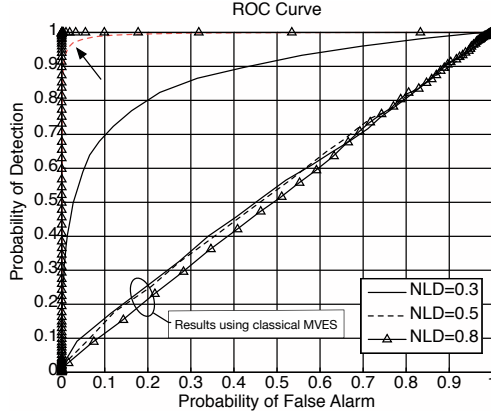


Figure 15: ROCs for different degrees of nonlinearity η_d and 50% of nonlinearly mixed pixels in the image. Endmember extraction using Algorithm 1.

5.2 SIMULATIONS

5.2.1 Simulations with an unknown endmember matrix \mathbf{M}

The simulations conducted in Section 4.2.4 (\mathbf{M} unknown) showed how the detection performance can degrade as the number of nonlinear pixels increases and as VCA loses accuracy in extracting the endmembers from the image. Figure 15 presents the results obtained with Algorithm 1 and classical MVSE for endmember extraction. For this experiment, we generated data with 50% of nonlinearly mixed pixels and different degrees of nonlinearity $\eta_d \in \{0.3, 0.5, 0.8\}$. The corresponding cases for $\eta_d = 0.5$ and 50% of nonlinearly mixed pixels are shown in red and pointed by arrows in Figure 14 and Figure 15. The poor results obtained using classical MVSE are also indicated. Comparing Fig. 11 and Fig. 15 shows that the results obtained with the iterative endmember extraction algorithm are very close to those obtained for a known endmember matrix \mathbf{M} (which can be considered as the reference detector).

Figure 16 illustrates a representative example of evolution obtained with the proposed iterative endmember extraction algorithm. These plots correspond to a simulation performed using 1000 synthetic samples, 500 being linearly mixed and 500 being nonlinearly mixed. The nonlinearly mixed pixels were created using the GBM (4.25) with $\eta_d = 0.5$. The data were

projected onto the space spanned by the columns of the current endmember matrix \mathbf{M} . They are represented as black dots. The current endmembers are shown as green dots. The true endmembers are shown as black circles at the vertices of the true simplex drawn with black lines. The data discarded at each iteration are shown within blue circles. Figure 16a shows the first iteration of Algorithm 1. Numerous nonlinear samples are outside the simplex and endmember are poorly estimated. The situation improves in Fig. 16b, which depicts the fourth iteration. Here, much less data lie outside the simplex, and two of the endmember estimates have improved significantly. Similar improvement can be noticed in the seventh iteration in Fig. 16c. The final result obtained after 10 iterations only is shown in Fig. 16d, where most of the nonlinear data were discarded and the endmember estimates are clearly close to the true endmembers.

5.2.2 Choosing the parameters r_f , N_{\max} , and ε

The implementation of Algorithm 1 requires the choice of parameters N_{\max} , ε , and r_f . We have found from several experiments that using $r_f \in [0.8, 0.9]$, $\varepsilon = 0.05$ and $N_{\max} = 10$ is a good choice for different scenarios. This section explores the sensitivity of the algorithm performance to variations of these parameter values about these choices. To this end we applied the algorithm to synthetic data with the following properties: 100 pixels, $R = 3$ endmembers, 50 pixels mixed with the LMM and 50 pixels mixed with the GBM with $\eta_d = 0.5$. The abundance vectors were sampled uniformly in the simplex, and WGN was added to the scene to produce an SNR of 21dB. The spectra used were the same used for the previous simulations, uniformly decimated by 5, resulting in 166 bands. The three parameters were chosen from the following sets: $N_{\max} \in [5, 10, 15]$, $\varepsilon \in [0.01, 0.05, 0.1]$, and $r_f \in [0.7, 0.8, 0.9]$. For each combination of parameters we performed $N_r = 900$ runs of Algorithm 1, and computed the RMSE of endmember estimation using (4.31) with the abundance vectors replaced with the endmembers. Table 7 shows the obtained results. The best results were obtained for each pair (r_f, ε) are highlighted in bold blue. These results show that the performance of the algorithm is not very sensitive to different parameter choices. They also show that choosing $N_{\max} < 10$ tends to increase the RMSE. Furthermore, it is clear that choosing $r_f < 0.8$ tends to require larger values of N_{\max} .

The choice of the parameters should be directed to prevent the algorithm from an early convergence with elimination of a large amount of linearly mixed pixels along with the nonlinearly mixed ones. This can be

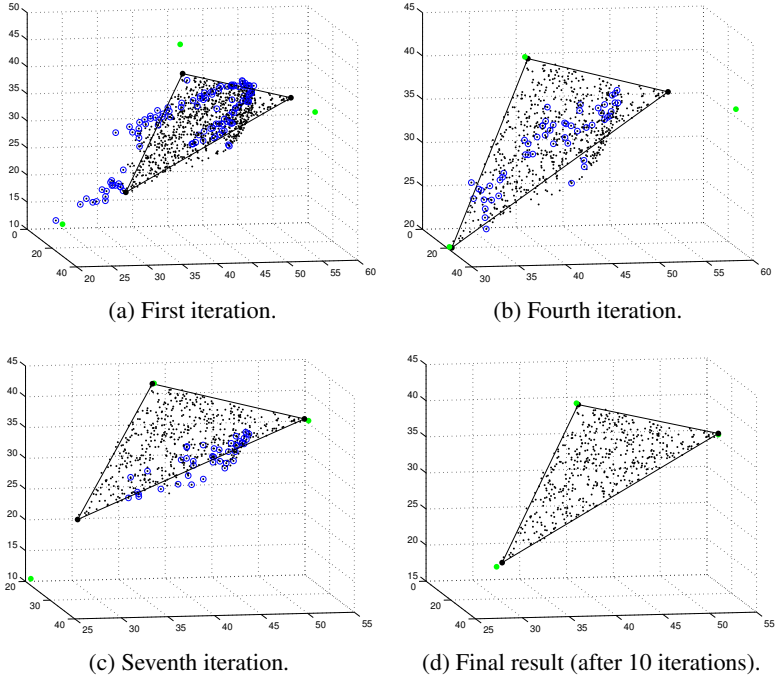


Figure 16: Graphical illustration of the endmember estimation process using the proposed iterative algorithm. The data set consists of 2000 pixels, with a proportion of 50% nonlinearly mixed pixels obtained with the GMB model and $\eta_d = 0.5$. Green dots are the current estimated endmembers, and black dots are the data projected onto the subspace spanned by the columns of the current matrix \mathbf{M} . The true endmembers are shown as black circles at the vertices of the true simplex drawn with black lines. The data discarded at the corresponding iteration are shown within blue circles.

Table 7: Mean RMSE for endmember estimation.

r_f	ε	$N_{\max} = 5$	$N_{\max} = 10$	$N_{\max} = 15$
0.7	0.01	0.0825 \pm 0.0384	0.0784 \pm 0.0378	0.0778 \pm 0.0408
	0.05	0.0817 \pm 0.0377	0.0788 \pm 0.0396	0.0778 \pm 0.0410
	0.1	0.0832 \pm 0.0383	0.0808 \pm 0.0437	0.0783 \pm 0.0398
0.8	0.01	0.0821 \pm 0.0401	0.0783 \pm 0.0416	0.0753 \pm 0.0393
	0.05	0.0805 \pm 0.0387	0.0766 \pm 0.0400	0.0778 \pm 0.0435
	0.1	0.0819 \pm 0.0416	0.0758 \pm 0.0367	0.0801 \pm 0.0406
0.9	0.01	0.0776 \pm 0.0428	0.0738 \pm 0.0399	0.0702 \pm 0.0394
	0.05	0.0764 \pm 0.0379	0.0741 \pm 0.0420	0.0744 \pm 0.042
	0.1	0.0787 \pm 0.0401	0.0785 \pm 0.0355	0.0780 \pm 0.0344

achieved by setting N_{\max} to a sufficiently large value, which controls both the maximum number of iterations and the increment of the detection threshold τ . From our experience with the proposed method, good results can be obtained as follows:

- Set r_f somewhere in the range [0.8 0.9] (Remark: a larger value would probably lead to early discarding of linearly mixed pixels).
- Set $\varepsilon \leq 0.05$, so that \mathbf{R}_{tmp} would contain basically linearly mixed pixels when the condition $T_{\max} - T_{\min} > \varepsilon$ is satisfied.
- Secure the algorithm stopping with mostly linearly mixed pixels if condition (b) cannot be satisfied by setting $N_{\max} \geq 10$.

5.2.3 Simulation with synthetic data extracted from a real scene

In this section we evaluate the performance of the proposed method using synthetic data that carries the characteristics of real data. While tests using real data are important, the use of synthetic data (for which the ground truth is known) is necessary for a more comprehensive evaluation. To conciliate both needs, we considered a scene corresponding to the alunite hill (depicted in Figure 18a) extracted from the 1997 AVIRIS scene from the Cuprite mining site in Nevada [119]. The chosen region is indicated in Figure 17. The alunite hill site has two interesting properties. First, it has a known number of endmembers ($R = 3$), i.e. alunite, muscovite, and kaolinite. Second,

Table 8: RMSE for the abundances in the alunite hill scene.

Algorithm	RMSE \pm STD (C. E. %)
FCLS	0.0797 \pm 0.0123 (-)
SK-Hype	0.0824 \pm 0.0059 (-)
detect-then-unmix	0.0671 \pm 0.0049 (3.83)

this scene has been accurately unmixed using linear mixing models [69]. To build the synthetic image we used the MVES to linearly unmix the pixels in the image. The reconstructed image considering the LMM is depicted in Figure 18b. The reconstructed image is clearly very similar to the original image, and thus carries its characteristics. To obtain a partly nonlinearly mixed image, we randomly selected 30% of the pixels from the reconstructed image and re-mixed them using the modified GBM model (4.25) with $\eta_d = 0.3$, but preserving the abundances. Finally, we added a WGN to each pixel with power adjusted to produce a 30dB SNR, which is typical for hyperspectral images. The resulting synthetic image is shown in Figure 18c. This is a partly nonlinearly mixed image for which we know the ground truth and that carries the characteristics of a real image.

We applied the proposed EEA to the image of Figure 18c and compared the endmember estimates with those obtained by applying the MVES and the VCA algorithms directly to the image. We considered $N_{\max} = 10$, $\varepsilon = 0.05$, and $r_f = 0.9$. The results are shown in Figure 19. It can be verified that the proposed method has led to the most accurate endmember estimates even for a moderate degree of nonlinearity. Figure 20 shows in black the real endmembers and the data projected into the column space of \mathbf{M} . The endmember estimates calculated by proposed EEA after 10 iterations are shown in blue. This figure clearly shows the challenging problem posed to the algorithm, as the chosen degree of nonlinearity introduces a relatively small detachment of the nonlinearly mixed pixels from the simplex. Table 8 presents the RMSE for the abundance vectors using the endmembers estimated with the proposed EEA (labelled “detect-then-unmix”) and with two alternative unmixing strategies: linear with the FCLS, and nonlinear with the SK-Hype. The improvement obtained using the proposed method is of the order of 18%. For a visual evaluation, Figure 21 compares the true nonlinearity map with the detection map. The white and gray pixels were correctly classified, and the black pixels were misclassified.

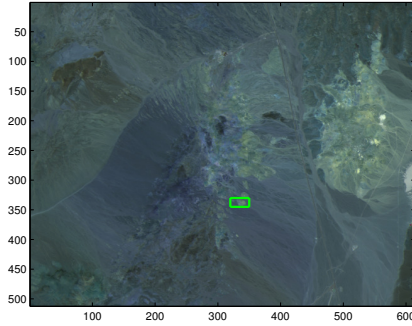
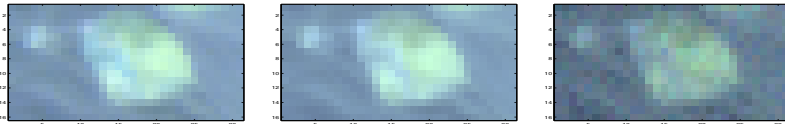
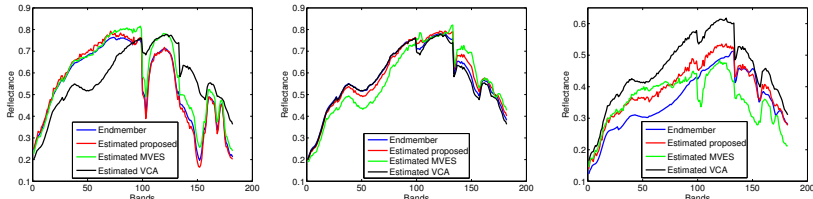


Figure 17: Cuprite mining site. The green box corresponds to the alunite hill scene.



(a) Alunite hill. (b) LMM. (c) GBM + WGN.

Figure 18: (a) Plot of the alunite hill with bands 30, 70 and 100. (b) Reconstruction of the scene using the LMM. (c) Adding 30 % of nonlinearly mixed pixels and WGN to give a 30dB SNR.



(a) Alunite. (b) Kaolinite. (c) Muscovite.

Figure 19: Endmember estimations for the nonlinearly mixed image with different extraction techniques.

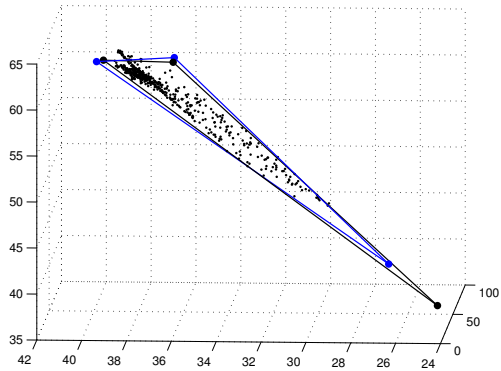


Figure 20: The black circles are the real endmembers, the black dots are the data projected in the columns of \mathbf{M} . The blue circles are the estimated endmembers with the proposed algorithm after 10 iterations. The simplex for the “true” and estimated endmembers are also drawn.

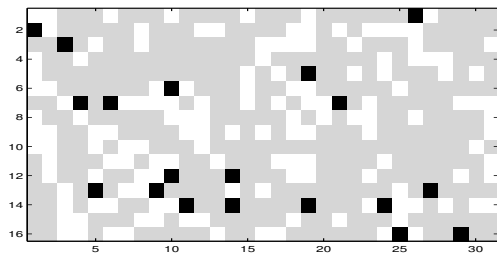


Figure 21: Detection map and true nonlinear map. Linearly mixed pixels in gray, nonlinearly mixed pixels in white, and misclassified pixels in black.

5.2.4 Real Data

5.2.4.1 Indian Pines

To test the proposed method using real images, we used the data set available at the Indian Pines test site in North-western Indiana [120]. This image was captured by the AVIRIS (Airborne Visible/Infrared Imaging Spectrometer). It has 145×145 samples over 220 contiguous bands with wavelengths ranging from 366 to 2497 nm. Prior to analysis, noisy and water absorption bands were removed resulting in a total of 200 bands that were uniformly decimated to 50 to speed up simulations. The data set has a ground truth map that divides the samples into 16 mutually exclusive classes. In Table 9, the classes are organized by numbers (1 to 16), and the number of samples of each class is shown. Note, however, that the number of samples in each class can vary considerably. Note also that some classes are composed of different materials. We can count 20 different materials if we consider grass as an isolated material for the whole image. We chose to count each grass (depending on the accompanying material) as a different material, leading to 22 endmembers. Figures 22a and 22b display images from the Indian Pines region constructed by selecting three different bands, while Fig. 23a presents the ground truth map for this image, where each class is represented by a different color. In Figure 23a, we also indice the class number for each area, where 0 represents the background, which is an unclassified area.

To perform the simulations, we divided the image into eight sub-images to work with smaller areas of the image and to deal with 3 to 4 endmembers at a time. To define these sub-images, we also paid attention to balance the number of samples per endmember. By looking at Figs 22a and 22b, we can note that some classes seem to have materials that are not accounted for in the available ground-truth information. For instance, this is the case for classes 5, 11 and 14. Therefore, we introduced extra endmembers for some of the sub-images. Table 10 describes how the sub-images were organized, showing the classes, materials, numbers of pixels and endmembers chosen for each of the eight sub-images.

For each sub-image, we estimated the endmembers as discussed in Section 5.1, with $N_{\max} = 10$, a relaxing factor initially set to $r_f = 0.8$, and incremented by $r_{\text{inc}} = (1 - r_f)/N_{\max} = 0.2$ at each of the 10 iterations. Then, we ran the detection algorithm with $\text{PFA} = 0.001$. Since we subdivided the real image into different sub-images, some of which have few pixels, we employed a more relaxed value of r_f when compared to previous simulations to avoid discarding too much data in the first few iterations. Moreover, natu-

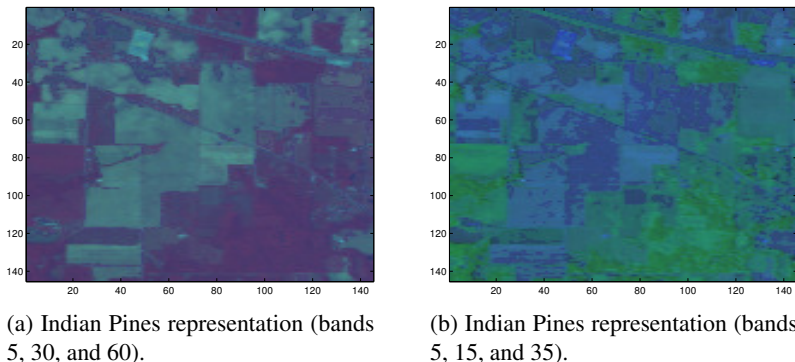


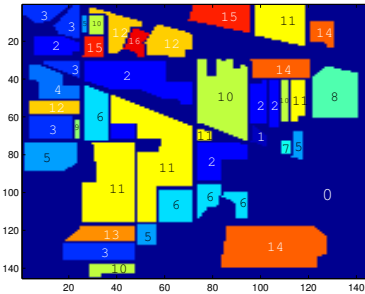
Figure 22: Indian Pines test site representation selecting 3 different bands in (a), and 3 other bands in (b).

ral phenomena such as endmember variability, wrong (or incomplete) ground truth and illumination factors (among others) tend to degrade the detection performance when dealing with real images, specially when considering non-linear algorithms which are more susceptible to overfitting. Thus, we have employed a smaller PFA to minimize incorrect detections of linearly mixed samples as nonlinearly mixed. We performed the unmixing step using FCLS for pixels detected as linearly mixed and SK-Hype for pixels detected as nonlinear mixtures. Figure 23b presents the detection map superimposed to the ground-truth classes, where black dots represent pixels detected as nonlinearly mixed.

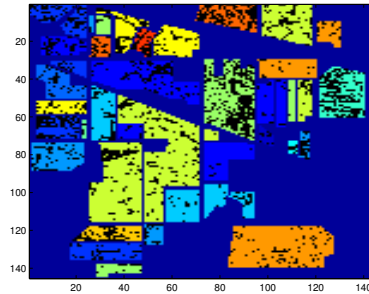
Comparing the detection map in Fig. 23b with Figs 22a and 22b, one can note similarities between the detection map and some patterns observed in the image representations. For instance, the black triangular shape in class 11 in Fig. 23b (centered about coordinate (40,80)) is just besides what seems to be a road or trail when looking to Figure 22a. Similarities can be found between contours of detected nonlinear regions in Fig. 23b and the corresponding regions in Figs 22a or 22b. Table 11 reports the RMSEs for the reconstruction error for each of the eight sub-images using three approaches, namely FCLS, SK-Hype, and detect-then-unmix. The results marked in bold blue correspond to the lowest RMSEs. For almost all sub-images, we note that the use of a nonlinear mixture detector improved the image reconstruction when compared to the pure linear or pure nonlinear unmixing strategies.

Table 9: Indian Pines classes by region.

Class number	Class	Num. of Samples
1	Alfalfa	46
2	Corn-notill	1428
3	Corn-mintill	830
4	Corn	237
5	Grass-pasture	483
6	Grass-trees	730
7	Grass-pasture-mowed	28
8	Hay-windrowed	478
9	Oats	20
10	Soybean-notill	972
11	Soybean-mintill	2455
12	Soybean-clean	593
13	Wheat	205
14	Woods	1265
15	Buildings-Grass-Trees-Drives	386
16	Stone-Steel-Towers	93



(a) Indian Pines ground truth.



(b) Indian Pines detection map.

Figure 23: Detection of nonlinearly mixed pixels in Indian Pines hyper-spectral image. Black pixels were detected as nonlinearly mixed ones by the proposed detector.

Table 10: Subimages organization

Subimg.	Classes	Materials	# of pxls.	# of endmem.
1	9 and 7	Oats and grass-pasture-mowed	48	3
2	1, 4 and 13	Alfafa, wheat and corn	488	3
3	16	Stone-steel-towers	93	3
4	15	Buildings-grass-trees-drives	386	4
5	5	Grass-Pasture	483	3
6	8 and 12	Hay-windrowed and Soybean-clean	1071	3
7	3,6 and 10	Corn-mintill, grass-trees and soybean-notill	2532	4
8	14 2 11	Woods, corn-notill, soybean-mintill	5148	4

Table 11: Indian Pines reconstruction error (RMSE) by subimage.

Subimg.	RMSE \pm STD		
	FCLS	SK-Hype	detect-then-unmix
1	0.0028627 \pm 6.6939e-06	0.0030332 \pm 6.0053e-06	0.0029083 \pm 6.5229e-06
2	0.0038963 \pm 1.2293e-05	0.003881 \pm 9.4813e-06	0.0038391 \pm 1.1505e-05
3	0.0044259 \pm 2.9087e-05	0.0035981 \pm 8.9722e-06	0.0035537 \pm 9.8622e-06
4	0.0040145 \pm 1.1417e-05	0.0039097 \pm 8.0165e-06	0.0038895 \pm 8.5058e-06
5	0.0030848 \pm 7.0516e-06	0.0032353 \pm 5.9761e-06	0.0030527 \pm 6.2275e-06
6	0.0039905 \pm 6.5627e-06	0.004055 \pm 7.1531e-06	0.0039644 \pm 6.6603e-06
7	0.0034804 \pm 5.8657e-06	0.0035049 \pm 5.9207e-06	0.0034552 \pm 5.9632e-06
8	0.0037665 \pm 7.5723e-06	0.0039314 \pm 7.3092e-06	0.0037531 \pm 7.4932e-06

5.2.4.2 Cuprite

This example applies the proposed EEA (Algorithm 1) to real data from a scene extracted from the Cuprite Mining site in Nevada (Figure 26a). This scene was captured by the AVIRIS instrument and has originally 224 bands. We removed the water absorption bands and decimated the data uniformly by a factor of 2, resulting in 94 bands. The decimation was carried out to speed up simulations. As reference spectra we selected 18 spectral signatures taken from the 1998 USGS spectral library. These spectral signatures were selected based on minerals reported to be present in the Cuprite Mining Field [23, 121, 122]. We estimated the number of endmembers using *Virtual Dimension* (VD) [113] with probability of false alarm $P_f = 10^{-4}$, resulting in $R = 5$ endmembers. We performed the endmember estimation using the proposed EEA (IEE), as well as VCA and MVES. We considered also a modification of Algorithm 1 where we replaced the proposed detector with the robust least-squares based detector presented in [2]. We refer to this method as LS for short. The parameter setting for the proposed EEA was $N_{\max} = 10$, $\varepsilon = 0.05$, and $r_f = 0.7$.¹ We searched the 18 USGS spectra for the best match (smaller spectral angle) with the endmembers extracted. The endmembers were identified as Sphene, Montmorillonite, Kaolinite, Dumortierite, Pyrope. These endmembers have strong components in this part of the Cuprite Mining Field [23]. Figure 24 shows the endmembers estimated with the proposed EEA (red line), with the LS (green lines), and the best matched signatures from USGS spectral library (blue lines). Table 12 lists the spectral angles, in radians, between the estimated and the library endmembers for the proposed EEA, LS, VCA, and MVES². Clearly, the proposed method presented good estimation performance, outperforming the other methods. Figure 25 presents the abundance maps for the unmixing process using the detect-then-unmix strategy with the GP detector. These abundance maps are in good agreement with abundance maps estimated in [23]. Figure 26 presents the reconstruction error (RMSE) for the Cuprite scene using the proposed EEA (Fig. 26b) and the VCA (Fig. 26c). In both cases the unmixing procedure was carried out using the SK-Hype algorithm. The darker tone dominating Figure 26b indicates a better fitting of the model when compared with Figure 26c. This result is corroborated by the smaller RMSE obtained using the

¹The parameter r_f for the LS detector case was modified to 1.2 to adjust the algorithm to the least-squares detector.

²Note that the mean spectral angle error used in [23] and [121] as a quality measure for the endmember estimation can be thought as a weighted mean projection of all the image vectors on the estimated endmembers, and therefore does not capture nonlinear relations between pixels and endmembers.

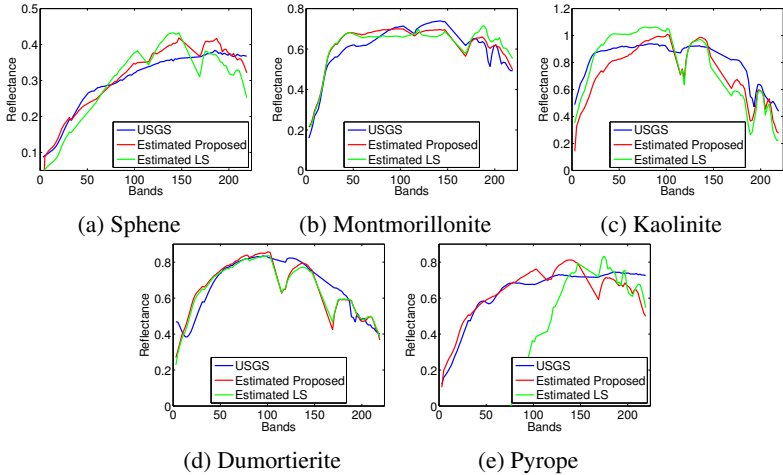


Figure 24: Estimated endmembers and USGS spectra.

Table 12: Spectral angles (in rad) between estimated and USGS spectra.

Endmemeber	IEE	LS	VCA	MVES
Spheene	0.0799	0.1498	0.3634	0.2457
Montmorillonite	0.0615	0.0852	0.0888	0.0762
Kaolinite	0.1471	0.1689	0.2022	0.2559
Dumortierite	0.1054	0.1008	0.0942	0.1422
Pyrope	0.1035	0.9792	0.1760	0.1588

proposed method ($\text{RMSE}_{\text{prop}} = 0.0040$, $\text{RMSE}_{\text{VCA}} = 0.0051$).

5.3 PRELIMINARY CONCLUSIONS

In this chapter an iterative algorithm was derived for endmember estimation as a pre-processing step for unsupervised unmixing problems. It was shown that the combined use of the detector presented in Chapter 4 and end-member estimation algorithm leads to better unmixing results when compared to state-of-the-art solutions. Simulations using different scenarios corroborate the conclusions.

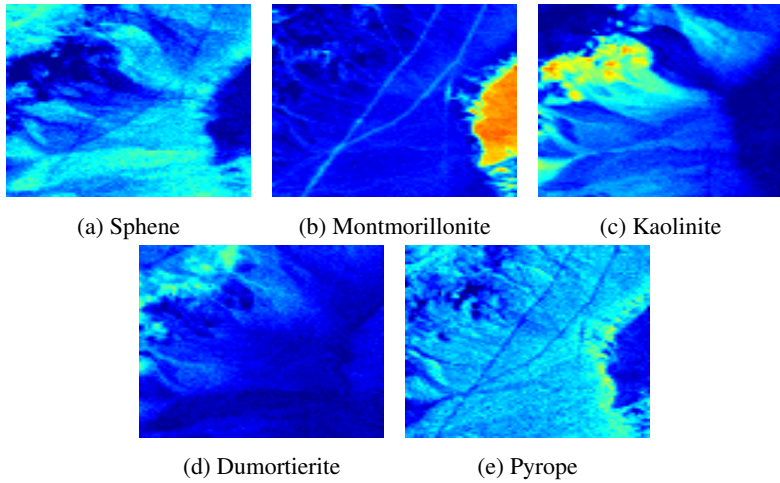


Figure 25: Abundance maps.

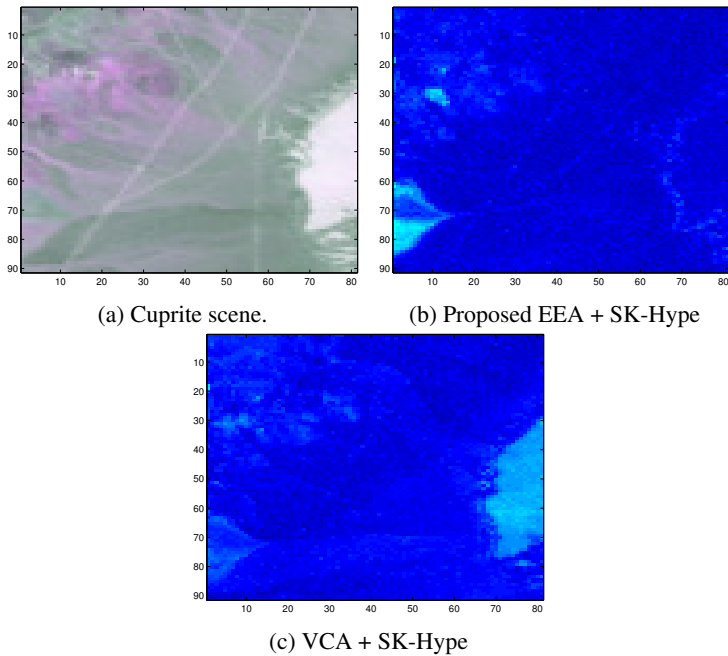


Figure 26: Cuprite scene and reconstruction errors.

6 BAND SELECTION IN RKHS

Band selection has been an active topic of research for classification of spectral patterns, see [86, 87, 88, 89, 90] and references therein. Subspace projection techniques [92, 23, 93] tend, however, to be preferred over BS [91, 123] for reducing the complexity of linear unmixing processes. They use the property that high-dimensional hyperspectral data are confined to a low-dimensional simplex in linearly-mixed images with only a few endmembers [6]. This assumption becomes invalid when nonlinear mixing phenomena are involved. Recently, in a preliminary work [30], we introduced a BS strategy method that employs the kernel k -means algorithm to identify clusters of spectral bands in the RKHS where nonlinear unmixing is performed. The HU results obtained were encouraging. One drawback of the approach in [30] is the need for an arbitrary choice of the order of the nonlinear model (the dimension of the dictionary). Given the order, band selection is performed based on the distances among different bands in the RKHS. Hence, the optimality of the solution is not driven by any direct measure of modeling accuracy. A new coherence-based method for BS in the RKHS was introduced in [31]. The coherence criterion is used to set the largest correlation between the basis kernel functions included in the unmixing model. We show that this BS approach is equivalent to search for a maximum clique in a graph, that is, the largest complete subgraph in this graph. Starting from a tentative dictionary cardinality, the proposed method determines both the dictionary size and its elements in order to satisfy the required coherence criterion. Using the maxCQL algorithm [124] to solve the maximum clique problem, the new method results in dictionaries of kernel functions, and thus spectral bands, that are less coherent than those obtained using kernel k -means initialized with dictionaries of the same size.

In this chapter, we present both methods. First, we review the kernel k -means approach. Then the strategy based on the so-called coherence criterion [125] and maximum clique search in a graph is presented. Although these two approaches are connected, they differ in their formulation and in the characteristics of the sets of bands they select. We also show the accuracy of the proposed methods with simulations using synthetic and real data. Finally, we present some preliminary conclusions.

6.1 REVISITING THE KERNEL FRAMEWORK

When employing kernelized methods for unmixing of hyperspectral images, standard mixing models are replaced by more flexible nonparametric or semi-parametric models. Thus, the ℓ -th band of a pixel observation can be modeled as

$$r_\ell = \psi(\mathbf{m}_{\lambda_\ell}) + n_\ell \quad (6.1)$$

with ψ a real-valued function in a RKHS \mathcal{H} that characterizes the nonlinear interactions between the endmembers, and n_ℓ an additive noise at the ℓ -th band. In order to estimate ψ in the least squares sense, it is possible to formulate a convex optimization problem as done in Section 3.1.3 and Appendix A.6 where the underlying function ψ can be written as

$$\psi = \sum_{j=1}^L \beta_j \kappa(\cdot, \mathbf{m}_{\lambda_j}) \quad (6.2)$$

where $\boldsymbol{\beta} = [\beta_1, \dots, \beta_L]$ is the vector of Lagrange multipliers, and can be found by solving the following linear system

$$(\mathbf{K} + \varepsilon \mathbf{I})\boldsymbol{\beta} = \mathbf{r} \quad (6.3)$$

where \mathbf{K} is the Gram matrix with entries $\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})$, with $i, j = 1, \dots, L$, and ε is the regularization parameter.

Although the formulation presented in Section 3.1.3 allows one to address an estimation problem in \mathcal{H} by solving the linear system in (6.3), this approach is computationally demanding since it involves the inversion of $L \times L$ matrices. Similarly, when considering the SK-Hype algorithm (Section 3.1.4), the linear system to be solved is given by

$$\left(\begin{array}{c|c} \mathbf{K}_u + \varepsilon \mathbf{I} & u\mathbf{M} \\ \hline u\mathbf{M}^\top & u\mathbf{I} \end{array} \right) \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} = \begin{pmatrix} \mathbf{r} \\ \mathbf{0} \end{pmatrix} \quad (6.4)$$

which requires a $(L+R) \times (L+R)$, $L \gg R$, matrix inverse. This issue is critical, as modern hyperspectral image sensors employ hundreds of contiguous bands with an ever increasing spatial resolution. Hence, it is of major interest to consider band selection techniques that lead to significant computational cost reduction without noticeable quality loss. Considering (6.2), a possible strategy is to focus on a reduced-order model of the form:

$$\psi = \sum_{j \in \mathcal{S}_D} \beta_j \kappa(\cdot, \mathbf{m}_{\lambda_j}) \quad (6.5)$$

where $\mathcal{S}_D \subset \{1, \dots, L\}$ is an M -element ($M < L$) subset of indexes. We shall call $\mathcal{D} = \{\kappa(\cdot, \mathbf{m}_{\lambda_j})\}_{j \in \mathcal{S}_D}$ the dictionary.

6.2 KERNEL K -MEANS FOR BAND SELECTION

Kernel k -means (KKM) is a direct extension of the k -means clustering algorithm [126]. It maps the input data $\mathbf{m}_{\lambda_\ell}$ into a RKHS \mathcal{H} , and groups their images $\kappa(\cdot, \mathbf{m}_{\lambda_\ell})$ into disjoint clusters $\mathcal{C}_1, \dots, \mathcal{C}_M$ based on their relative distance in \mathcal{H} . Since determining centroids in \mathcal{H} is intractable, KKM calculates distances using the reproducing property, see Definition (10) in the appendix Section A.3.2.

Given a cluster \mathcal{C}_k enclosing points $\{\kappa(\cdot, \mathbf{m}_{\lambda_\ell})\}_{\ell \in \mathcal{C}_k}$, its centroid is defined as

$$\mathbf{v}_k = \frac{1}{N_k} \sum_{i \in \mathcal{C}_k} \kappa(\cdot, \mathbf{m}_{\lambda_i}) \quad (6.6)$$

where N_k is the number of points in \mathcal{C}_k . The squared distance of any point $\kappa(\cdot, \mathbf{m}_{\lambda_\ell})$ to \mathbf{v}_k is computed as

$$\begin{aligned} \|\kappa(\cdot, \mathbf{m}_{\lambda_\ell}) - \mathbf{v}_k\|_{\mathcal{H}}^2 &= \kappa(\mathbf{m}_{\lambda_\ell}, \mathbf{m}_{\lambda_\ell}) \\ &\quad - \frac{1}{N_k} \sum_{i \in \mathcal{C}_k} \kappa(\mathbf{m}_{\lambda_\ell}, \mathbf{m}_{\lambda_i}) \\ &\quad + \frac{1}{N_k^2} \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} \kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j}) \end{aligned} \quad (6.7)$$

and the clustering error to minimize is defined as

$$E(\mathbf{v}_1, \dots, \mathbf{v}_M) = \sum_{k=1}^M \sum_{\ell \in \mathcal{C}_k} \|\kappa(\cdot, \mathbf{m}_{\lambda_\ell}) - \mathbf{v}_k\|_{\mathcal{H}}^2. \quad (6.8)$$

Each cluster \mathcal{C}_k is then represented by the band ℓ_k corresponding to the closest point to its centroid \mathbf{v}_k :

$$\ell_k = \arg \min_{\ell \in \mathcal{C}_k} \|\kappa(\cdot, \mathbf{m}_{\lambda_\ell}) - \mathbf{v}_k\|_{\mathcal{H}}^2. \quad (6.9)$$

The global kernel k -means (GKKM) algorithm uses the principles described above for incremental clustering [126]. GKKM does not suffer from poor convergence to local minima and produces near-optimal solutions that are robust to cluster initialization. A fast GKKM (FGKKM) version that performs

a unique KKM run and greatly reduces the complexity of the algorithm can also be used.

Algorithm 2 presents the pseudo code for the KMM BS algorithm. It receives as inputs the Gram matrix \mathbf{K} and the desired number of bands M in the final dictionary. In line 3, M clusters are found using the FGKKM. Then, for each cluster the vector $\mathbf{m}_{\lambda_\ell}$ minimizing (6.9) is found (line 6) and included in the index set $\mathcal{I}_{\mathcal{D}}$ in line 7. Finally, the algorithm returns the index set of selected bands $\mathcal{I}_{\mathcal{D}}$ in line 9.

Algorithm 2: FGKKM Band Selection (KKMBS)

Input : The $L \times L$ Gram matrix \mathbf{K} and the desired number of bands M .

Output: Selected band indexes $\mathcal{I}_{\mathcal{D}}$.

- 1 Initialization: $\mathcal{I}_{\mathcal{D}} = \{\emptyset\}$;
 - 2 % Find clusters indices
 - 3 $[\mathcal{C}_1, \dots, \mathcal{C}_M] = \text{FGKKM}(\mathbf{K}, M)$;
 - 4 % Find the vectors $\kappa(\cdot, \mathbf{m}_{\lambda_\ell})$ closest to the centroids in $\mathcal{C}_1, \dots, \mathcal{C}_M$
 - 5 **for** $k = 1$ **to** M **do**
 - 6 $l_k = \arg \min_{\ell \in \mathcal{C}_k} \|\kappa(\cdot, \mathbf{m}_{\lambda_\ell}) - \mathbf{v}_k\|_{\mathcal{H}}^2$;
 - 7 Insert l_k into $\mathcal{I}_{\mathcal{D}}$;
 - 8 **end**
 - 9 **return** $\mathcal{I}_{\mathcal{D}}$
-

6.3 COHERENCE-BASED BAND SELECTION

6.3.1 Coherence criterion for dictionary selection

Coherence is a parameter of fundamental interest for characterizing dictionaries of atoms in linear sparse approximation problems [127]. It was first introduced as a heuristic quantity for Matching Pursuit in [128]. Formal studies followed in [129], and were enriched for Basis Pursuit in [130, 131].

Consider a set of kernel functions $\{\kappa(\cdot, \mathbf{m}_{\lambda_\ell})\}_{\ell=1, \dots, M}$ in \mathcal{H} . The defi-

inition of coherence was extended to RKHS as [125]:

$$\begin{aligned}\mu &= \max_{i \neq j} |\langle \kappa(\cdot, \mathbf{m}_{\lambda_i}), \kappa(\cdot, \mathbf{m}_{\lambda_j}) \rangle_{\mathcal{H}}| \\ &= \max_{i \neq j} |\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})|\end{aligned}\tag{6.10}$$

where κ is a unit-norm kernel. Otherwise, replace $\kappa(\cdot, \mathbf{m}_{\lambda_i})$ with

$$\kappa(\cdot, \mathbf{m}_{\lambda_i}) / \sqrt{\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_i})}$$

in (6.10). Parameter μ is the largest absolute value of the off-diagonal entries in the Gram matrix. It reflects the largest cross correlation in the dictionary $\{\kappa(\cdot, \mathbf{m}_{\lambda_\ell})\}_\ell$, and is equal to zero for every orthonormal basis. A dictionary is said to be incoherent when its coherence μ is small. Although its definition is rather simple, coherence possesses important properties [125]. In particular, it can be shown that the kernel functions in the dictionary $\mathcal{D} = \{\kappa(\cdot, \mathbf{m}_{\lambda_\ell})\}_{\ell=1, \dots, M}$ are linearly independent if $(M - 1)\mu < 1$. This sufficient condition illustrates that the coherence (6.10) provides valuable information on a dictionary at low computational cost. Other properties are discussed in [125].

Kernel-based dictionary learning methods usually consider approximate linear dependence conditions to evaluate whether a candidate kernel function $\kappa(\cdot, \mathbf{m}_{\lambda_i})$ can be reasonably well represented by a combination of the kernel functions that are already in the dictionary \mathcal{D} . To avoid excessive computational complexity, a greedy dictionary learning method has been introduced in [125]. It consists of inserting the candidate $\kappa(\cdot, \mathbf{m}_{\lambda_i})$ into the dictionary \mathcal{D} provided its coherence is still below a given threshold μ_0 , namely,

$$\max_{j \in \mathcal{D}} |\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})| \leq \mu_0\tag{6.11}$$

where μ_0 is a parameter $[0, 1[$ determining both the maximum coherence in \mathcal{D} and its cardinality $|\mathcal{D}|$. Using coherence criterion for BS allows to explicitly limit the correlation of kernel functions in the dictionary. This contrasts with the kernel k -means strategy, which starts from a number of dictionary elements prescribed by the user without taking the coherence of kernel functions into consideration.

The coherence criterion (6.11) was proposed within the context of parameter estimation from streaming data. The design of the dictionary follows a greedy strategy. The first kernel function is selected arbitrarily, and each new candidate kernel function is tested using (6.11) to determine if it deserves being included in the dictionary. This procedure is appropriate for online ap-

plications because of its minimal computational cost. However, alternatives should be sought which may lead to more effective solutions in batch mode applications.

6.3.2 Band selection as a maximum clique problem

Consider a set of kernel functions $\{\kappa(\cdot, \mathbf{m}_{\lambda_\ell})\}_{\ell=1, \dots, L}$. Determining a subset \mathcal{D} with a prescribed coherence level can be viewed as a two-step procedure. The first step aims at listing all the pairs of functions that satisfy the coherence rule (6.11). This can be performed by constructing a $L \times L$ binary matrix \mathbf{B} with entries defined as:

$$\mathbf{B}_{ij} = \begin{cases} 1 & \text{if } |\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})| \leq \mu_0 \\ 0 & \text{otherwise.} \end{cases} \quad (6.12)$$

The second step consists of finding in \mathbf{B} , up to a simultaneous reordering of its rows and columns, the largest submatrix of only ones. This problem can be recast as determining a maximum clique in an undirected graph $\mathcal{G} = \{V, E\}$, where each vertex ℓ of $V = \{1, \dots, L\}$ corresponds to a candidate function $\kappa(\cdot, \mathbf{m}_{\lambda_\ell})$, and edges in $E \subseteq V \times V$ connecting the vertices are defined by the adjacency matrix \mathbf{B} . Two vertices are said to be adjacent if they are connected by an edge. A complete subgraph of \mathcal{G} is one whose vertices are pairwise adjacent. The maximal clique problem (MCP) consists of finding the maximal complete subgraph of \mathcal{G} [32]. This problem is NP-Complete [132]. Figure 27 illustrates this problem within the context of BS. This figure shows for instance that the coherence of $\kappa(\cdot, \mathbf{m}_{\lambda_1})$ and $\kappa(\cdot, \mathbf{m}_{\lambda_4})$ is lower than the preset threshold μ_0 , and the coherence of $\kappa(\cdot, \mathbf{m}_{\lambda_1})$ and $\kappa(\cdot, \mathbf{m}_{\lambda_2})$ is larger than μ_0 . This graph has one maximum clique defined by the set of vertices $\mathcal{S}_{\mathcal{D}} = \{1, 3, 4, 5\}$, which means that the coherence of the dictionary $\mathcal{D} = \{\kappa(\cdot, \mathbf{m}_{\lambda_j})\}_{j \in \mathcal{S}_{\mathcal{D}}}$ is lower than μ_0 and it has maximum cardinality. A vast literature exists on maximum clique problems (MCP), see [33] and references therein. The next section reviews the main algorithms for MCP.

6.3.3 The maximum clique problem

MCP has a wide range of practical applications arising in a number of domains such as bioinformatics, coding theory, economics, social network analysis, etc. Given its theoretical importance and practical interests, considerable efforts have been devoted for deriving exact and heuristic algorithms.

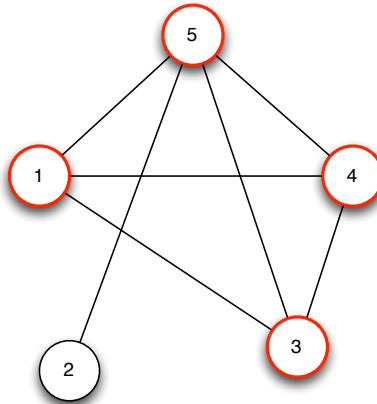


Figure 27: The maximum clique problem (MCP).

Efficient exact methods have been designed mainly based on the branch-and-bound (B&B) framework. Dynamic bounds on the clique size are used to prune (or discard) branches during search, and then dramatically reducing the search space [133]. Although algorithms are now much faster and efficient than their past counterparts [134], the inherent complexity of exact methods can still lead to a prohibitive computation time when large problems are addressed [33]. To handle problems whose optimal solutions cannot be reached within a reasonable time, various heuristic and metaheuristic algorithms have been derived with the purpose of providing sub-optimal solutions in an acceptable time. In this thesis, however, we shall focus on exact algorithms since our application concerns small graphs with a number of vertices equal to the number of bands.

Since the introduction of the Carraghan and Pardalos (CP) exact algorithm [133], many refinements have been proposed to improve its performance with a focus on two main issues. The first one is to tighten the upper bound on the maximum clique during search for the purpose of more efficient subtree pruning. The second one is to improve the branching rule, and then select the most promising vertices to expand candidate cliques. In [33], the authors classify the exact MCP algorithms into four groups, depending on their strategies for pruning and branching. The first group solves sub-clique problems for each vertex with iterative deepening and pruning strategies. Examples are the CP algorithm [133] and its improved version [135]. Both algorithms are sensitive to the order of vertices, which can result in drastically different execution times for a given graph [135]. A second group is based

on vertex coloring techniques [136]. The most prominent algorithms in this group use B&B strategies based on subgraph coloring. Examples of algorithms are BT and the recent MCQ, MCR, MaxCliqueDyn, BB-MaxClique, among others [33]. The third group improves the basic CP by tightening candidate sets via the removal of vertices that cannot be used to extend the current clique to a maximum clique. Along this line, three B&B algorithms, denoted DF, χ and χ +DF were proposed in [137]. The fourth group consists of the exact methods based on MaxSAT [124], which improve the techniques based on vertex coloring. The MaxCLQ algorithm proposed in [124] is considered to be very effective and solved the DIMACS problem (p_hat1000-3) for the first time [33]. A complex approach (ILS&MaxCLQ) that combines different algorithms such as the MaxCLQ, MCS and the ILS, was recently proposed [138]. A comparative discussion on exact methods is presented in [33]. The MaxCLQ and ILS&MaxCLQ were the only methods to solve all the presented problems, with the smallest CPU times for the former.

6.3.4 Coherence-based BS algorithms

We shall now introduce kernel BS algorithms based on the coherence criterion. As a baseline for performance comparisons, we consider first a greedy strategy that consists of testing candidate kernel functions sequentially and inserting them into the dictionary if coherence stays below a threshold value μ_0 . Next, we propose an exact strategy based on MCP solving.

6.3.4.1 Automatic parameter settings

Before describing the kernel BS methods, we briefly present a procedure for automatic parameter setting. It allows to set the coherence threshold μ_0 and Gaussian kernel bandwidth σ^2 given a desired number of elements in the dictionary.

Let \mathbf{K}_σ be the $L \times L$ Gram matrix whose (i, j) -th entry is defined by $\kappa_\sigma(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j})$, where κ_σ denotes the Gaussian kernel (A.27) parameterized by the bandwidth σ^2 , *i.e.*,

$$\kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{m}_{\lambda_i} - \mathbf{m}_{\lambda_j}\|^2\right).$$

Let \mathcal{D} be an M -element dictionary with coherence μ and index set $\mathcal{I}_\mathcal{D}$. Then, as shown in [125], a sufficient condition for linear independence of the M

elements of \mathcal{D} is given by $(M-1)\mu < 1$. We write:

$$\mu < \frac{1}{(M-1)}. \quad (6.13)$$

The objective is to build a dictionary with (approximately) M linearly independent elements. We thus propose to set the coherence threshold μ_0 as:

$$\mu_0 = \frac{1}{(M-1)} \quad (6.14)$$

and adjust σ^2 to obtain a Gram matrix \mathbf{K}_σ whose entries are close to μ_0 in some sense. Indeed, on the one hand, if all the off-diagonal entries of \mathbf{K}_σ are smaller than μ_0 , then \mathcal{D} contains the L available elements. On the other hand, if all the off-diagonal entries of \mathbf{K}_σ are greater than μ_0 , then \mathcal{D} should be composed of only one element. Therefore, we propose to adjust σ^2 such that $\mathbb{E}\{(\mathbf{K}_{\sigma_{ij}})_{(i \neq j)}\} = \mu_0$, where $\mathbb{E}\{\cdot\}$ is the expected value and can be approximated as

$$\mathbb{E}\{(\mathbf{K}_{\sigma_{ij}})_{(i \neq j)}\} \approx \frac{2}{L^2 - L} \sum_{i=1}^{L-1} \sum_{j=i+1}^L \mathbf{K}_{\sigma_{ij}}. \quad (6.15)$$

Then, we set σ^2 as the solution of the following optimization problem:

$$\begin{aligned} \sigma^2 = \arg \min_{\sigma^2} & \left(\frac{2}{L^2 - L} \sum_{i=1}^{L-1} \sum_{j=i+1}^L [\mathbf{K}_{1_{ij}}]^{1/\sigma^2} - \mu_0 \right)^2 \\ \text{s. t. } & \sigma^2 \in \mathbb{R}^+. \end{aligned} \quad (6.16)$$

where $\mathbf{K}_1 = \mathbf{K}_\sigma$ is the Gram matrix for $\sigma = 1$. Finally, we determine $\mathbf{K}_\mathcal{D}$ as the largest sub-matrix of \mathbf{K}_σ whose all off-diagonal entries satisfy (6.11). We emphasize that since $\mathbf{K}_{\sigma_{ij}} \leq 1$, (6.15) is a decreasing function of σ^{-2} , and thus (6.16) has a unique solution.

6.3.4.2 Algorithms

In this section we present the two band selection algorithms using the greedy and clique approaches that will be used in Section 6.4.

The greedy coherence-based approach is presented in Algorithm 3. The inputs to Algorithm 3 are the desired number M of bands in the final dictionary, and the $L \times L$ Gaussian kernel Gram matrix with $\sigma = 1$ and en-

tries $\mathbf{K}_{1ij} = \kappa(\mathbf{m}_{\lambda_i}, \mathbf{m}_{\lambda_j}) = \exp(-0.5\|\mathbf{m}_{\lambda_i} - \mathbf{m}_{\lambda_j}\|^2)$. It returns the index of selected bands and the Gaussian kernel bandwidth σ^2 . Initialization occurs in line 1, where the index set $\mathcal{I}_{\mathcal{G}}$ is initialized with the first spectral band index, the number N_b of bands in the dictionary is set to one, and the coherence threshold μ_0 is adjusted according to (6.14). Next, σ^2 is determined by solving problem (6.16) in line 2, and the Gram matrix \mathbf{K}_{σ} is computed with the optimum σ^2 in line 3. From line 4 to line 13 the algorithm sequentially tests all the $L-1$ remaining bands using condition (6.11). Breaking the parts down, in line 5 a zero vector \mathbf{c} of length N_b is created, and the off diagonal terms $(\ell, \mathcal{I}_{\mathcal{G}_j})$ of the Gram matrix \mathbf{K}_{σ} are stored in \mathbf{c} . If the maximum absolute value of the entries of \mathbf{c} is less than the coherence threshold (line 9), then the ℓ -th band index is added to $\mathcal{I}_{\mathcal{G}}$, and N_b is incremented by one (lines 10 and 11). Finally, the algorithm returns the complete set of selected bands and the kernel bandwidth in line 14.

Algorithm 3: Greedy Coherence-based Band Selection (GCBS)

Input : The $L \times L$ Gram matrix $\mathbf{K}_1 = (\mathbf{K}_{\sigma})_{\sigma=1}$, and the desired number M of atoms.

Output: The indices $\mathcal{I}_{\mathcal{G}}$ of selected atoms, and the Gaussian kernel bandwidth σ^2 .

- 1 Initialization: $\mathcal{I}_{\mathcal{G}} = \{1\}$, $N_b = 1$, $\mu_0 = 1/(M-1)$;
- 2 Find σ^2 solving (6.16);
- 3 Compute \mathbf{K}_{σ} using σ^2 obtained in line 2;
- 4 **for** $\ell := 2$ **to** L **do**
- 5 $\mathbf{c} := \mathbf{0}_{N_b \times 1}$;
- 6 **for** $j := 1$ **to** N_b **do**
- 7 $\mathbf{c}_j := \mathbf{K}_{\sigma_{\ell, \mathcal{I}_{\mathcal{G}_j}}}$;
- 8 **end**
- 9 **if** $\max(|\mathbf{c}_j|) \leq \mu_0$ **then**
- 10 Insert ℓ into $\mathcal{I}_{\mathcal{G}}$;
- 11 $N_b := N_b + 1$;
- 12 **end**
- 13 **end**
- 14 **return** $\mathcal{I}_{\mathcal{G}}$, σ^2 ;

The clique coherence-based band selection method is described in Algorithm 4. Similarly to Algorithm 3, the inputs are \mathbf{K}_1 and M . The adjacency

matrix \mathbf{B} in initialized with zeros (line 1), the vertices vector V with the indices of all available wavelengths, μ_0 following (6.14), and $\mathcal{I}_{\mathcal{G}}$ as an empty set. The kernel bandwidth is computed in line 2, and the Gram matrix is computed for the optimum σ^2 in line 3. Through line 4 to 10 every entry of the upper diagonal part of \mathbf{B} is set according to (6.12). In line 11 the *MaxCLQ* algorithm is used to find the indices of the maximum clique in the graph. These indices are assigned to the dictionary index set $\mathcal{I}_{\mathcal{G}}$, which is returned in line 12 together with the kernel bandwidth.

Algorithm 4: Clique Coherence-based BS (CCBS)

Input : The $L \times L$ Gram matrix $\mathbf{K}_1 = (\mathbf{K}_{\sigma})_{\sigma=1}$, and the desired number M of atoms.

Output: The indices $\mathcal{I}_{\mathcal{G}}$ of selected atoms, and the Gaussian kernel bandwidth σ^2 .

- 1 Initialization: $\mathbf{B} := \mathbf{0}_{L \times L}$, $V = \{1, \dots, L\}$, $\mu_0 = 1/(M - 1)$,
 $\mathcal{I}_{\mathcal{G}_c} = \{\emptyset\}$;
 - 2 Find σ^2 solving (6.16);
 - 3 \mathbf{K}_{σ} using σ^2 obtained in line 2;
 - 4 **for** $i := 1$ **to** $L - 1$ **do**
 - 5 **for** $j := i + 1$ **to** L **do**
 - 6 **if** $[\mathbf{K}_{\sigma}]_{ij} \leq \mu_0$ **then**
 - 7 $\mathbf{B}_{ij} := 1$;
 - 8 **end**
 - 9 **end**
 - 10 **end**
 - 11 $\mathcal{I}_{\mathcal{G}} := \text{MaxCLQ}(V, \mathbf{B})$;
 - 12 **return** $\mathcal{I}_{\mathcal{G}}$, σ^2 ;
-

Note that M is used in Algorithm 4 and Algorithm 3 as a design parameter, which is required to obtain the coherence threshold and the Gaussian kernel bandwidth. The number N_b of bands in the final dictionary can differ from M .

6.4 SIMULATIONS

6.4.1 Simulation with synthetic data

This section presents simulation results using synthetic data to illustrate the performance of the proposed band selection methods under controlled conditions for which the abundance values are known. We constructed synthetic images using two sets of endmembers. The first set had 8 endmembers extracted from the spectral library of the ENVI software and correspond to the spectral signatures of minerals present in the Cuprite mining field in Nevada. The minerals are alunite, calcite, epidote, kaolinite, buddingtonite, almandine, jarosite and lepidolite, and their spectra consisted of 420 contiguous bands, covering wavelengths from $0.3951\mu\text{m}$ to $2.56\mu\text{m}$. The second set was extracted from the Pavia University data acquired by the ROSIS spectrometer. It has 610×340 pixels with 103 bands over the spectral range of 430–680 nm (Figure 28, left). The data also has a ground truth labelling 42776 pixels (out of the 207400) into 9 classes labeled asphalt, meadows, gravel, trees, painted metal sheets, bare soil, bitumen, self-blocking bricks and shadows (Figure 28, right). We extracted the endmembers from this data set using the vertex component analysis algorithm (VCA [23]), and considering only the labeled pixels. We constructed four 2000-pixel hyperspectral images ($N = 2000$), each using 8 endmembers ($R = 8$) from the Cuprite or Pavia data, and the GBM and PNMM (2.21) mixing models (see Section 2.2.2 and 2.2.3) with $\gamma_{i,i} = \gamma = 1$ and $\xi = 0.7$, respectively. The abundances were obtained by uniformly sampling from the simplex, *i.e.*, obeying the positivity and sum-to-one constraints. WGN was added to all images with power adjusted to produce a 21dB SNR. We consider the root mean square error (RMSE) in abundance estimation

$$\text{RMSE} = \sqrt{\frac{1}{NR} \sum_{n=1}^N \|\alpha_n - \alpha_n^*\|^2} \quad (6.17)$$

and the CPU time required for both BS (when applicable) and unmixing (averaged over 100 unmixings of the same HIs) to compare the different BS strategies. All unmixings were performed using a Gaussian kernel and considering either the full set of bands or smaller sets selected using the BS strategies presented in this chapter. SK-Hype was implemented for the full set of bands. The kernel bandwidth for SK-Hype was selected among the values $\sigma_{\text{skp}} \in \{0.5\sigma, \sigma, 2\sigma, 10\sigma, 20\sigma\}$ to obtain the minimum RMSE, where σ is the solution of (6.16), for $M = 30$. The global kernel k-means band selec-

tion (KKMBS) algorithm (Section 6.2) implementation requires the number of bands to be fixed *a priori*. To circumvent this issue, we considered a selection approach based on the Akaike Information Criterion and given by [139]

$$M = \arg \min_M [E(v_1, \dots, v_M) + \lambda M] \quad (6.18)$$

where the parameter λ controls the complexity of the model, and needs to be found empirically. The kernel bandwidth σ_{kkm} also needs to be selected for KKMBS. A grid search was performed using a small part (200 pixels) of the synthetic image to find λ and σ_{kkm} that would lead to a good RMSE performance. The parameters were chosen among the values $\lambda \in \{2, 4, 6\}$ and $\sigma_{\text{kkm}} \in \{0.5\sigma, \sigma, 2\sigma, 10\sigma, 20\sigma\}$, again with σ being the solution of (6.16), for $M = 30$. The parameter set leading to the best performance in terms of RMSE for the abundances was then selected. It is important to notice that, in general, the abundance ground truth is not available from real data. Thus, the RMSE in abundance estimation could not be used in design as a measure to select model parameters. Hence, the SK-Hype and KKMBS designs used in this comparison are based on a quasi-optimal choice of parameters for these methods, which could not be determined in practice. The proposed design for the BS methods, however, can be employed in practical applications.

BS with the CCBS and GCBS algorithms was performed using $M \in \{5, 10, 20, 30\}$, with parameters μ_0 and σ adjusted using the methodology presented in Section 6.3.4.1. We emphasize that this parameter setting strategy assumes no *prior* knowledge about the abundance ground truth.

The simulation results are summarized in Tables 13 to 16. In these tables, the first column shows the BS strategy considered prior to unmixing. SK-Hype in this column indicates the solution without BS. The symbol “(r)” besides CCBS or GCBS means that we have randomized the order of the bands prior to applying the BS strategy. The second column shows the obtained RMSE and the standard deviation (STD) in abundance estimation. The third column lists the average CPU time elapsed in the (BS + unmixing) process. Column four shows the number of selected bands N_b , and last column shows the coherence of the final dictionary.

Tables 13 and 14 show the results for HIs built with Cuprite endmembers and using, respectively, the GBM and the PNMM mixing models. Note that the RMSE obtained using the BS algorithms are very close to those obtained using all bands. Nevertheless the reduction in number of bands obtained through BS is at least tenfold. The computational complexity advantage of the BS methods is evidenced by the required average CPU time, which show reductions by factors ranging from 50 to 110, depending on the algorithm and parameter settings. Note also that the number of bands in the

final dictionary tends to be larger than the value M used to initialize the algorithms. This increase in the anticipated number of bands is obtained to optimize the dictionary coherence, what is not possible in the KKMBS algorithm. As expected, the number of bands remained the same for the clique algorithm (CCBS) for each value of M , and the slight changes in the RMSE results indicate that the maximum clique is not unique. For the greedy approach (GCBS), however, different numbers of bands are obtained at each execution due to initial randomization, and the results in terms of RMSE and CPU time vary slightly. In general, randomization did not have any significant impact on the results. Finally, one should note from these tables that the coherence-based algorithms produced dictionaries with coherence close to μ_0 , and 2 to 23 times smaller than the coherence obtained using KKMBS.

Tables 15 and 16 show the results for the HIs created with the Pavia endmembers using the GBM and PNMM respectively. Although the results in Tables 15 and 16 follows the same pattern that the results in Tables 13 and 14, we highlight that for the Pavia HIs the number of available bands is 103 in contrast to the 420 used in the previous example. This explains the smaller improvement in the Av. Time (Average Time) when using the BS algorithms which is about 3 to 4 times smaller than using all the bands. Another difference in the results is that using the BS algorithms, and its reasoning for setting μ_0 and σ^2 , the best results in terms of RMSE were obtained by the proposed method CCBS with $M = 30$ in both Tables. When concerning the number of bands, the final N_b were closer to M than in the previous example. For the coherence of the final dictionary the same pattern obtained in Tables 13 and 14 repeats for the Pavia HIs.

6.4.2 Simulation with real data

When working with real data ground truth for the fractional abundances are rarely available. Thus, we compare the abundance estimation results obtained using a full band approach and using the proposed band selection strategy. First, the data is unmixed using the SK-Hype algorithm using all the available spectral bands, what yields the estimated abundances $\alpha_n^{\text{skp}}, n = 1, \dots, N$. The unmixing is then done for all of the BS methods presented in this chapter. Generically denominating the BS-based estimated abundances as $\alpha_n^{\text{bs}}, n = 1, \dots, N$, the RMSE between the SK-Hype abundances and those obtained using a given BS algorithm is computed as

$$\text{RMSE} = \sqrt{\frac{1}{NR} \sum_{n=1}^N \|\alpha_n^{\text{skp}} - \alpha_n^{\text{bs}}\|^2}. \quad (6.19)$$

Table 13: RMSE. 100 runs, 2000 pxl., 8 endmembers (Cuprite), SNR=21dB, GBM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).

Strategy	RMSE \pm STD	Av. Time	N_b	μ
SK-Hype	0.0680 \pm 0.0028	301.08 \pm 17.93	420	-
KKMBS	0.0664 \pm 0.0026	25.40 \pm 0.22	36	0.5893
$M = 5, \mu_0 = 0.25, \sigma = 0.2548$				
CCBS	0.0687 \pm 0.0028	3.10 \pm 0.14	10	0.2482
CCBS (r)	0.0687 \pm 0.0028	3.13 \pm 0.12	10	0.2482
GCBS	0.0724 \pm 0.0031	2.91 \pm 0.02	8	0.2482
GCBS (r)	0.0721 \pm 0.0030	3.15 \pm 0.15	7.13 \pm 0.97	0.2331
$M = 10, \mu_0 = 0.1111, \sigma = 0.1320$				
CCBS	0.0678 \pm 0.0027	2.85 \pm 0.13	16	0.1108
CCBS (r)	0.0679 \pm 0.0027	2.89 \pm 0.17	16	0.1108
GCBS	0.0685 \pm 0.0028	2.57 \pm 0.02	16	0.1104
GCBS (r)	0.0688 \pm 0.0028	2.65 \pm 0.06	13.09 \pm 1.10	0.0996
$M = 20, \mu_0 = 0.0526, \sigma = 0.0965$				
CCBS	0.0659 \pm 0.0026	2.96 \pm 0.15	21	0.0520
CCBS (r)	0.0660 \pm 0.0026	3.01 \pm 0.17	21	0.0520
GCBS	0.0670 \pm 0.0027	2.59 \pm 0.02	20	0.0525
GCBS (r)	0.0678 \pm 0.0027	2.67 \pm 0.08	15.95 \pm 1.13	0.0467
$M = 30, \mu_0 = 0.0345, \sigma = 0.0503$				
CCBS	0.0637 \pm 0.0024	5.54 \pm 0.22	42	0.0339
CCBS (r)	0.0637 \pm 0.0024	5.74 \pm 0.18	42	0.0339
GCBS	0.0637 \pm 0.0024	3.32 \pm 0.04	41	0.0344
GCBS (r)	0.0644 \pm 0.0025	2.83 \pm 0.07	33.39 \pm 1.43	0.0326

Table 14: RMSE. 100 runs, 2000 pxl., 8 endmembers (Cuprite), SNR=21dB, PNMM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).

Strategy	RMSE \pm STD	Av. Time	N_b	μ
SK-Hype	0.0728 ± 0.0030	277.03 ± 4.30	420	-
KKMBS	0.0729 ± 0.0030	25.52 ± 0.18	36	0.7760
$M = 5, \mu_0 = 0.25, \sigma = 0.2548$				
CCBS	0.0748 ± 0.0031	2.99 ± 0.10	10	0.2482
CCBS (r)	0.0749 ± 0.0031	3.12 ± 0.18	10	0.2482
GCBS	0.0764 ± 0.0032	2.85 ± 0.06	8	0.2482
GCBS (r)	0.0776 ± 0.0033	2.99 ± 0.15	7.13 ± 0.97	0.2331
$M = 10, \mu_0 = 0.1111, \sigma = 0.1320$				
CCBS	0.0746 ± 0.0031	2.85 ± 0.19	16	0.1108
CCBS (r)	0.0745 ± 0.0031	2.84 ± 0.14	16	0.1108
GCBS	0.0757 ± 0.0032	2.57 ± 0.04	16	0.1104
GCBS (r)	0.0757 ± 0.0031	2.64 ± 0.10	13.09 ± 1.10	0.0996
$M = 20, \mu_0 = 0.0526, \sigma = 0.0965$				
CCBS	0.0735 ± 0.0029	2.87 ± 0.12	21	0.0520
CCBS (r)	0.0737 ± 0.0029	2.96 ± 0.17	21	0.0520
GCBS	0.0753 ± 0.0031	2.55 ± 0.03	20	0.0525
GCBS (r)	0.0753 ± 0.0031	2.56 ± 0.04	15.95 ± 1.13	0.0467
$M = 30, \mu_0 = 0.0345, \sigma = 0.0503$				
CCBS	0.0740 ± 0.0029	5.41 ± 0.18	42	0.0339
CCBS (r)	0.0740 ± 0.0029	5.62 ± 0.19	42	0.0339
GCBS	0.0737 ± 0.0029	3.24 ± 0.04	41	0.0344
GCBS (r)	0.0742 ± 0.0030	2.74 ± 0.07	33.39 ± 1.43	0.0326

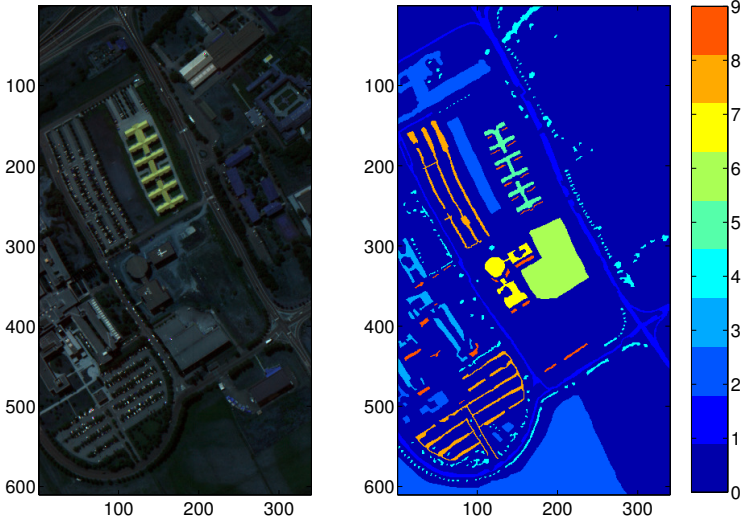


Figure 28: Pavia University. In the left, the Pavia University HI is represented using the bands 5, 30, and 50. In the right, the classified areas are labelled from 1 to 9, while 0 corresponds to unclassified areas.

The images used are shown in Figure 28 and Figure 29. The first image is the scene from the Pavia University described in Section 6.4.1. It has 207400 pixels and the endmembers were also extracted using VCA, see Section 6.4.1. The second image is a scene from the Cuprite mining field site in Nevada, acquired by the AVIRIS instrument. It has originally 224 spectral bands, from which we have removed the water absorption bands, resulting in 188 bands. This scene has 7371 pixels and previous analysis identified five minerals (Sphene, Montmorillonite, Kaolinite, Dumortierite, and Pyrope) to have strong components in this particular region [3]. The endmember matrix was extracted using the VCA algorithm [23].

Tables 17 and 18 show the abundance RMSE results obtained using (6.19). For both tables, the RMSE performance is compatible to that obtained using synthetic images, and the savings in computational complexity can be inferred from the CPU time reduction by a factor of at least 13 (for $M = 30$) for the Cuprite scene and at least 3 (for $M = 30$) for the Pavia scene. In comparing CCBS and GCBS with KKMBS one should note the significant reduction obtained in dictionary coherence for the same model complexity (N_b).

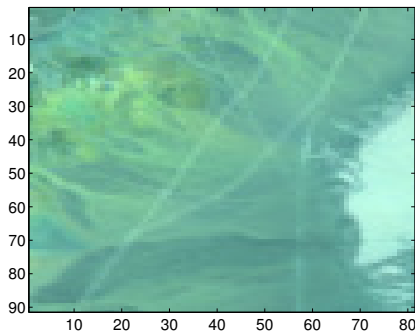


Figure 29: Cuprite scene used in [3].

6.5 PRELIMINARY CONCLUSIONS

In this chapter we have proposed two methods for nonlinear unmixing of hyperspectral images, which employ band selection directly in the reproducing kernel Hilbert space (RKHS). The first method employs the KKM algorithm to find clusters in the RKHS where each cluster centroids are associated to the closest mapped spectral vector. The second method is centralized and based on the coherence criterion, which incorporates a measure of the quality of the dictionary in the RKHS for the nonlinear unmixing. We have shown that this BS approach is equivalent to solving a maximum clique problem (MCP). Contrary to competing methods that do not include an efficient choice of the model parameters, the CCBS requires only an initial guess on the number of selected bands. Simulation results employing both synthetic and real data illustrate the quality of the unmixing results obtained with the proposed methods, which leads to abundance estimations as accurate as those obtained using the full-band SK-Hype method, at a small fraction of the computational cost.

Table 15: RMSE. 100 runs, 2000 pxl., 8 endmembers (Pavia), SNR=21dB, GBM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).

Strategy	RMSE \pm STD	Av. Time	N_b	μ
SK-Hype	0.0810 \pm 0.0035	15.2468 \pm 0.3231	103	-
KKMBS	0.0852 \pm 0.0038	5.69 \pm 0.01	5	0.5347
$M = 5, \mu_0 = 0.25, \sigma = 0.2385$				
CCBS	0.0845 \pm 0.0037	4.62 \pm 0.05	6	0.2402
CCBS (r)	0.0845 \pm 0.0037	4.64 \pm 0.05	6	0.2395
GCBS	0.0848 \pm 0.0037	4.54 \pm 0.02	6	0.2338
GCBS (r)	0.0862 \pm 0.0038	5.02 \pm 0.21	4.89 \pm 0.37	0.1812
$M = 10, \mu_0 = 0.1111, \sigma = 0.1$				
CCBS	0.0813 \pm 0.0035	3.51 \pm 0.04	12	0.1098
CCBS (r)	0.0813 \pm 0.0035	3.53 \pm 0.05	12	0.1098
GCBS	0.0824 \pm 0.0035	3.65 \pm 0.03	12	0.1080
GCBS (r)	0.0832 \pm 0.0036	3.76 \pm 0.12	9.58 \pm 0.75	0.0907
$M = 20, \mu_0 = 0.0526, \sigma = 0.0498$				
CCBS	0.0795 \pm 0.0034	3.43 \pm 0.04	20	0.0383
CCBS (r)	0.0794 \pm 0.0034	3.45 \pm 0.04	20	0.0437
GCBS	0.0795 \pm 0.0034	3.49 \pm 0.02	20	0.0499
GCBS (r)	0.0804 \pm 0.0035	3.45 \pm 0.07	16.55 \pm 0.88	0.0408
$M = 30, \mu_0 = 0.0345, \sigma = 0.0353$				
CCBS	0.0784 \pm 0.0034	3.68 \pm 0.03	25	0.0314
CCBS (r)	0.0784 \pm 0.0033	3.68 \pm 0.04	25	0.0311
GCBS	0.0787 \pm 0.0034	3.67 \pm 0.03	25	0.0300
GCBS (r)	0.0790 \pm 0.0034	3.54 \pm 0.06	21.09 \pm 1.02	0.0282

Table 16: RMSE. 100 runs, 2000 pxl., 8 endmembers (Pavia), SNR=21dB, PNMM, SK-Hype. μ_0 computed using Equation (6.14) for a given M , and σ is found solving problem (6.16).

Strategy	RMSE \pm STD	Av. Time	N_b	μ
SK-Hype	0.0839 ± 0.0035	14.6747 ± 0.3073	103	-
KKMBS	0.0878 ± 0.0038	5.31 ± 0.02	5	0.5347
$M = 5, \mu_0 = 0.25, \sigma = 0.2385$				
CCBS	0.0861 ± 0.0037	4.34 ± 0.04	6	0.2402
CCBS (r)	0.0861 ± 0.0037	4.34 ± 0.05	6	0.2395
GCBS	0.0877 ± 0.0038	4.17 ± 0.02	6	0.2338
GCBS (r)	0.0882 ± 0.0039	4.56 ± 0.23	4.89 ± 0.37	0.1812
$M = 10, \mu_0 = 0.1111, \sigma = 0.1$				
CCBS	0.0835 ± 0.0035	3.27 ± 0.03	12	0.1098
CCBS (r)	0.0835 ± 0.0035	3.25 ± 0.04	12	0.1098
GCBS	0.0852 ± 0.0035	3.32 ± 0.01	12	0.1080
GCBS (r)	0.0857 ± 0.0036	3.38 ± 0.08	9.58 ± 0.75	0.0907
$M = 20, \mu_0 = 0.0526, \sigma = 0.0498$				
CCBS	0.0817 ± 0.0034	3.22 ± 0.04	20	0.0383
CCBS (r)	0.0817 ± 0.0034	3.23 ± 0.05	20	0.0437
GCBS	0.0817 ± 0.0034	3.27 ± 0.02	20	0.0499
GCBS (r)	0.0828 ± 0.0035	3.24 ± 0.05	16.55 ± 0.88	0.0408
$M = 30, \mu_0 = 0.0345, \sigma = 0.0353$				
CCBS	0.0804 ± 0.0033	3.43 ± 0.05	25	0.0314
CCBS (r)	0.0803 ± 0.0033	3.45 ± 0.03	25	0.0311
GCBS	0.0806 ± 0.0033	3.48 ± 0.05	25	0.0300
GCBS (r)	0.0810 ± 0.0034	3.33 ± 0.06	21.09 ± 1.02	0.0282

Table 17: Cuprite image. RMSE between the abundances estimated with SK-Hype (all bands) and BS + SK-Hype.

Strategy	RMSE \pm STD	CPU Time	N_b	μ
SK-Hype	-	282.42	188	-
KKMBS	0.0777 ± 0.0036	19.289	13	0.8162
$M = 5, \mu_0 = 0.25, \sigma = 0.0963$				
CCBS	0.0805 ± 0.0038	18.4835	9	0.2495
GCBS	0.0833 ± 0.0040	17.7114	9	0.2483
$M = 10, \mu_0 = 0.1111, \sigma = 0.0489$				
CCBS	0.0659 ± 0.0027	15.2023	16	0.1090
GCBS	0.0695 ± 0.0029	14.5721	15	0.1090
$M = 20, \mu_0 = 0.0526, \sigma = 0.0260$				
CCBS	0.0477 ± 0.0015	17.0942	25	0.0471
GCBS	0.0484 ± 0.0015	16.9595	25	0.0493
$M = 30, \mu_0 = 0.0345, \sigma = 0.0178$				
CCBS	0.0378 ± 0.0010	20.6932	35	0.0333
GCBS	0.0395 ± 0.0011	20.4790	34	0.0300

Table 18: Pavia University image. RMSE between the abundances estimated with SK-Hype (all bands) and BS + SK-Hype.

Strategy	RMSE \pm STD	CPU Time	N_b	μ
SK-Hype	-	1740.47	103	-
KKMBS	0.0446 ± 0.0015	568.10	13	0.5066
$M = 5, \mu_0 = 0.25, \sigma = 0.2492$				
CCBS	0.0659 ± 0.0037	513.21	6	0.2499
GCBS	0.0650 ± 0.0036	533.48	6	0.2499
$M = 10, \mu_0 = 0.1111, \sigma = 0.1017$				
CCBS	0.0435 ± 0.0016	495.13	12	0.1024
GCBS	0.0500 ± 0.0023	497.92	12	0.1019
$M = 20, \mu_0 = 0.0526, \sigma = 0.0503$				
CCBS	0.0301 ± 0.0008	488.67	21	0.0433
GCBS	0.0309 ± 0.0009	488.66	21	0.0472
$M = 30, \mu_0 = 0.0345, \sigma = 0.0336$				
CCBS	0.0260 ± 0.0007	535.64	26	0.0336
GCBS	0.0263 ± 0.0007	538.63	26	0.0336

7 CONCLUSIONS

In this thesis we tackled different issues within the complete unmixing process of nonlinearly mixed hyperspectral images. We contributed to the detection of nonlinearly mixed pixels, endmember estimation when a considerable part of the HI is nonlinearly mixed, and band selection in the RKHS.

In Chapter 4 we presented a nonparametric method for detecting nonlinear mixtures in hyperspectral images. The performance of the detector was studied for supervised and unsupervised unmixing problems. Furthermore, we showed that the improvement in the unmixing performance obtained when using the proposed detector is statistically consistent. Additionally, a degree of mixture nonlinearity based on the relative energies of the linear and nonlinear contributions to the mixing process was defined to quantify the importance of the linear and nonlinear model counterparts. Such definition is important for a proper evaluation of the relative performances of different nonlinear mixture detection strategies.

In Chapter 5 an iterative algorithm was derived for endmember estimation as a pre-processing step for unsupervised unmixing problems. It was shown that the combined use of the detector presented in Chapter 4 and endmember estimation algorithm leads to better unmixing results when compared to state-of-the-art solutions. Simulations using different scenarios corroborate the conclusions.

In Chapter 6 we have proposed two methods for nonlinear unmixing of hyperspectral images, which employ band selection directly in the reproducing kernel Hilbert space (RKHS). The first method employs the kernel k-means (KKM) algorithm to find clusters in the RKHS where each cluster centroids are associated to the closest mapped spectral vector. The second method is centralized and based on the coherence criterion, which incorporates a measure of the quality of the dictionary in the RKHS for the nonlinear unmixing. We have shown that this BS approach is equivalent to solving a maximum clique problem (MCP). Contrary to competing methods that do not include an efficient choice of the model parameters, the CCBS requires only an initial guess on the number of selected bands. Simulation results employing both synthetic and real data illustrate the quality of the unmixing results obtained with the proposed methods, which leads to abundance estimations as accurate as those obtained using the full-band SK-Hype method, at a small fraction of the computational cost.

7.1 FUTURE WORK

This work concentrated in kernel methods to solve the unmixing problem in HIs. Thus, different possibilities were explored relating kernel methods and spectral unmixing itself. In this context, we highlight four problems that naturally arise as possible continuations of the research work in the thesis. They are:

- analyze the impact of the band selection strategy discussed in Chapter 6 on the performance of the nonlinearity detector discussed in Chapter 4;
- investigate the possibilities of using different detection strategies applied to the detection of nonlinearly mixed pixels detection;
- follow the track opened in Section 6.3.4.1 and try to better understand the relation between the kernel parameters, coherence, the accuracy of the abundance estimation, and the universal property of the Gaussian kernel to improve the methodology for designing the system (including kernel's) parameters;
- consider total least-squares, as opposed to least-squares, kernelized approach. This makes sense since spectral unmixing uncertainties are present in both endmembers and pixel observations. Thus, the kernel-TLS is a natural choice to treat the problem when nonlinearity is present.

7.2 PUBLICATIONS

During the period of this work we produced the following papers:

- T. Imbiriba, J.C.M. Bermudez, J.-Y. Tournet, and C. Richard. Detection of nonlinear mixtures using Gaussian processes: application to hyperspectral imaging. In *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 7949-7953, May 2014.
- T. Imbiriba, J.C.M. Bermudez, C. Richard, J.-Y. Tournet. Nonparametric detection of nonlinearly mixed pixels and endmember estimation in hyperspectral images. *IEEE Transactions on Image Processing*, v. 25, n. 3, p. 1136–1151, March 2016. ISSN 1057-7149.
- T. Imbiriba, J.C.M. Bermudez, C. Richard, J.-Y. Tournet. Band selection in RKHS for fast nonlinear unmixing of hyperspectral images. In: 2015 23rd European Signal Processing Conference (EUSIPCO). 2015. p. 1651–1655.

- T. Imbiriba, J.C.M. Bermudez, C. Richard. Band selection for nonlinear unmixing of hyperspectral images as a maximal clique problem. *IEEE Transactions on Image Processing* (Submitted). March 2016.

7.3 SOURCE CODE

Matlab source codes and datasets which replicate the simulations presented in this work are available at https://github.com/talesim/NP_NL_Det_EE_HI/archive/master.zip (Chapters 4 and 5) and https://github.com/talesim/cliqeu_BS/archive/master.zip (Chapter 6).

APPENDIX A – Convex Optimization in RKHS

A.1 CONVEX FUNCTIONS

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is *convex* if $\mathbf{dom} f$ is a convex set, and if for all $\mathbf{x}, \mathbf{y} \in \mathbf{dom} f$, and $\theta \in [0, 1]$, we have

$$f(\theta \mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}). \quad (\text{A.1})$$

A function f is *strictly convex* if strict inequality holds in (A.1) whenever $\mathbf{x} \neq \mathbf{y}$ and $0 < \theta < 1$. We say f is *concave* if $-f$ is convex, and *strictly concave* if $-f$ is strictly convex.

For an affine function we always have equality in (A.1), so all affine (and therefore also linear) functions are both convex and concave. Conversely, any function that is convex and concave is affine [140, pg. 67].

A.2 THE LAGRANGE DUAL PROBLEM

Consider the following constrained optimization problem

$$\begin{aligned} & \text{minimize} && f_0(\mathbf{x}) \\ & \text{subject to} && f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & && h_i(\mathbf{x}) = 0, \quad i = 1, \dots, p. \end{aligned} \quad (\text{A.2})$$

with $\mathbf{x} \in \mathbb{R}^n$. We assume its domain $\mathcal{D} = (\bigcap_{i=1}^m \mathbf{dom} f_i) \cap (\bigcap_{i=1}^p \mathbf{dom} h_i)$ is nonempty, and denote the optimal value of (A.2) by $p^* = f_0(\mathbf{x}^*)$.

It is possible to formulate a (dual) problem for which the optimal solution d^* is a lower bound for the problem (A.2), i.e. $d^* \leq p^*$. This dual problem is formulated considering the Lagrangean $L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ associated to problem (A.2), which is given by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p v_i h_i(\mathbf{x}), \quad (\text{A.3})$$

where $\lambda_i, i = 1, \dots, m$ and $v_i, i = 1, \dots, p$, are Lagrange multipliers. The solution of the dual problem yields to lower bounds on the optimal value p^* if we consider $\boldsymbol{\lambda} \geq 0$ and any \mathbf{v} . To better understand this first note that for $\boldsymbol{\lambda} \geq 0$ we have $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v}) \leq f_0(\mathbf{x})$. Thus, the Lagrange dual problem can finally be formulated as

$$\begin{aligned} & \text{maximize} && g(\boldsymbol{\lambda}, \mathbf{v}) \\ & \text{subject to} && \boldsymbol{\lambda} \geq 0, \end{aligned} \quad (\text{A.4})$$

for which the dual function is given by $g(\boldsymbol{\lambda}, \mathbf{v}) = \inf_{\mathbf{x} \in \mathcal{D}} (L(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v}))$ and has optimal value $d^* = g(\boldsymbol{\lambda}^*, \mathbf{v}^*) \leq p^*$.

Note that the problem (A.4) is concave whether the problem (A.2) is convex or not. This is true because $g(\boldsymbol{\lambda}, \mathbf{v})$ is the pointwise infimum of a family of affine functions of $(\boldsymbol{\lambda}, \mathbf{v})$ [140, pg. 216]. Also note, that $d^* \leq p^*$ holds since $f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}^*) + \sum_{i=1}^p v_i h_i(\mathbf{x}^*) \leq f_0(\mathbf{x}^*)$ considering the equality constrains $h_i(\mathbf{x}^*) = 0, i = 1, \dots, p, \lambda_i \geq 0$ and $f_i(\mathbf{x}^*) \leq 0, i = 1, \dots, m$. The equality between the solutions of the primal (A.2) and associated dual problem (A.4) is achieved if the problem has *strong duality* [140, pg. 226]. Convex problems with affine equality constrains as

$$\begin{aligned} & \text{minimize} && f_0(\mathbf{x}) \\ & \text{subject to} && f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & && \mathbf{Ax} = \mathbf{b}. \end{aligned} \tag{A.5}$$

with f_0, \dots, f_m convex, usually (but not always!) have strong duality. For problems in the form (A.5), strong duality can be verified using Slater's condition: $\exists \mathbf{x} \in \mathbf{relint} \mathcal{D}$ ¹ such that

$$f_i(\mathbf{x}) < 0, i = 1, \dots, m, \quad \mathbf{Ax} = \mathbf{b}.$$

Such point is sometimes called strictly feasible, since the inequality constrains hold strict inequalities. If k inequality constrains are affine, the the Slater's condition can be redefined to: $\exists \mathbf{x} \in \mathbf{relint} \mathcal{D}$ with

$$f_i(\mathbf{x}) \leq 0, i = 1, \dots, k, \quad f_i(\mathbf{x}) < 0, i = k + 1, \dots, m, \quad \mathbf{Ax} = \mathbf{b}.$$

Strong duality implies $g(\boldsymbol{\lambda}^*, \mathbf{v}^*) = f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) + \sum_{i=1}^p v_i^* h_i(\mathbf{x}^*) = f_0(\mathbf{x}^*)$, and hence that

$$\sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) = 0,$$

since each term in the sum is non-negative,

$$\lambda_i^* f_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m.$$

Thus, we have $\lambda_i^* > 0 \implies f_i(\mathbf{x}^*) = 0$, or equivalently, $f_i(\mathbf{x}^*) < 0 \implies \lambda_i^* = 0$.

Now, assuming that f_0, \dots, f_m and h_1, \dots, h_p are differentiable, and since \mathbf{x}^* minimizes $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \mathbf{v}^*)$ in \mathbf{x} , then the gradient of the Lagrangean

¹The *relative interior* of a set $C \subseteq \mathbb{R}^n$, denoted $\mathbf{relint} C$, is defined as the interior related to the affine set $\mathbf{aff} C$, i.e., $\mathbf{relint} C = \{\mathbf{x} \in C \mid \mathcal{B}(\mathbf{x}, r) \cap \mathbf{aff} C \subseteq C \text{ for any } r > 0\}$.

must vanish in \mathbf{x}^*

$$\nabla f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(\mathbf{x}^*) + \sum_{i=1}^p \mathbf{v}_i^* \nabla h_i(\mathbf{x}^*) = 0.$$

Thus, we can write the Karush-Kuhn-Tucker (KKT) optimality conditions as

$$\begin{aligned} f_i(\mathbf{x}^*) &\leq 0, & i = 1, \dots, m \\ h_i(\mathbf{x}^*) &= 0, & i = 1, \dots, p \\ \lambda_i^* &\geq 0, & i = 1, \dots, m \\ \lambda_i^* f_i(\mathbf{x}^*) &= 0, & i = 1, \dots, m \end{aligned} \quad (\text{A.6})$$

$$\nabla f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(\mathbf{x}^*) + \sum_{i=1}^p \mathbf{v}_i^* \nabla h_i(\mathbf{x}^*) = 0.$$

If the problem is convex, then the KKT conditions are also sufficient conditions, *i.e.*, if a candidate point $(\tilde{\mathbf{x}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{v}})$ satisfies the KKT conditions, then this point is global optimum for the primal and dual problems. Therefore, $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \mathbf{v}^*) = (\tilde{\mathbf{x}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{v}})$.

A.2.1 Example for the regularized LS

Consider an input-output sequence $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, with $\mathbf{x} \in \mathbb{R}^D$ and scalar y , governed by the model $y_i = f(\mathbf{x}_i) + n_i$, where n_i is an additive noise independent of \mathbf{x}_i , the function $f(\mathbf{x}_i) = \boldsymbol{\varphi}(\mathbf{x}_i)^\top \mathbf{w}$, linear in $\mathbf{w} \in \mathbb{R}^M$, and considering the mapping $\boldsymbol{\varphi} : \mathbb{R}^D \rightarrow \mathbb{R}^M$, we can write the primal convex problem as

$$\begin{aligned} \min \quad & \frac{1}{2\mu} \sum_{i=1}^N e_i^2 + \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & e_i = y_i - \boldsymbol{\varphi}(\mathbf{x}_i)^\top \mathbf{w}, \quad i = 1, \dots, N, \end{aligned} \quad (\text{A.7})$$

where μ is the regularization parameter. We can write the Lagrangean for (A.7) as

$$\mathcal{L}(\mathbf{w}, \mathbf{e}, \boldsymbol{\beta}) = \frac{1}{2\mu} \sum_{i=1}^N e_i^2 + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \beta_i (e_i - y_i + \boldsymbol{\varphi}(\mathbf{x}_i)^\top \mathbf{w}). \quad (\text{A.8})$$

Since $\mathcal{L}(\mathbf{w}, \mathbf{e}, \boldsymbol{\beta})$ is convex with respect to the primal variables (\mathbf{w}, \mathbf{e}) , we have $\nabla_{(\mathbf{w}, \mathbf{e})} \mathcal{L}(\mathbf{w}^*, \mathbf{e}^*, \boldsymbol{\beta}) = 0$. Thus, we can find the arguments $(\mathbf{w}^*, \mathbf{e}^*)$ which

minimize (A.8) by doing

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial e_i} &= \frac{1}{\mu} e_i - \beta_i = 0 \implies e_i^* = \mu \beta_i \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^N \beta_i \boldsymbol{\varphi}(\mathbf{x}_i) = 0 \implies \mathbf{w}^* = \boldsymbol{\Phi}^\top \boldsymbol{\beta}\end{aligned}\tag{A.9}$$

where $\boldsymbol{\Phi} = [\boldsymbol{\varphi}(\mathbf{x}_1), \dots, \boldsymbol{\varphi}(\mathbf{x}_N)]^\top$. Replacing $(\mathbf{w}^*, \mathbf{e}^*)$ in (A.8) we obtain

$$\begin{aligned}g(\boldsymbol{\beta}) &= \mathcal{L}(\mathbf{w}^*, \mathbf{e}^*, \boldsymbol{\beta}) \\ &= \frac{1}{2\mu} \mu^2 \boldsymbol{\beta}^\top \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\beta} - \mu \boldsymbol{\beta}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{y} - \boldsymbol{\beta}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\beta} \\ &= -\frac{1}{2} \boldsymbol{\beta}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \mu \mathbf{I}) \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{y}.\end{aligned}\tag{A.10}$$

Thus, we can formulate the following convex unconstrained dual problem

$$\boldsymbol{\beta}^* = \arg \max_{\boldsymbol{\beta}} -\frac{1}{2} \boldsymbol{\beta}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \mu \mathbf{I}) \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{y},\tag{A.11}$$

where the global optimum solution $d^* = p^*$ is achieved by solving (A.11), resulting in

$$\boldsymbol{\beta}^* = (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \mu \mathbf{I})^{-1} \mathbf{y}.\tag{A.12}$$

Now, for any vector \mathbf{x} , we can write the predictive form of $f(\mathbf{x})$ as

$$f(\mathbf{x}) = \boldsymbol{\varphi}(\mathbf{x})^\top \mathbf{w}^* = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Phi}^\top \boldsymbol{\beta}^* = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Phi}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \mu \mathbf{I})^{-1} \mathbf{y}.\tag{A.13}$$

Interestingly, the dual formulation allows us to work in a space whose dimension is N (number of data), as $\boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Phi}^\top$ is a $1 \times N$ vector and $(\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \mu \mathbf{I})$ is a $N \times N$ matrix, of inner products $\boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\varphi}(\mathbf{x}')$ in \mathbb{R}^M , while in the original formulation we were working in a space with dimension M . Another characteristic of the dual representation is that the coefficient vector \mathbf{w} can be calculated in a implicit manner by solving the dual problem (A.11) for $\boldsymbol{\beta}$. This allow us to find a solution entirely in \mathbb{R}^n . Also note that if we use other traditional forms for solving the optimization problem (A.7), such as [59]

$$\min \frac{1}{2\mu} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|^2 + \frac{1}{2} \|\mathbf{w}\|^2,\tag{A.14}$$

whose solution is

$$\mathbf{w}^* = \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{y} \quad (\text{A.15})$$

and has the following predictive form

$$f(\mathbf{x}) = \boldsymbol{\varphi}(\mathbf{x})^\top \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{y} \quad (\text{A.16})$$

which can no longer be written as a function of N -dimensional matrices and vectors, neither as a function of inner products in the space of $\boldsymbol{\varphi}$. Although it is necessary to invert $N \times N$ matrices (which can be really large matrices) when using the dual form, this formulation becomes particularly interesting when M is very large, specially for cases in which M is infinite. Thus the dual formulation can be very attractive as it will be seen in the *kernel* based formulation.

A.3 RKHS

In this section we will discuss about reproducing kernel Hilbert spaces (RKHSs), kernels and their relation to feature spaces.

A.3.1 Important concepts and definitions from Functional Analysis

Before dealing with RKHS we will introduce a few key concepts from functional analysis in the form of definitions and theorems. Those theorems not properly demonstrated here are accompanied by references that have their proof with the proper rigour. The theory of RKHS and functional analysis, usually assumes vector spaces (or linear spaces) of complex functions built over abstract fields. Here, however, we will consider only real functional spaces \mathcal{F} defined over real fields.

Definition 1 (Metric space, metric). *A metric space is a pair (\mathcal{F}, d) where \mathcal{F} is a set and d is a metric on \mathcal{F} (or distance function on \mathcal{F}), that is, a function defined on $\mathcal{F} \times \mathcal{F}$ such that for all $f, g, h \in \mathcal{F}$ we have:*

- (i) d is real-valued, finite and nonnegative.
- (ii) $d(f, g) = 0$ if and only if $f = g$.
- (iii) $d(f, g) = d(g, f)$ (Symmetry).
- (iv) $d(f, g) \leq d(f, h) + d(h, g)$ (**Triangle inequality**);

Definition 2 (Cauchy sequence). A sequence $\{f_n\}_{n=1}^{\infty}$ of elements of a metric space \mathcal{F} with metric d is said Cauchy (or fundamental) if for all $\varepsilon > 0$ there exists a $N = N(\varepsilon) \in \mathbb{N}$, such that

$$d(f_m, f_n) < \varepsilon \quad \text{for all } m, n > N.$$

Definition 3 (Complete space). A Space \mathcal{F} is complete if all Cauchy sequences in \mathcal{F} converge (i.e., the sequence has a limit and this limit is an element of \mathcal{F}).

Definition 4 (Normed spaces, Banach spaces). A normed space² \mathcal{F} is a vector space with a norm defined on it. A Banach space is a complete normed (complete in the metric defined by the norm; see (A.17)). Here a norm on a linear space \mathcal{F} is a real-valued function $\|\cdot\|_{\mathcal{F}} : \mathcal{F} \rightarrow \mathbb{R}$ whose value at an $f \in \mathcal{F}$ is denoted by

$$\|f\|_{\mathcal{F}} \quad (\text{read "norm of } f\text{"})$$

and which has the properties

- (i) $\|f\|_{\mathcal{F}} \geq 0$
- (ii) $\|f\|_{\mathcal{F}} = 0 \Leftrightarrow f = 0$
- (iii) $\|\alpha f\|_{\mathcal{F}} = |\alpha| \|f\|_{\mathcal{F}}$
- (iv) $\|f + g\|_{\mathcal{F}} \leq \|f\|_{\mathcal{F}} + \|g\|_{\mathcal{F}}$ (Triangle inequality);

here f and g are arbitrary vectors in \mathcal{F} and α is any scalar.

A norm on \mathcal{F} defines a metric d on \mathcal{F} which is given by

$$d(f, g) = \|f - g\|_{\mathcal{F}}, \quad (f, g, \in \mathcal{F}) \quad (\text{A.17})$$

and is called metric induced by the norm. The normed space just defined is denoted by $(\mathcal{F}, \|\cdot\|_{\mathcal{F}})$ or simply by \mathcal{F} .

Definition 5 (Hilbert space). An inner product space (or pre-Hilbert space) is a vector \mathcal{H} with an inner product defined on \mathcal{H} . A Hilbert space is a complete inner product space (complete in the metric defined by the inner product; cf. (A.19) below). Here, an inner product on \mathcal{H} is a mapping $\mathcal{H} \times \mathcal{H}$ into the scalar field \mathbb{R} of \mathcal{H} ; that is, with every pair of vectors $f, g \in \mathcal{H}$ there is associated a scalar which is written

$$\langle f, g \rangle_{\mathcal{H}}$$

²Also called *normed vector space* or *normed linear space*.

and is called the inner product (or scalar product) of f and g , such that for all vectors $f, g, h \in \mathcal{H}$ and scalars α we have

- (i) $\langle f + g, h \rangle_{\mathcal{H}} = \langle f, h \rangle_{\mathcal{H}} + \langle g, h \rangle_{\mathcal{H}}$
- (ii) $\langle \alpha f, g \rangle_{\mathcal{H}} = \alpha \langle f, g \rangle_{\mathcal{H}}$
- (iii) $\langle f, g \rangle_{\mathcal{H}} = \overline{\langle g, f \rangle_{\mathcal{H}}} = \langle g, f \rangle_{\mathcal{H}}$ (the bar denotes complex conjugation.)
- (iv) $\langle f, f \rangle_{\mathcal{H}} \geq 0$
 $\langle f, f \rangle_{\mathcal{H}} = 0 \Leftrightarrow f = 0.$

Note that in (iii) the conjugate can be ignored since we are assuming only fields of real numbers.

An inner product on \mathcal{H} defines a norm on \mathcal{H} given by

$$\|f\|_{\mathcal{H}} = \sqrt{\langle f, f \rangle_{\mathcal{H}}} \quad (\text{A.18})$$

and a metric on \mathcal{H} given by

$$d(f, g) = \|f - g\|_{\mathcal{H}} = \sqrt{\langle f - g, f - g \rangle_{\mathcal{H}}}. \quad (\text{A.19})$$

Hence inner product spaces are normed spaces, and Hilbert spaces are Banach spaces [98]. An inner product and the corresponding norm satisfy the **Cauchy-Schwarz inequality** [98, pg. 137], therefore, we have

$$|\langle f, g \rangle_{\mathcal{H}}| \leq \|f\|_{\mathcal{H}} \|g\|_{\mathcal{H}}. \quad (\text{A.20})$$

A functional is an operator whose range is on the real line \mathbb{R} or in the complex plane \mathbb{C} . Here, we deal only with spaces defined on real fields and, therefore, we present the following definition for linear functionals.

Definition 6 (Continuous/bounded linear functional). *A linear functional Ω is a linear mapping from the vector space \mathcal{H} to the scalar field \mathbb{R} :*

$$\Omega : \mathcal{H} \rightarrow \mathbb{R}.$$

If Ω is bounded, then there exists a real number $\lambda > 0$ such that for all $f \in \mathcal{H}$

$$|\Omega[f]| \leq \lambda \|f\|_{\mathcal{H}}.$$

Furthermore, the norm of Ω is

$$\|\Omega\|_{\mathcal{H}'} = \sup_{f \in \mathcal{H}, f \neq 0} \frac{|\Omega[f]|}{\|f\|_{\mathcal{H}}} = \sup_{f \in \mathcal{H}, \|f\|=1} |\Omega[f]|$$

or

$$|\Omega[f]| \leq \|\Omega\|_{\mathcal{H}'} \|f\|_{\mathcal{H}},$$

where \mathcal{H}' is the space of bounded linear functionals (Definition 7).

Theorem 1 (Continuity and boundedness [98, pg.104]). *A linear functional Ω with domain in a normed space \mathcal{H} is continuous if and only if Ω is bounded.*

The set of all bounded linear functionals form a normed space (Banach space [98]) called *dual space*. As a Hilbert space is also a Banach space, we present a definition for dual space considering only Hilbert spaces. However, this definition can be directly carried out to the more general case of Banach spaces.

Definition 7. [Dual space \mathcal{H}'] *Let \mathcal{H} be a Hilbert space. Then the set of all bounded linear functionals on Ω em \mathcal{H} constitutes a Hilbert space with norm defined by*

$$\|\Omega\|_{\mathcal{H}'} = \sup_{f \in \mathcal{H}, f \neq 0} \frac{|\Omega[f]|}{\|f\|_{\mathcal{H}}} = \sup_{f \in \mathcal{H}, \|f\|_{\mathcal{H}}=1} |\Omega[f]|$$

which is called **dual space** of \mathcal{H} and is denoted by \mathcal{H}' .

Definition 8 (Definite positive functions). *A function $h : X \times X \rightarrow \mathbb{R}$ is called positive definite if, $\forall n \in \mathbb{N}$, $\forall \alpha_1, \dots, \alpha_n \in \mathbb{R}$, and $\forall \mathbf{x}_1, \dots, \mathbf{x}_n \in X$, we have*

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j h(\mathbf{x}_i, \mathbf{x}_j) \geq 0. \quad (\text{A.21})$$

Furthermore, h is said **strictly positive definite** if, for all mutually distinct $\forall \mathbf{x}_1, \dots, \mathbf{x}_n \in X$, the equality in (A.21) holds only when all $\alpha_i = \dots = \alpha_n = 0$. Finally, h is said **symmetric** if $h(\mathbf{x}, \mathbf{x}') = h(\mathbf{x}', \mathbf{x})$ for all $\mathbf{x}, \mathbf{x}' \in X$.

Note that in order to a real function h be positive definite it must be symmetric [141, pg.14, Lema 4]. In the literature there is no consensus about Definition 8. Often the term *positive definite* is applied to strictly positive definite functions, and the term *positive semidefinite* to positive definite functions.

Theorem 2 (Riesz's representation). *Let \mathcal{H} be a Hilbert space and $f \in \mathcal{H}$ a element of this space. Every bounded linear functional Ω on \mathcal{H} can be represented in terms of the inner product, namely,*

$$\Omega[f] = \langle f, h_\Omega \rangle_{\mathcal{H}}$$

where $h_\Omega \in \mathcal{H}$ depends on Ω , is uniquely determined by Ω and has norm

$$\|h_\Omega\|_{\mathcal{H}} = \|\Omega\|_{\mathcal{H}'}$$

The proof for Theorem 2 can be found in [98, pg. 189].

Example 1 (Representation theorem in \mathbb{R}^n). *Consider the case where $\mathcal{H} = \mathbb{R}^n$, and a vector $\mathbf{x} \in \mathbb{R}^n$. For any bounded linear functional $\Omega : \mathcal{H} \rightarrow \mathbb{R}$ there exists a vector $\mathbf{y}_\Omega \in \mathbb{R}^n$ such that (Riesz)*

$$\Omega(\mathbf{x}) = \langle \mathbf{x}, \mathbf{y}_\Omega \rangle_{\mathcal{H}} = \mathbf{x}^\top \mathbf{y}_\Omega = \sum_{i=1}^n y_i x_i$$

which is the form of all linear functional in \mathbb{R}^n . Also note that the vector (of coefficients of the linear functional) \mathbf{y}_Ω is uniquely determined by the functional Ω since $\sum_{i=1}^n y_i x_i$ is the only form of writing a linear functional in \mathbb{R}^n .

A.3.2 RKHS and reproducing kernels

Let \mathcal{H} be a Hilbert space of functions mapping a nonempty set $X \subset \mathbb{R}^n$ to the field real numbers \mathbb{R} . We write the inner product on \mathcal{H} as $\langle f, g \rangle_{\mathcal{H}}$ and the associated norm as $\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}}$. Note that since \mathcal{H} is a space of functions on X , there is for every $\mathbf{x} \in X$ a very special functional on \mathcal{H} : the one that assigns to each $f \in \mathcal{H}$ its value at \mathbf{x} .

Definition 9 (Evaluation functional). *Let \mathcal{H} be a Hilbert space of functions $f : X \rightarrow \mathbb{R}$ defined on a nonempty set $X \subset \mathbb{R}^n$. For a fixed $\mathbf{x} \in X$, the map $\delta_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}$, $\delta_{\mathbf{x}} : f \mapsto f(\mathbf{x})$ is called the (Dirac) evaluation functional at \mathbf{x} .*

Note that the evaluation functionals are always linear: for $f, g \in \mathcal{H}$ and $\alpha, \beta \in \mathbb{R}$, $\delta_{\mathbf{x}}[\alpha f + \beta g] = (\alpha f + \beta g)(\mathbf{x}) = \alpha f(\mathbf{x}) + \beta g(\mathbf{x}) = \alpha \delta_{\mathbf{x}}[f] + \beta \delta_{\mathbf{x}}[g]$.

Definition 10 (Reproducing kernel and RKHS). *Let \mathcal{H} be a Hilbert space of real functions $f : X \rightarrow \mathbb{R}$ defined on the nonempty set X .*

(i) a function $\kappa : X \times X \rightarrow \mathbb{R}$ is called a **reproducing kernel** of \mathcal{H} if

$\kappa(\cdot, \mathbf{x}) \in \mathcal{H}$ for all $\mathbf{x} \in X$ and the **reproducing property**

$$f(\mathbf{x}) = \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}$$

holds for all $f \in \mathcal{H}$ and all $\mathbf{x} \in X$.

(ii) The space \mathcal{H} is called **reproducing kernel Hilbert space (RKHS)** on X if for all $\mathbf{x} \in X$ the evaluation functional $\delta_{\mathbf{x}}$ is continuous (or equivalently bounded): $\forall \mathbf{x} \in X$ there exists a scalar $\lambda_{\mathbf{x}} \geq 0$ such that $|f(\mathbf{x})| = |\delta_{\mathbf{x}}[f]| \leq \lambda_{\mathbf{x}} \|f\|_{\mathcal{H}}$.

Note that that if \mathcal{H} is a RKHS, then we can write (using the Riesz's theorem 2)

$$\delta_{\mathbf{x}}[f] = f(\mathbf{x}) = \langle f, h_{\delta_{\mathbf{x}}} \rangle_{\mathcal{H}},$$

where $h_{\delta_{\mathbf{x}}} \in \mathcal{H}$ is a function uniquely determined by $\delta_{\mathbf{x}}$, and that condition (i) in Definition 10 implies $h_{\delta_{\mathbf{x}}} = \kappa(\cdot, \mathbf{x})$. Thus

$$f(\mathbf{x}) = \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}$$

and

$$\delta_{\mathbf{y}}[\kappa(\cdot, \mathbf{x})] = \kappa(\mathbf{x}, \mathbf{y}) = \langle \kappa(\cdot, \mathbf{x}), h_{\delta_{\mathbf{y}}} \rangle_{\mathcal{H}} = \langle \kappa(\cdot, \mathbf{x}), \kappa(\cdot, \mathbf{y}) \rangle_{\mathcal{H}},$$

showing the reproducing property. Therefore we say that $\kappa(\cdot, \mathbf{x}) = h_{\delta_{\mathbf{x}}}$ is the representer of the evaluation functional $\delta_{\mathbf{x}}$, or just *representer of the evaluation* at \mathbf{x} . The fact that $\kappa(\cdot, \mathbf{x})$ is the representer of the evaluation at \mathbf{x} implies that κ is uniquely determined by the evaluation functional, and consequently each RKHS has just one reproducing kernel.

Theorem 3 (Existence of the reproducing kernel). *\mathcal{H} is a RKHS (i.e. its evaluation functionals $\delta_{\mathbf{x}}$ are continuous linear operators) if and only if \mathcal{H} has a reproducing kernel [96].*

Proof. Given that a Hilbert space \mathcal{H} has a reproducing kernel κ with the reproducing property $\langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = f(\mathbf{x})$, then

$$\begin{aligned} |\delta_{\mathbf{x}}[f]| &= |f(\mathbf{x})| \\ &= |\langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}| \\ &\leq \|\kappa(\cdot, \mathbf{x})\|_{\mathcal{H}} \|f\|_{\mathcal{H}} \quad (\text{Cauchy-Schwarz (A.20)}) \\ &= \langle \kappa(\cdot, \mathbf{x}), \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}^{1/2} \|f\|_{\mathcal{H}} \\ &= \kappa(\mathbf{x}, \mathbf{x})^{1/2} \|f\|_{\mathcal{H}}. \end{aligned}$$

consequently, $\delta_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}$ is a bounded linear operator. □

Theorem 4. *If a reproducing kernel κ exists it is unique [96].*

Proof. In fact, if another reproducing kernel $\tilde{\kappa}$ existed, we would have for some $\mathbf{x} \in X$

$$\begin{aligned}
 0 < \|\kappa(\cdot, \mathbf{x}) - \tilde{\kappa}(\cdot, \mathbf{x})\|_{\mathcal{H}}^2 &= \langle \kappa(\cdot, \mathbf{x}) - \tilde{\kappa}(\cdot, \mathbf{x}), \kappa(\cdot, \mathbf{x}) - \tilde{\kappa}(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} \\
 &= \langle \kappa(\cdot, \mathbf{x}) - \tilde{\kappa}(\cdot, \mathbf{x}), \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} - \langle \kappa(\cdot, \mathbf{x}) - \tilde{\kappa}(\cdot, \mathbf{x}), \tilde{\kappa}(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} \\
 &= \kappa(\mathbf{x}, \mathbf{x}) - \tilde{\kappa}(\mathbf{x}, \mathbf{x}) - \kappa(\mathbf{x}, \mathbf{x}) + \tilde{\kappa}(\mathbf{x}, \mathbf{x}) \quad (\text{reproducing prop.}) \\
 &= 0.
 \end{aligned}
 \tag{A.22}$$

□

RKHSs also have the important property that norm convergence implies pointwise convergence [141, pg.10].

Corollary 1 (Convergence implies pointwise convergence). *If two functions converge in RKHS norm, then they converge at every point, i.e., if $\lim_{n \rightarrow \infty} \|f_n - f\|_{\mathcal{H}} = 0$, then $\lim_{n \rightarrow \infty} f_n(\mathbf{x}) = f(\mathbf{x})$, $\forall \mathbf{x} \in X$.*

Proof. For any $\mathbf{x} \in X$,

$$\begin{aligned}
 |f_n(\mathbf{x}) - f(\mathbf{x})| &= |\delta_{\mathbf{x}}[f_n] - \delta_{\mathbf{x}}[f]| \\
 &= |\langle f_n, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} - \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}| \\
 &= |\langle f_n - f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}| \\
 &\leq \|\kappa(\cdot, \mathbf{x})\|_{\mathcal{H}} \|f_n - f\|_{\mathcal{H}} \quad (\text{Cauchy-Schwarz (A.20)}), \\
 &= \|\delta_{\mathbf{x}}\|_{\mathcal{H}'} \|f_n - f\|_{\mathcal{H}} \quad (\text{Theorem 2}),
 \end{aligned}$$

where $\|\delta_{\mathbf{x}}\|_{\mathcal{H}'}$ is the norm of the evaluation functional (which is bounded by definition on the RKHS). □

The last step of the above proof used the fact that $\kappa(\cdot, \mathbf{x})$ is the representer of the evaluation at \mathbf{x} whose consequence to norm is presented in Theorem 2.

Theorem 5 (Every reproducing kernel is a positive definite function).

Proof.

$$\begin{aligned}
 \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle \kappa(\cdot, \mathbf{x}_i), \kappa(\cdot, \mathbf{x}_j) \rangle_{\mathcal{H}} \\
 &= \left\langle \sum_{i=1}^n \alpha_i \kappa(\cdot, \mathbf{x}_i), \sum_{j=1}^n \alpha_j \kappa(\cdot, \mathbf{x}_j) \right\rangle_{\mathcal{H}} \\
 &= \left\| \sum_{i=1}^n \alpha_i \kappa(\cdot, \mathbf{x}_i) \right\|_{\mathcal{H}}^2 \geq 0.
 \end{aligned}$$

□

Theorem 5 admits a converse presented by the following theorem [96].

Theorem 6 (Moore-Aronszajn). *To every positive definite function κ on $X \times X$ there exists only one Hilbert space \mathcal{H} of functions on X with κ as reproducing kernel. The subspace $\mathcal{H}_0 \subset \mathcal{H}$ spanned by the functions $\{\kappa(\cdot, \mathbf{x})\}_{\mathbf{x} \in X}$ is dense in \mathcal{H} and \mathcal{H} is the set of functions on X which are pointwise limits of Cauchy sequences in \mathcal{H}_0 with inner product*

$$\langle f, g \rangle_{\mathcal{H}_0} = \sum_i^n \sum_j^m \alpha_i \beta_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \tag{A.23}$$

where $f = \sum_i^n \alpha_i \kappa(\cdot, \mathbf{x}_i)$ and $g = \sum_j^m \beta_j \kappa(\cdot, \mathbf{x}_j)$.

The proof of Theorem 6 is long and complex. Therefore, the following discussion addresses only the main steps to conduct the proof. The interested reader can consult [141] to a complete proof of Theorem 6 as well as the original references [97, 95].

First it can be shown that Equation (A.23) is a valid inner product on \mathcal{H}_0 . However, \mathcal{H}_0 is not a complete space yet, and therefore is not a RKHS. Nevertheless it is possible to complete \mathcal{H}_0 forming a Hilbert space \mathcal{H} . \mathcal{H} is a space of functions f for which there exists a Cauchy sequence f_n in \mathcal{H}_0 converging pointwise to f , with inner product defined as

$$\langle f, g \rangle_{\mathcal{H}} = \lim_{n \rightarrow \infty} \langle f_n, g_n \rangle_{\mathcal{H}_0}. \tag{A.24}$$

implying that \mathcal{H} is unique (except for isometries [98, pg. 41]) and that its topology is induced from \mathcal{H}_0 . It can also be shown that κ is the reproducing

kernel of \mathcal{H} , since for any function $f \in \mathcal{H}$ we have

$$\begin{aligned} \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} &= \lim_{n \rightarrow \infty} \langle f_n, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_0} \\ &= \lim_{n \rightarrow \infty} f_n(\mathbf{x}) \\ &= f(\mathbf{x}), \end{aligned} \tag{A.25}$$

and that the evaluation functional in \mathcal{H} is continuous, and that pointwise convergence implies norm convergence, presenting the same characteristics and properties already presented to RKHSs.

Theorem 7 (Feature space). *A symmetric function $\kappa : X \times X \rightarrow \mathbb{R}$ is a reproducing kernel, or a positive definite function, if and only if there exists a map φ from X to some space of convergent sequences $\ell^2(A)$ (where A is the index set) such that*

$$\begin{aligned} \forall (\mathbf{x}, \mathbf{y}) \in X \times X \quad \kappa(\mathbf{x}, \mathbf{y}) &= \langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle_{\ell^2(A)} \\ &= \sum_{\alpha \in A} (\varphi(\mathbf{x}))_{\alpha} (\varphi(\mathbf{y}))_{\alpha}. \end{aligned}$$

Below we present the very interesting definition of universal kernel [99, corollary 4.52 pg. 152].

Definition 11 (Universal kernel). *Let $C(X)$ be the space of continuous functions $f : X \rightarrow \mathbb{R}$. A continuous kernel κ on a compact metric space X is called **universal** if the RKHS \mathcal{H} of κ is dense in $C(X)$, i.e., for every function $g \in C(X)$ and all $\varepsilon > 0$ there exists an $f \in \mathcal{H}$ such that*

$$\|f - g\|_{\infty} \leq \varepsilon.$$

Universal kernels produce RKHSs rich enough to provide arbitrarily accurate function approximations for all distributions. This guarantee learning in the absence of assumptions on the data-generating distribution. However, this flexibility also carries the danger of overfitting. Some examples of universal kernels are presented below.

Example 2 (Examples of universal kernels). *Let X be a compact subset (i.e. closed and bounded [98, pg. 77, Theorem 2.5-3]) of \mathbb{R}^n , $\sigma > 0$ and $\alpha > 0$. The the following kernels on X are universal:*

- kernel exponential:

$$\kappa(\mathbf{x}, \mathbf{x}') := \exp \{ \langle \mathbf{x}, \mathbf{x}' \rangle \} \tag{A.26}$$

- *kernel Gaussian (or RBF):*

$$\kappa(\mathbf{x}, \mathbf{x}') := \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2 \right\} \quad (\text{A.27})$$

- *kernel binomial:*

$$\kappa(\mathbf{x}, \mathbf{x}') := (1 - \langle \mathbf{x}, \mathbf{x}' \rangle)^{-\alpha} \quad (\text{A.28})$$

where for the last kernel we additionally assume $X \subset \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < 1\}$.

A.4 DIRECTIONAL DERIVATIVES ON HILBERT SPACES

Directional derivatives and gradients are essential to several optimization techniques. This section introduces the concept of directional derivatives for Hilbert spaces.

Let J be a functional

$$\begin{aligned} J : \mathcal{H} &\longrightarrow \mathbb{R} \\ f &\longmapsto J(f) \end{aligned} \quad (\text{A.29})$$

where \mathcal{H} is a Hilbert space. If for two elements $f, g \in \mathcal{H}$, the limit

$$\partial_g J(f) = \lim_{\varepsilon \rightarrow 0} \frac{J(f + \varepsilon g) - J(f)}{\varepsilon} \quad (\text{A.30})$$

exists, then $\partial_g J(f)$ is called directional derivative of J at f in the direction of g . If the limit (A.30) exists for all $g \in \mathcal{H}$, then J is said to be directionally differentiable at f [142, pg. 38]. Since the directional derivative is a linear functional, we can use Theorem 2 (Riesz's representation) to represent it as an inner product on \mathcal{H} . Thus, the gradient $\nabla J(f)$ of J at $f \in \mathcal{H}$, if it exists, satisfies [143, pg. 139]

$$\partial_g J(f) = \langle \nabla J(f), g \rangle_{\mathcal{H}}, \quad \text{para todo } g \in \mathcal{H}. \quad (\text{A.31})$$

We now present some illustrative examples.

Example 3 (Gradient on \mathbb{R}^n). Consider the set of vectors $\mathbf{e}_i = [0, \dots, 1, \dots, 0]^\top$, $i = 1, \dots, n$, forming an orthonormal basis in \mathbb{R}^n and the functional $J : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto J(\mathbf{x})$, such that $J(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$ with $\mathbf{a}, \mathbf{x} \in \mathbb{R}^n$. The directional derivative

at \mathbf{x} in the direction of \mathbf{e}_i can be written, using Riesz's theorem, as

$$\begin{aligned}\partial_{\mathbf{e}_i} J(\mathbf{x}) &= \langle \mathbf{e}_i, \nabla J(\mathbf{x}) \rangle \\ &= \mathbf{e}_i^\top \nabla J(\mathbf{x}) \\ &= \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_1} 0 + \dots + \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_i} 1 + \dots + \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_n} 0 \\ &= \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_i} = a_i,\end{aligned}$$

where used $\nabla J(\mathbf{x}) = [\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_1}, \dots, \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_i}, \dots, \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial x_n}]^\top = \mathbf{a}$ in the third line. Note, however, that we can achieve the same result above using the definition of directional derivative in (A.30)

$$\begin{aligned}\partial_{\mathbf{e}_i} J(\mathbf{x}) &= \lim_{\varepsilon \rightarrow 0} \frac{J(\mathbf{x} + \varepsilon \mathbf{e}_i) - J(\mathbf{x})}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\mathbf{a}^\top \mathbf{x} + \varepsilon \mathbf{a}^\top \mathbf{e}_i - \mathbf{a}^\top \mathbf{x}}{\varepsilon} \\ &= \mathbf{a}^\top \mathbf{e}_i = \langle \mathbf{a}, \mathbf{e}_i \rangle \\ &= a_i.\end{aligned}$$

Note that in the third line we can conclude that $\nabla J(\mathbf{x}) = \mathbf{a}$.

Example 4 ($J(f) = \|f\|_{\mathcal{H}}^2$).

$$\begin{aligned}\partial_g J(f) &= \lim_{\varepsilon \rightarrow 0} \frac{\|f + \varepsilon g\|_{\mathcal{H}}^2 - \|f\|_{\mathcal{H}}^2}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\|f\|_{\mathcal{H}}^2 + 2\varepsilon \langle f, g \rangle_{\mathcal{H}} + \varepsilon^2 \|g\|_{\mathcal{H}}^2 - \|f\|_{\mathcal{H}}^2}{\varepsilon} \quad (\text{A.32}) \\ &= \langle 2f, g \rangle_{\mathcal{H}}.\end{aligned}$$

Thus the gradient is $\nabla J(f) = 2f$.

Example 5 ($J(f) = f(\mathbf{x})$, $f \in \mathcal{H}$). Using the reproducing property of the RKHS we have $f(\mathbf{x}) = \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}$, thus

$$\begin{aligned}\partial_g J(f) &= \lim_{\varepsilon \rightarrow 0} \frac{\langle f + \varepsilon g, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} - \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} + \varepsilon \langle g, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} - \langle f, \kappa(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}}{\varepsilon} \quad (\text{A.33}) \\ &= \langle \kappa(\cdot, \mathbf{x}), g \rangle_{\mathcal{H}},\end{aligned}$$

and $\nabla J(f) = \kappa(\cdot, \mathbf{x})$.

A.5 THE REPRESENTER THEOREM

Here we present a generalized version of the Wahba's representer theorem [144].

Theorem 8 (Nonparametric representer theorem). *Let X be a nonempty set, κ a real-valued kernel on $X \times X$, a training sample $\{\mathbf{x}_i, y_i\}_{i=1}^N \in X \times \mathbb{R}$, a strictly monotonically increasing real-valued function g on $[0, \infty)$, an arbitrary cost function $c : (X \times \mathbb{R}^2)^N \rightarrow \mathbb{R} \cup \{\infty\}$, and a class of functions*

$$\mathcal{F} = \left\{ f : X \rightarrow \mathbb{R} \mid f(\cdot) = \sum_{i=1}^{\infty} \beta_i \kappa(\cdot, \mathbf{x}_i), \beta_i \in \mathbb{R}, \mathbf{x}_i \in X, \|f\|_{\mathcal{H}} < \infty \right\}. \quad (\text{A.34})$$

Then any $f \in \mathcal{F}$ minimizing the regularized risk functional

$$c((\mathbf{x}_1, y_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_N, y_N, f(\mathbf{x}_N))) + g(\|f\|_{\mathcal{H}}) \quad (\text{A.35})$$

admits a representation of the form

$$f(\cdot) = \sum_{i=1}^N \alpha_i \kappa(\cdot, \mathbf{x}_i). \quad (\text{A.36})$$

The proof for Theorem 8 is relatively simple and can be found in [145].

A.6 LS-SVM FOR REGRESSION (LS-SVR)

We return now to the regularized LS problem presented in Section A.2.1. Let us consider an input-output sequence $\{\mathbf{x}_i, y_i\}_{i=1}^N$, with $\mathbf{x} \in \mathbb{R}^D$ and scalar y . Here, however, let us consider a more general (possibly nonlinear) relation $y_i = \psi(\mathbf{x}_i) + n_i$, where n_i is an additive noise independent of \mathbf{x}_i and $\psi \in \mathcal{H}$ is a function of a given functional Hilbert space \mathcal{H} . Thus, we can formulate the following convex (regularized) optimization problem, also called *least-squares support vector machine* (LS-SVM), as

$$\min \frac{1}{2} \|\psi\|_{\mathcal{H}}^2 + \frac{1}{2\mu} \sum_{i=1}^N e_i^2 \quad (\text{A.37})$$

such that $e_i = y_i - \psi(\mathbf{x}_i)$, $i = 1, \dots, N$,

where $\|\cdot\|_{\mathcal{H}} = \langle \cdot, \cdot \rangle_{\mathcal{H}}$ is the norm on \mathcal{H} . Using Lagrange multipliers β_ℓ , we present the Lagrangean function as

$$\mathcal{L}(\boldsymbol{\psi}, \mathbf{e}, \boldsymbol{\beta}) = \frac{1}{2} \|\boldsymbol{\psi}\|_{\mathcal{H}}^2 + \frac{1}{2\mu} \sum_{i=1}^N e_i^2 - \sum_{i=1}^N \beta_i (e_i - y_i + \boldsymbol{\psi}(\mathbf{x}_i)). \quad (\text{A.38})$$

Analogously to the LS case, the Lagrangean function is convex with respect to the primal variables $\boldsymbol{\psi}$ and e_i . Thus we have $\nabla_{\boldsymbol{\psi}, \mathbf{e}} \mathcal{L}(\boldsymbol{\psi}, \mathbf{e}, \boldsymbol{\beta}) = \mathbf{0}$, and

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\psi}} &= \boldsymbol{\psi} - \sum_{i=1}^N \beta_i \boldsymbol{\kappa}(\cdot, \mathbf{x}_i) = 0 \implies \boldsymbol{\psi}^* = \sum_{i=1}^N \beta_i \boldsymbol{\kappa}(\cdot, \mathbf{x}_i) \\ \frac{\partial \mathcal{L}}{\partial e_i} &= \frac{1}{\mu} e_i - \beta_i = 0 \implies e_i^* = \mu \beta_i \end{aligned} \quad (\text{A.39})$$

where we used (A.31) for the functional derivative of $\boldsymbol{\psi}$, and the results (A.32) and (A.33). Note that in the first equation in (A.39) we derived directly the result of the *representer theorem* 8. Replacing the results found in (A.39) in the Lagrangean function (A.38) we have

$$\begin{aligned} g(\boldsymbol{\beta}) &= \mathcal{L}(\boldsymbol{\psi}^*, \mathbf{e}^*, \boldsymbol{\beta}) \\ &= -\frac{1}{2} \boldsymbol{\beta}^\top (\mathbf{K} + \mu \mathbf{I}) \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{y}, \end{aligned} \quad (\text{A.40})$$

where \mathbf{K} is the Gram matrix whose (i, j) -th entry is defined by $\boldsymbol{\kappa}(\mathbf{x}_i, \mathbf{x}_j)$. Now we can state the following dual problem

$$\boldsymbol{\beta}^* = \arg \max_{\boldsymbol{\beta}} -\frac{1}{2} \boldsymbol{\beta}^\top (\mathbf{K} + \mu \mathbf{I}) \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{y}, \quad (\text{A.41})$$

from which we can find the optimal solution making $\nabla_{\boldsymbol{\beta}} \mathcal{L} = \mathbf{0}$,

$$\boldsymbol{\beta}^* = (\mathbf{K} + \mu \mathbf{I})^{-1} \mathbf{y}. \quad (\text{A.42})$$

In this context, for a test sample \mathbf{x} we can write the kernelized predictive form replacing (A.42) in the expression for $\boldsymbol{\psi}^*$ in (A.39)

$$\begin{aligned} \boldsymbol{\psi}^*(\mathbf{x}_*) &= \boldsymbol{\kappa}_*^\top \boldsymbol{\beta}^* \\ &= \boldsymbol{\kappa}_*^\top (\mathbf{K} + \mu \mathbf{I})^{-1} \mathbf{y}, \end{aligned} \quad (\text{A.43})$$

where $\boldsymbol{\kappa}_*^\top = [\boldsymbol{\kappa}(\mathbf{x}_*, \mathbf{x}_1), \dots, \boldsymbol{\kappa}(\mathbf{x}_*, \mathbf{x}_N)]$.

Note that the chosen formulation in (A.37) and its dual formulation (A.41)

brought us some interesting characteristics. Among these characteristics we can highlight two: directly derive the representer theorem (first equation of (A.39)), and find solutions that are functions of inner products in the feature space, which we can calculate implicitly using kernels. This last characteristic allowed us to work in the data space, with N -dimensional vectors and matrices (*kernel trick*) avoiding the problem known as the *curse of dimensionality*. As an example, the Gaussian kernel represents the inner product in an infinite-dimensional RKHS. Therefore it would be impossible to work directly in such a large space.

APPENDIX B – Gaussian Process Regression

This appendix shows the application of Gaussian processes regression to solve nonlinear problems. But first let us return to the linear regression problem presented in Section A.2.1, but now, within a Bayesian formulation.

B.1 REVISITING THE LINEAR REGRESSION

Consider a sequence of input-output pairs $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ to be modeled following $y_i = f(\mathbf{x}_i) + n_i$, where $\mathbf{x} \in \mathbb{R}^D$ is the input, the output y is a scalar, n_i is a zero mean white Gaussian noise, independent from \mathbf{x} , with variance σ_n^2 . Suppose f to be a linear function of $\mathbf{w} \in \mathbb{R}^M$, given by $f(\mathbf{x}_i) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_i)$, where $\boldsymbol{\phi} : \mathbb{R}^D \rightarrow \mathbb{R}^M$. Considering a Bayesian approach, and assuming the values of \mathbf{w} to float around zero, let's consider an isotropic Gaussian PDF as a prior distribution for the weight vectors

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I}) \quad (\text{B.1})$$

where σ_w^2 is the variance of this distribution.

Adopting a vector notation we can write

$$\mathbf{y} = \mathbf{f} + \mathbf{n}, \quad (\text{B.2})$$

where $\mathbf{f} = \mathbf{f}(\mathbf{X}) = \boldsymbol{\Phi} \mathbf{w}$ (the matrix \mathbf{X} was omitted to lighten the notation), $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$, and $\boldsymbol{\Phi} = [\boldsymbol{\phi}(\mathbf{x}_1), \dots, \boldsymbol{\phi}(\mathbf{x}_N)]^\top$. To determine the distribution of \mathbf{y} , note that, for a given \mathbf{X} , $\boldsymbol{\Phi} \mathbf{w}$ is a linear combination of Gaussian random variables and that the noise vector \mathbf{n} is also Gaussian. Therefore, \mathbf{y} has a Gaussian PDF [146, pg. 464] with moments given by

$$\mathbb{E}\{\mathbf{y}\} = \mathbb{E}\{\mathbf{f}\} + \mathbb{E}\{\mathbf{n}\} = \mathbf{0} \quad (\text{B.3})$$

and

$$\begin{aligned} \text{cov}\{\mathbf{y}\} &= \mathbb{E}\{\mathbf{y}\mathbf{y}^\top\} \\ &= \boldsymbol{\Phi} \mathbb{E}\{\mathbf{w}\mathbf{w}^\top\} \boldsymbol{\Phi}^\top + \mathbb{E}\{\boldsymbol{\Phi} \mathbf{w} \mathbf{n}^\top\} + \mathbb{E}\{\mathbf{n} \mathbf{w}^\top \boldsymbol{\Phi}^\top\} + \mathbb{E}\{\mathbf{n} \mathbf{n}^\top\} \\ &= \sigma_w^2 \boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma_n^2 \mathbf{I}. \end{aligned} \quad (\text{B.4})$$

Thus, the distribution of the observations \mathbf{y} , given the data matrix \mathbf{X} , is given by

$$\mathbf{y} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma_n^2 \mathbf{I}) \quad (\text{B.5})$$

where $\mathbf{K} = \sigma_w^2 \boldsymbol{\Phi} \boldsymbol{\Phi}^\top$.

For many regression problems a predictive distribution is desired. This

predictive distribution allows one to “predict” the value of the function $f_* = f(\mathbf{x}_*)$ for a new input \mathbf{x}_* . Thus, following analogous steps considered above, it is easy to show that f_* also has a Gaussian distribution given by

$$f_* | \mathbf{x}_* \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\kappa}_{**}), \quad (\text{B.6})$$

where $\boldsymbol{\kappa}_{**} = \boldsymbol{\varphi}(\mathbf{x}_*)^\top \boldsymbol{\varphi}(\mathbf{x}_*)$. Assuming that \mathbf{y} and f_* are jointly Gaussian [147, pg. 257], the joint PDF for the vector $\mathbf{z} = [\mathbf{y}^\top, f_*]^\top$ has a Gaussian distribution with moments

$$\mathbb{E}\{\mathbf{z}\} = \begin{bmatrix} \mathbb{E}\{\mathbf{y}\} \\ \mathbb{E}\{f_*\} \end{bmatrix} = \mathbf{0} \quad (\text{B.7})$$

and

$$\text{cov}\{\mathbf{z}\} = \mathbb{E}\{\mathbf{z}\mathbf{z}^\top\} = \begin{bmatrix} \mathbb{E}\{\mathbf{y}\mathbf{y}^\top\} & \mathbb{E}\{\mathbf{y}f_*\} \\ \mathbb{E}\{f_*\mathbf{y}^\top\} & \mathbb{E}\{f_*^2\} \end{bmatrix} = \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \boldsymbol{\kappa}_* \\ \boldsymbol{\kappa}_*^\top & \boldsymbol{\kappa}_{**} \end{bmatrix}, \quad (\text{B.8})$$

where $\boldsymbol{\kappa}_* = [\boldsymbol{\varphi}(\mathbf{x}_1)^\top \boldsymbol{\varphi}(\mathbf{x}_*), \dots, \boldsymbol{\varphi}(\mathbf{x}_N)^\top \boldsymbol{\varphi}(\mathbf{x}_*)]^\top$. Thus, we can write the distribution for the vector \mathbf{z} as¹

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \boldsymbol{\kappa}_* \\ \boldsymbol{\kappa}_*^\top & \boldsymbol{\kappa}_{**} \end{bmatrix}\right). \quad (\text{B.9})$$

The property that for jointly Gaussian random variables, the mean vector of the joint distribution can be written stacking the mean vectors of the two marginals, and that the covariance matrix can be written as in (B.8) is sometimes referred to in the literature as *marginalization of property* [100].

The predictive (or posterior) distribution of f_* is then obtained by conditioning (B.9) on the observations \mathbf{y}

$$f_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_* \sim \mathcal{N}\left(\boldsymbol{\kappa}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}, \boldsymbol{\kappa}_{**} - \boldsymbol{\kappa}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \boldsymbol{\kappa}_*\right) \quad (\text{B.10})$$

where the identity presented in [100, (A.6) pg. 200] was used. The expansion for the multivariate case is straightforward. Thus, for new input matrix \mathbf{X}_* the predictive distribution of \mathbf{f}_* is given by

$$\mathbf{f}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_* \sim \mathcal{N}\left(\mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}, \mathbf{K}_{**} - \mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{K}_*\right) \quad (\text{B.11})$$

with $[\mathbf{K}_*]_{ij} = \boldsymbol{\varphi}(\mathbf{x}_i)^\top \boldsymbol{\varphi}(\mathbf{x}_j)$ and $[\mathbf{K}_{**}]_{ij} = \boldsymbol{\varphi}(\mathbf{x}_{*i})^\top \boldsymbol{\varphi}(\mathbf{x}_{*j})$.

Finally, the minimum mean squared error (MMSE) estimator can be

¹Note that Equation (B.9) is in fact the conditional distribution of \mathbf{z} given the data \mathbf{X} and \mathbf{x}_* . In a more rigorous notation one would write $\mathbf{z} | \mathbf{X}, \mathbf{x}_*$.

obtained by taking the mean of the distribution (B.11)

$$\hat{\mathbf{f}}_*^{MMSE} = \mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}. \quad (\text{B.12})$$

The Bayesian formulation also allows us to estimate the parameters $\boldsymbol{\theta} = [\sigma_n^2, \sigma_w^2]^\top$ intuitively by maximizing the *log marginal likelihood* of (B.9), that is, maximizing the logarithm of the distribution $\mathbf{y}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma_n^2 \mathbf{I})$. Thus, one can formulate the following optimization problem

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) \quad (\text{B.13})$$

where

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{y}^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K} + \sigma_n^2 \mathbf{I}| - \frac{N}{2} \log(2\pi).$$

However, problem (B.13) is not convex and may contain local minima.

We emphasize here that in the approach presented in this section, data are explicitly mapped by the function $\boldsymbol{\phi}$ into a M -dimensional space. It is also important to note that the solutions for both the scalar and multivariate cases are always written in function of internal products $\boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\phi}(\mathbf{x}')$. Thus, $\mathbf{K} = \sigma_w^2 \boldsymbol{\Phi} \boldsymbol{\Phi}^\top$ can be seen as the Gram matrix with kernel $\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j) = \sigma_w^2 \boldsymbol{\phi}(\mathbf{x}_i)^\top \boldsymbol{\phi}(\mathbf{x}_j)$. In Section A.6, we used kernel functions to implicitly map the data into a high (even infinite) dimensional feature space and compute inner products of the mapped data by evaluation a real function in the input space. Such approach added great flexibility to the solution. Next, we will demonstrate how to consider similar strategies but within a Bayesian formulation.

B.2 GAUSSIAN PROCESS REGRESSION

In Section A.6 we converted a linear problem into a nonlinear one by working directly in a functional space generated by kernels, i.e., the RKHS. Following the same idea, we will transform the Bayesian linear regression approach presented above in a kernelized nonlinear regression solution by working directly in the functional space. For this, distributions will be considered directly over functions belonging to the RKHS.

Consider again the same problem presented in Section A.6, for which a input-output sequence $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, with $\mathbf{x} \in \mathbb{R}^D$ and scalar output y , is modeled by the following nonlinear relation $y_i = \boldsymbol{\psi}(\mathbf{x}_i) + n_i$, where n_i is zero mean WGN with power σ_n^2 and independent of \mathbf{x}_i and $\boldsymbol{\psi} \in \mathcal{H}$ is a function belong-

ing to the RKHS \mathcal{H} .

Lets assume a Gaussian prior distribution for $\boldsymbol{\psi}$,

$$\boldsymbol{\psi}|\mathbf{x} \sim \mathcal{N}(0, \boldsymbol{\kappa}), \quad (\text{B.14})$$

where $\boldsymbol{\kappa} = \boldsymbol{\kappa}(\mathbf{x}, \mathbf{x})$ and $\boldsymbol{\kappa}(\cdot, \cdot)$ is any kernel function. If the multivariate (vector) case is considered, then $\mathbf{y} = \boldsymbol{\psi} + \mathbf{n}$, and the prior distribution becomes

$$\boldsymbol{\psi}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad (\text{B.15})$$

with $\mathbf{K}_{ij} = \boldsymbol{\kappa}(\mathbf{x}_i, \mathbf{x}_j)$, and

$$\mathbf{y}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma_n^2 \mathbf{I}). \quad (\text{B.16})$$

Following the same reasoning used in Section B.1, lets consider a new input \mathbf{x}_* and assume $\boldsymbol{\psi}_*$ as its “predicted” output. Then, the joint distribution for $\mathbf{z} = [\mathbf{y}^\top, \boldsymbol{\psi}_*^\top]^\top$ is given by

$$\begin{bmatrix} \mathbf{y} \\ \boldsymbol{\psi}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \boldsymbol{\kappa}_* \\ \boldsymbol{\kappa}_*^\top & \boldsymbol{\kappa}_{**} \end{bmatrix}\right) \quad (\text{B.17})$$

where $\boldsymbol{\kappa}_* = [\boldsymbol{\kappa}(\mathbf{x}_1, \mathbf{x}_*), \dots, \boldsymbol{\kappa}(\mathbf{x}_N, \mathbf{x}_*)]^\top$ e $\boldsymbol{\kappa}_{**} = \boldsymbol{\kappa}(\mathbf{x}_*, \mathbf{x}_*)$. Thus, the conditional, or predictive, distribution can be written as

$$\boldsymbol{\psi}_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_* \sim \mathcal{N}\left(\boldsymbol{\kappa}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}, \boldsymbol{\kappa}_{**} - \boldsymbol{\kappa}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \boldsymbol{\kappa}_*\right) \quad (\text{B.18})$$

for the univariate case, and for multiple inputs \mathbf{X}_* as

$$\boldsymbol{\psi}_*|\mathbf{y}, \mathbf{X}, \mathbf{X}_* \sim \mathcal{N}\left(\mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}, \mathbf{K}_{**} - \mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{K}_*\right). \quad (\text{B.19})$$

Finally, the MMSE estimator can be obtained by considering the expected value of the conditional distribution (B.19)

$$\hat{\boldsymbol{\psi}}_*^{MMSE} = \mathbf{K}_*^\top [\mathbf{K} + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}. \quad (\text{B.20})$$

The kernel hyperparameters $\boldsymbol{\theta}_k$ and the noise power σ_n^2 can be estimated by maximizing the log marginal likelihood. Thus, for $\boldsymbol{\theta} = [\boldsymbol{\theta}_k^\top, \sigma_n^2]^\top$, we have

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) \quad (\text{B.21})$$

where

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2}\mathbf{y}^\top [\mathbf{K} + \sigma_n^2\mathbf{I}]^{-1} \mathbf{y} - \frac{1}{2}\log|\mathbf{K} + \sigma_n^2\mathbf{I}| - \frac{N}{2}\log(2\pi).$$

REFERENCES

- [1] SWAYZE, G. et al. Ground-truthing aviris mineral mapping at cuprite, nevada. In: *Summaries of the Third Annual JPL Airborne Geoscience Workshop*. [S.l.: s.n.], 1992. v. 1, p. 47–49.
- [2] ALTMANN, Y. et al. A robust test for nonlinear mixture detection in hyperspectral images. In: *2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. [S.l.: s.n.], 2013.
- [3] IMBIRIBA, T. et al. Nonparametric detection of nonlinearly mixed pixels and endmember estimation in hyperspectral images. *IEEE Transactions on Image Processing*, v. 25, n. 3, p. 1136–1151, March 2016. ISSN 1057-7149.
- [4] LANDGREBE, D. The evolution of landsat data analysis. *Photogrammetric Engineering and Remote Sensing*, LXIII, n. 7, p. 859–867, 1997.
- [5] LANDGREBE, D. Hyperspectral image data analysis. *IEEE Signal Processing Magazine*, v. 19, n. 1, p. 17–28, 2002.
- [6] KESHAVA, N.; MUSTARD, J. Spectral unmixing. *IEEE Signal Processing Magazine*, v. 19, n. 1, p. 44–57, 2002.
- [7] BIOUCAS-DIAS, J. M. et al. Hyperspectral remote sensing data analysis and future challenges. *IEEE Geoscience and Remote Sensing Magazine*, v. 1, n. 2, p. 6–36, 2013.
- [8] DOBIGEON, N. et al. Nonlinear unmixing of hyperspectral images: Models and algorithms. *IEEE Signal Processing Magazine*, v. 31, n. 1, p. 82–94, Jan 2014.
- [9] SONG, C. Sciencedirect.com - remote sensing of environment - spectral mixture analysis for subpixel vegetation fractions in the urban environment: How to incorporate endmember variability? *Remote Sensing of Environment*, 2005.
- [10] ZARE, A.; HO, K. C. Endmember variability in hyperspectral analysis: Addressing spectral variability during spectral unmixing. *IEEE Signal Processing Magazine*, v. 31, p. 95, January 2014.

- [11] KEEF, J.; ARIZONA, T. U. of. Hyper-spectral sensor calibration extrapolated from multi-spectral measurements. the university of arizona. *Optical Sciences*, 2008.
- [12] IMBIRIBA, T. et al. Detection of nonlinear mixtures using gaussian processes: Application to hyperspectral imaging. In: *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. [S.l.: s.n.], 2014. p. 7949–7953.
- [13] CHEN, J.; RICHARD, C.; HONEINE, P. Nonlinear unmixing of hyperspectral data based on a linear-mixture/nonlinear-fluctuation model. *IEEE Transactions on Signal Processing*, v. 61, p. 480–492, Jan 2013.
- [14] CHEN, J.; RICHARD, C.; HONEINE, P. Nonlinear unmixing of hyperspectral images with multi-kernel learning. In: *Proc. IEEE WHISPERS*. [S.l.: s.n.], 2012. p. 1–4.
- [15] CHEN, J.; RICHARD, C.; HONEINE, P. Estimating abundance fractions of materials in hyperspectral images by fitting a post-nonlinear mixing model. In: *Proc. IEEE WHISPERS*. [S.l.: s.n.], 2013. p. 1–4.
- [16] CHEN, J.; RICHARD, C.; HONEINE, P. Nonlinear estimation of material abundances in hyperspectral images with ℓ_1 -norm spatial regularization. *IEEE Transactions on Geoscience and Remote Sensing*, v. 52, n. 5, p. 2654–2665, May 2014. ISSN 0196-2892.
- [17] ALTMANN, Y. et al. Supervised nonlinear spectral unmixing using a postnonlinear mixing model for hyperspectral imagery. *IEEE Transactions on Image Processing*, v. 21, n. 6, p. 3017–3025, June 2012.
- [18] HALIMI, A. et al. Nonlinear Unmixing of Hyperspectral Images Using a Generalized Bilinear Model. *IEEE Transactions on Geoscience and Remote Sensing*, v. 49, n. 11, p. 4153–4162, Nov. 2011.
- [19] BROADWATER, J. et al. Kernel fully constrained least squares abundance estimates. In: *2007 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. [S.l.: s.n.], 2007. p. 4041–4044.

- [20] BROADWATER, J.; BANERJEE, A. A comparison of kernel functions for intimate mixture models. In: *2009 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. [S.l.: s.n.], 2009. p. 1–4.
- [21] ALTMANN, Y. et al. Residual component analysis of hyperspectral images—application to joint nonlinear unmixing and nonlinearity detection. *IEEE Transactions on Image Processing*, v. 23, p. 2148–2158, May 2014.
- [22] BOARDMAN, J. Automatic spectral unmixing of AVIRIS data using convex geometry concepts. In: *Proc. AVIRIS workshop*. [S.l.: s.n.], 1993. v. 1, p. 11–14. ISBN 0780314972.
- [23] NASCIMENTO, J. M. P.; BIOUCAS-DIAS, J. M. Vertex Component Analysis: A fast algorithm to unmix hyperspectral data. *IEEE Transactions on Geoscience and Remote Sensing*, v. 43, n. 4, p. 898–910, April 2005. ISSN 0196-2892.
- [24] CHANG, C. et al. A new growing method for simplex-based endmember extraction algorithm. *IEEE Transactions on Geoscience and Remote Sensing*, IEEE, v. 44, n. 10, p. 2804–2819, 2006. ISSN 0196-2892.
- [25] CHAN, T.-H. et al. A convex analysis-based minimum-volume enclosing simplex algorithm for hyperspectral unmixing. *Signal Processing, IEEE Transactions on*, IEEE, v. 57, n. 11, p. 4418–4432, Nov 2009.
- [26] CHAN, T.-H. et al. A simplex volume maximization framework for hyperspectral endmember extraction. *IEEE Transactions on Geoscience and Remote Sensing*, v. 49, n. 11, p. 4177–4193, May 2011.
- [27] HEYLEN, R.; BURAZEROVIC, D.; SCHEUNDERS, P. Non-linear spectral unmixing by geodesic simplex volume maximization. *IEEE Journal of Selected Topics in Signal Processing*, v. 5, n. 3, p. 534–542, 2011.
- [28] ALTMANN, Y. et al. Nonlinear spectral unmixing of hyperspectral images using gaussian processes. *IEEE Transactions on Signal Processing*, v. 61, p. 2442–2453, May 2013.

- [29] CHANG, C.-I. *Hyperspectral data processing - Algorithm design and analysis*. [S.l.]: Wiley, 2014.
- [30] IMBIRIBA, T. et al. Band selection in RKHS for fast nonlinear unmixing of hyperspectral images. In: *2015 23rd European Signal Processing Conference (EUSIPCO)*. [S.l.: s.n.], 2015. p. 1651–1655.
- [31] Imbiriba, T.; Bermudez, J. C. M.; Richard, C. Technical Report: Band selection for nonlinear unmixing of hyperspectral images as a maximal clique problem. *ArXiv e-prints*, mar. 2016.
- [32] PARDALOS, P. M.; XUE, J. The maximum clique problem. *Journal of global Optimization*, Springer, v. 4, n. 3, p. 301–328, 1994.
- [33] WU, Q.; HAO, J.-K. A review on algorithms for maximum clique problems. *European Journal of Operational Research*, Elsevier, v. 242, n. 3, p. 693–709, 2015.
- [34] BOREL, C. C.; GERSTL, S. A. W. Nonlinear spectral mixing models for vegetative and soil surfaces. *Remote Sensing of Environment*, v. 47, n. 3, p. 403–416, Jan 1994.
- [35] HAPKE, B. Bidirectional reflectance spectroscopy, 1, Theory. *Journal of Geophysical Research*, v. 86, n. B4, p. 3039–3054, 1981.
- [36] HAPKE, B. *Theory of Reflectance and Emittance Spectroscopy*. [S.l.]: Cambridge University Press, 1993.
- [37] SOMERS, B. et al. Nonlinear hyperspectral mixture analysis for tree cover estimates in orchards. *Remote Sensing of Environment*, v. 113, n. 6, p. 1183–1193, February 2009.
- [38] FAN, W. et al. Comparative study between a new nonlinear model and common linear model for analysing laboratory simulated-forest hyperspectral data. *International Journal of Remote Sensing*, v. 30, n. 11, p. 2951–2962, 2009.
- [39] NASCIMENTO, J. M. P.; BIOUCAS-DIAS, J. M. Nonlinear mixture model for hyperspectral unmixing. In: INTERNATIONAL SOCIETY FOR OPTICS AND PHOTONICS. *SPIE Europe Remote Sensing*. [S.l.], 2009. p. 74770I–74770I.
- [40] HALIMI, A. et al. Unmixing hyperspectral images using the generalized bilinear model. In: *2011 IEEE International Geoscience*

- and Remote Sensing Symposium (IGARSS)*. [S.l.]: IEEE, 2011. p. 1886–1889.
- [41] JUTTEN, C.; KARHUNEN, J. Advances in nonlinear blind source separation. In: *Proc. International Symposium on Independent Component Analysis and Blind Signal Separation (ICA)*. [S.l.: s.n.], 2003. p. 245–256.
- [42] ALTMANN, Y.; DOBIGEON, N.; TOURNERET, J.-Y. Nonlinearity detection in hyperspectral images using a polynomial post-nonlinear mixing model. *IEEE Transactions on Image Processing*, v. 22, n. 4, p. 1267–1276, 2013.
- [43] ALTMANN, Y. et al. Supervised nonlinear spectral unmixing using a polynomial post nonlinear model for hyperspectral imagery. In: *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. [S.l.: s.n.], 2011. p. 1009–1012.
- [44] ALTMANN, Y. et al. A polynomial post nonlinear model for hyperspectral image unmixing. In: *2011 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. Vancouver, Canada: [s.n.], 2011.
- [45] DOBIGEON, N.; TOURNERET, J.-Y.; CHANG, C. I. Semi-Supervised Linear Spectral Unmixing Using a Hierarchical Bayesian Model for Hyperspectral Imagery. *IEEE Transactions on Signal Processing*, v. 56, n. 7, p. 2684–2695, 2008.
- [46] NASH, D. B.; CONEL, J. E. Spectral reflectance systematics for mixtures of powdered hypersthene, labradorite, and ilmenite. *Journal of Geophysical Research*, Wiley Online Library, v. 79, n. 11, p. 1615–1621, 1974.
- [47] MUSTARD, J. F.; PIETERS, C. M. Photometric phase functions of common geologic minerals and applications to quantitative analysis of mineral mixture reflectance spectra. *Journal of Geophysical Research: Solid Earth (1978–2012)*, Wiley Online Library, v. 94, n. B10, p. 13619–13634, 1989.
- [48] DRAINE, B. T. The discrete-dipole approximation and its application to interstellar graphite grains. *The Astrophysical Journal*, v. 333, p. 848–872, 1988.

- [49] SHKURATOV, Y. et al. A model of spectral albedo of particulate surfaces: Implications for optical properties of the moon. *Icarus*, Elsevier, v. 137, n. 2, p. 235–246, 1999.
- [50] MA, W.-K. et al. A signal processing perspective on hyperspectral unmixing: Insights from remote sensing. *IEEE Signal Processing Magazine*, v. 31, n. 1, p. 67–81, January 2014.
- [51] NASCIMENTO, J. M. P.; DIAS, J. M. B. Does independent component analysis play a role in unmixing hyperspectral data? *Geoscience and Remote Sensing, IEEE Transactions on*, v. 43, p. 175–187, Jan 2005.
- [52] MA, W.-K. et al. Signal and image processing in hyperspectral remote sensing. *IEEE Signal Processing Magazine*, v. 31, n. 1, p. 22–23, January 2014.
- [53] HEYLEN, R.; PARENTE, M.; GADER, P. A review of nonlinear hyperspectral unmixing methods. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, IEEE, v. 7, n. 6, p. 1844–1868, June 2014.
- [54] GILLIS, N.; VAVASIS, S. A. Fast and robust recursive algorithms for separable nonnegative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, v. 36, n. 4, 2014.
- [55] WINTER, M. E. N-FINDR: an algorithm for fast autonomous spectral end-member determination in hyperspectral data. In: *Proc. SPIE Spectrometry V*. [S.l.: s.n.], 1999. v. 3753, p. 266–277.
- [56] NEVILLE, R. A. et al. Automatic endmember extraction from hyperspectral data for mineral exploration. In: *Proc. 21st Canadian Symp. Remote Sens.* [S.l.: s.n.], 1999. p. 21–24.
- [57] BIOUCAS-DIAS, J. M.; PLAZA, A. Hyperspectral unmixing: Geometrical, statistical, and sparse regression-based approaches. In: *Proc. SPIE Image and Signal Processing for Remote Sensing XVI*. [S.l.: s.n.], 2010. v. 7830, p. 78300A1–78300A15.
- [58] STRANG, G. *Linear Algebra and its Applications*. San Diego: Harcourt Brace Jovanovich College Publishers, 1988.
- [59] KAY, S. M. *Fundamentals of Statistical Signal Processing: Estimation Theory*. Upper Saddle River: Prentice Hall, 1993.

- [60] BISHOP, C. M. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. [S.l.]: Springer-Verlag New York, Inc., 2006.
- [61] KESHAHA, N. et al. An algorithm taxonomy for hyperspectral unmixing. In: INTERNATIONAL SOCIETY FOR OPTICAL ENGINEERING (SPIE). [S.l.], 2000.
- [62] HONEINE, P.; RICHARD, C. Geometric Unmixing of Large Hyperspectral Images: A Barycentric Coordinate Approach. *IEEE Transactions on Geoscience and Remote Sensing*, v. 50, n. 6, p. 2185–2195, June 2012.
- [63] HEINZ, D.; CHANG, C.-I.; ALTHOUSE, M. L. G. Fully constrained least-squares based linear unmixing [hyperspectral image classification]. In: *IEEE 1999 International Geoscience and Remote Sensing Symposium, 1999. IGARSS '99 Proceedings*. [S.l.: s.n.], 1999. v. 2, p. 1401–1403.
- [64] CRAIG, M. D. Minimum-volume transforms for remotely sensed data. *IEEE Transactions on Geoscience and Remote Sensing*, v. 32, n. 3, p. 542–552, maio 1994.
- [65] BERMAN, M. et al. Ice: A statistical approach to identifying endmembers in hyperspectral images. *IEEE Transactions on Geoscience and Remote Sensing*, Citeseer, v. 42, n. 10, p. 2085–2095, 2004.
- [66] MIAO, L.; QI, H. Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization. *IEEE Transactions on Geoscience and Remote Sensing*, v. 45, n. 3, p. 765–777, 2007.
- [67] BIOUSCAS-DIAS, J. M. et al. Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches. *Selected Topics in Applied Earth Observations and Remote Sensing, IEEE Journal of*, IEEE, v. 5, n. 2, p. 354–379, 2012.
- [68] ROBERT, C. P.; CASELLA, G. *Monte Carlo statistical methods*. [S.l.]: Springer Verlag, 2004.
- [69] DOBIGEON, N. et al. Joint bayesian endmember extraction and linear unmixing for hyperspectral imagery. *IEEE Transactions on*

Signal Processing, IEEE, v. 57, n. 11, p. 4355–4368, 2009. ISSN 1053-587X.

- [70] DOBIGEON, N. et al. Bayesian separation of spectral sources under non-negativity and full additivity constraints. *Signal Processing*, Elsevier, v. 89, n. 12, p. 2657–2669, 2009.
- [71] DOBIGEON, N.; MOUSSAOUI, S.; COULON, M. Subspace-based Bayesian blind source separation for hyperspectral imagery. In: IEEE (Ed.). *Proc. 3rd IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. [S.l.]: IEEE Press, 2009. p. 372–375.
- [72] HEINZ, D. C.; CHANG, C.-I. Fully constrained least squares linear spectral mixture analysis method for material quantification in hyperspectral imagery. *Geoscience and Remote Sensing, IEEE Transactions on*, v. 39, n. 3, p. 529–545, 2001.
- [73] ALTMANN, Y.; DOBIGEON, N.; TOURNERET, J.-Y. Unsupervised post-nonlinear unmixing of hyperspectral images using a hamiltonian monte carlo algorithm. *Image Processing, IEEE Transactions on*, IEEE, v. 23, n. 6, p. 2663–2675, 2014.
- [74] HEYLEN, R.; SCHEUNDERS, P. Calculation of geodesic distances in nonlinear mixing models: Application to the generalized bilinear model. *Geoscience and Remote Sensing Letters, IEEE, IEEE*, v. 9, n. 4, p. 644–648, 2012.
- [75] GUILFOYLE, K. J.; ALTHOUSE, M. L.; CHANG, C.-I. A quantitative and comparative analysis of linear and nonlinear spectral mixture models using radial basis function neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, v. 39, n. 10, p. 2314–2318, 2001.
- [76] VAPNIK, V. N. *The nature of statistical learning theory*. New York, NY: Springer, 1995.
- [77] SCHÖLKOPF, B.; SMOLA, A. J. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. [S.l.]: The MIT Press, 2001.
- [78] NGUYEN, N. H. et al. Hyperspectral image unmixing using manifold learning methods: derivations and comparative tests. In:

- 2012 *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. [S.l.: s.n.], 2012.
- [79] LICCIARDI, G. A. et al. Unsupervised nonlinear spectral unmixing by means of nlpca applied to hyperspectral imagery. In: *IEEE. 2012 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. [S.l.], 2012. p. 1369–1372.
- [80] RAY, T. W.; MURRAY, B. C. Nonlinear spectral mixing in desert vegetation. *Remote Sensing of Environment*, v. 55, n. 1, p. 59–64, 1996.
- [81] HAN, T.; GOODENOUGH, D. G. Investigation of nonlinearity in hyperspectral imagery using surrogate data methods. *IEEE Transactions on Geoscience and Remote Sensing*, v. 46, p. 2840–2847, Oct 2008.
- [82] THEILER, J. et al. Testing for nonlinearity in time series: the method of surrogate data. *Physica D*, v. 58, p. 77–94, Mar 1992.
- [83] KUGIUMTZIS, D. Test your surrogate data before you test for nonlinearity. *Physical Review E, APS*, v. 60, n. 3, p. 2808, 1999.
- [84] ALTMANN, Y. et al. Residual component analysis of hyperspectral images for joint nonlinear unmixing and nonlinearity detection. In: *IEEE International Conference on Acoustics, Speech and Signal Processing, 2014. ICASSP 2014*. [S.l.: s.n.], 2014. p. 3166–3170.
- [85] CHANG, C.-I.; LIU, K.-H. Progressive band selection of spectral unmixing for hyperspectral imagery. *IEEE Transactions on Geoscience and Remote Sensing*, v. 52, n. 4, p. 2002–2017, 2014.
- [86] DU, Q.; YANG, H. Similarity-based unsupervised band selection for hyperspectral image analysis. *IEEE Geoscience and Remote Sensing Letters*, IEEE, v. 5, n. 4, p. 564–568, 2008.
- [87] ESTÉVEZ, P. et al. Normalized mutual information feature selection. *IEEE Transactions on Neural Networks*, IEEE, v. 20, n. 2, p. 189–201, 2009.
- [88] MARTÍNEZ-USÓ, A. et al. Clustering-based hyperspectral band selection using information measures. *IEEE Transactions on Geoscience and Remote Sensing*, IEEE, v. 45, n. 12, p. 4158–4171, 2007.

- [89] FENG, J. et al. Hyperspectral band selection based on trivariate mutual information and clonal selection. *IEEE Transactions on Geoscience and Remote Sensing*, v. 52, n. 7, p. 4092–4105, 2014.
- [90] FENG, J. et al. Mutual-information-based semi-supervised hyperspectral band selection with high discrimination, high information, and low redundancy. *IEEE Transactions on Geoscience and Remote Sensing*, v. 53, p. 2956–2969, 2015.
- [91] CHANG, C.-I.; WANG, S. Constrained band selection for hyperspectral imagery. *IEEE Transactions on Geoscience and Remote Sensing*, v. 44, n. 6, p. 1575–1585, 2006.
- [92] BIOUCAS-DIAS, J. M.; NASCIMENTO, J. M. P. Hyperspectral Subspace Identification. *IEEE Transactions on Geoscience and Remote Sensing*, v. 46, n. 8, p. 2435–2445, 2008.
- [93] BIOUCAS-DIAS, J. M.; PLAZA, A. An overview on hyperspectral unmixing: Geometrical, statistical, and sparse regression based approaches. *2011 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, p. 1135–1138, 2011.
- [94] MERCER, J. Functions of positive and negative type and their connection with the theory of integral equations. *Philos. Trans. Roy. Soc. London Ser. A*, v. 209, p. 415–446, 1909.
- [95] MOORE, E. H. On properly positive hermitian matrices. *Bull. American Mathematical Society*, v. 23, p. 59, 1916.
- [96] ARONSZAJN, N. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, v. 68, 1950.
- [97] ARONSZAJN, N. La théorie des noyaux reproduisants et ses applications première partie. *Mathematical Proceedings of the Cambridge Philosophical Society*, v. 39, p. 133–153, 10 1943. ISSN 1469-8064.
- [98] KREYSZIG, E. *Introductory functional analysis with applications*. [S.l.]: wiley New York, 1989. v. 81.
- [99] STEINWART, I.; CHRISTMANN, A. *Support vector machines*. [S.l.]: Springer, 2008.

- [100] RASMUSSEN, C. E.; WILLIAMS, C. K. I. *Gaussian Processes for Machine Learning*. [S.l.]: The MIT Press, 2006.
- [101] RASMUSSEN, C. E.; NICKISCH, H. Gaussian processes for machine learning (gpml) toolbox. *The Journal of Machine Learning Research*, JMLR. org, v. 11, p. 3011–3015, 2010.
- [102] SUYKENS, J. A. K. et al. *Least Squares Support Vector Machines*. Singapore: World Scientific, 2002.
- [103] KADRI, H. et al. *General Framework for Nonlinear Functional Regression with Reproducing Kernel Hilbert Spaces*. [S.l.], 2009.
- [104] CHEN, J. et al. Nonlinear unmixing of hyperspectral data with partially linear least-squares support vector regression. In: *2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. [S.l.: s.n.], 2013. p. 2174–2178.
- [105] BOYD, S.; VANDENBERGHE, L. *Convex Optimization*. Cambridge: University Press, 2004.
- [106] KAY, S. M. *Fundamentals of statistical signal processing: detection theory*. [S.l.]: Prentice-hall, 1998.
- [107] DUDA, R. O.; HART, P. E.; STORK, D. G. *Pattern classification*. [S.l.]: John Wiley & Sons, 2012.
- [108] THEODORIDIS, S.; KOUTROUMBAS, K. *Pattern recognition*. [S.l.]: Academic press, 2003.
- [109] JOHNSON, N. J.; KOTZ, S.; BALAKRISHNAN, N. *Continuous Univariate Distributions*. [S.l.]: Wiley-Interscience, 1994. v. 2.
- [110] JOHNSON, N. J.; KOTZ, S.; BALAKRISHNAN, N. *Continuous Univariate Distributions*. [S.l.]: Wiley-Interscience, 1994. v. 1.
- [111] INC., R. R. S. *ENVI User's guide Version 4.0*. Boulder, CO 80301 USA: [s.n.], 2013.
- [112] SÁ, J. P. M. D. *Applied statistics using SPSS, STATISTICA and MATLAB*. [S.l.]: Springer, 2003.
- [113] CHANG, C.-I.; DU, Q. Estimation of number of spectrally distinct signal sources in hyperspectral imagery. *IEEE Transactions on Geoscience and Remote Sensing*, v. 42, n. 3, p. 608–619, 2004.

- [114] BIOUCAS-DIAS, J. M.; NASCIMENTO, J. M. P. Hyperspectral subspace identification. *IEEE Transactions on Geoscience and Remote Sensing*, v. 46, n. 8, p. 2435–2445, 2008.
- [115] HALIMI, A. et al. Estimating the intrinsic dimension of hyperspectral images using a noise-whitened eigengap approach. *IEEE Transactions on Geoscience and Remote Sensing*, PP, n. 99, p. 1–11, 2016. ISSN 0196-2892.
- [116] AMBIKAPATHI, A. et al. Chance-constrained robust minimum-volume enclosing simplex algorithm for hyperspectral unmixing. *IEEE Transactions on Geoscience and Remote Sensing*, IEEE, v. 49, n. 11, p. 4194–4209, 2011.
- [117] MIAO, L.; QI, H. Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization. *IEEE Transactions on Geoscience and Remote Sensing*, IEEE, v. 45, n. 3, p. 765–777, 2007.
- [118] HENDRIX, E. M. T. et al. On the minimum volume simplex enclosure problem for estimating a linear mixing model. *Journal of Global Optimization*, Springer, p. 1–14, 2013.
- [119] CLARK, R. N. et al. Imaging spectroscopy: Earth and planetary remote sensing with the usgs tetracorder and expert systems. *Journal of Geophysical Research: Planets (1991–2012)*, Wiley Online Library, v. 108, n. E12, 2003.
- [120] DOPIDO, I. et al. Unmixing prior to supervised classification of remotely sensed hyperspectral images. *IEEE Geoscience and Remote Sensing Letters*, v. 8, n. 4, p. 760–764, 2011.
- [121] GENG, X. et al. Optimizing the endmembers using volume invariant constrained model. *IEEE Transactions on Image Processing*, v. 24, n. 11, p. 3441–3449, Nov 2015. ISSN 1057-7149.
- [122] AMMANOUIL, R. et al. Blind and fully constrained unmixing of hyperspectral images. *IEEE Transactions on Image Processing*, v. 23, n. 12, p. 5510–5518, 2014.
- [123] CHANG, C.-I.; LIU, K.-H. Progressive band selection of spectral unmixing for hyperspectral imagery. *IEEE Transactions on Geoscience and Remote Sensing*, v. 52, n. 4, p. 2002–2017, 2014.

- [124] LI, C.-M.; QUAN, Z. An efficient branch-and-bound algorithm based on maxsat for the maximum clique problem. In: *AAAI*. [S.l.: s.n.], 2010. v. 10, p. 128–133.
- [125] RICHARD, C.; BERMUDEZ, J. C. M.; HONEINE, P. Online prediction of time series data with kernels. *IEEE Transactions on Signal Processing*, v. 57, n. 3, p. 1058–1067, March 2009.
- [126] TZORTZIS, G. F.; LIKAS, A. C. The global kernel k-means algorithm for clustering in feature space. *IEEE Transactions on Neural Networks*, v. 20, p. 1181–1194, 2009.
- [127] TROPP, J. A. Greed is good: Algorithmic results for sparse approximation. *IEEE Transactions on Information Theory*, IEEE, v. 50, n. 10, p. 2231–2242, 2004.
- [128] MALLAT, S. G.; ZHANG, Z. Matching pursuits with time-frequency dictionaries. *IEEE Transactions on Signal Processing*, IEEE, v. 41, n. 12, p. 3397–3415, 1993.
- [129] DONOHO, D. L.; HUO, X. Uncertainty principles and ideal atomic decomposition. *IEEE Transactions on Information Theory*, IEEE, v. 47, n. 7, p. 2845–2862, 2001.
- [130] ELAD, M.; BRUCKSTEIN, A. M. A generalized uncertainty principle and sparse representation in pairs of bases. *IEEE Transactions on Information Theory*, IEEE, v. 48, n. 9, p. 2558–2567, 2002.
- [131] DONOHO, D. L.; ELAD, M. Optimally sparse representation in general (nonorthogonal) dictionaries via l_1 minimization. *Proceedings of the National Academy of Sciences*, National Acad Sciences, v. 100, n. 5, p. 2197–2202, 2003.
- [132] KARP, R. M. *Reducibility among combinatorial problems*. [S.l.]: Springer, 1972.
- [133] CARRAGHAN, R.; PARDALOS, P. M. An exact algorithm for the maximum clique problem. *Operations Research Letters*, Elsevier, v. 9, n. 6, p. 375–382, 1990.
- [134] LI, C.-M.; FANG, Z.; XU, K. Combining maxsat reasoning and incremental upper bound for the maximum clique problem. In: IEEE.

2013 IEEE 25th International Conference on Tools with Artificial Intelligence (ICTAI). [S.l.], 2013. p. 939–946.

- [135] ÖSTERGÅRD, P. R. J. A fast algorithm for the maximum clique problem. *Discrete Applied Mathematics*, Elsevier, v. 120, n. 1, p. 197–207, 2002.
- [136] BIGGS, N. Some heuristics for graph coloring. *Graph Colourings*, Longman, New York, p. 87–96, 1990.
- [137] FAHLE, T. Simple and fast: Improving a branch-and-bound algorithm for maximum clique. In: *Algorithms–ESA 2002*. [S.l.]: Springer, 2002. p. 485–498.
- [138] MASLOV, E.; BATSYN, M.; PARDALOS, P. M. Speeding up branch and bound algorithms for solving the maximum clique problem. *Journal of Global Optimization*, Springer, v. 59, n. 1, p. 1–21, 2014.
- [139] MANNING, C. D.; RAGHAVAN, P.; SCHÜTZE, H. *Introduction to information retrieval*. [S.l.]: Cambridge university press, Cambridge, 2008. v. 1.
- [140] BOYD, S.; VANDENBERGHE, L. *Convex optimization*. [S.l.]: Cambridge university press, 2009.
- [141] BERLINET, A.; THOMAS-AGNAN, C. *Reproducing kernel Hilbert spaces in probability and statistics*. [S.l.]: Springer, 2004. v. 3.
- [142] JAHN, J. *Introduction to the Theory of Nonlinear Optimization*. Third. [S.l.]: Springer, 2007.
- [143] COLEMAN, R. *Calculus on Normed Vector Spaces*. [S.l.]: Springer, 2010.
- [144] WAHBA, G. *Spline models for observational data*. [S.l.]: Siam, 1990. v. 59.
- [145] SCHÖLKOPF, B.; HERBRICH, R.; SMOLA, A. J. A generalized representer theorem. In: SPRINGER. *Computational learning theory*. [S.l.], 2001. p. 416–426.
- [146] KAY, S. M. *Intuitive Probability and Random Processes using MATLAB®*. [S.l.]: Springer, 2006.

- [147] PAPOULIS, A.; PILLAI, S. U. *Probability, Random Variables and Stochastic Processes*. Fourth. [S.l.]: Mc Graw Hill, 2006.