

Prepared by Will Wright

May 20, 2019

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Introduction and contributions . . . . .	3
1.2	Notation . . . . .	4
<b>2</b>	<b>Noisy phase retrieval: a brief survey of applications and methods</b>	<b>8</b>
2.1	Mathematical model . . . . .	8
2.2	Applications and experimental models . . . . .	9
2.3	Survey of noisy phase retrieval methods . . . . .	11
2.3.1	Alternating projection methods . . . . .	12
2.3.2	Structured optimization methods . . . . .	14
2.3.3	Unstructured optimization methods . . . . .	16
2.4	Justification for optimizing the PLGD model . . . . .	22
<b>3</b>	<b>Gauge duality theory and the PhaseLift gauge dual model</b>	<b>27</b>
3.1	Introduction . . . . .	27
3.2	Definitions and primal-dual pairs . . . . .	28
3.3	Theory for linearly and nonlinearly constrained models . . . . .	32
3.4	Optimality conditions and primal recovery for the PLGD model . . . . .	40
<b>4</b>	<b>A first-order method for the PLGD model</b>	<b>45</b>
4.1	Introduction . . . . .	45
4.2	Algorithm and implementation details . . . . .	45
4.3	Efficient recovery for noiseless phase retrieval . . . . .	51
<b>5</b>	<b>New termination conditions for the first-order PLGD algorithm</b>	<b>53</b>
5.1	Introduction . . . . .	53
5.2	Experimental models and potential residuals . . . . .	53
5.3	Stagnation of the first-order PLGD algorithm . . . . .	56
5.4	New termination conditions . . . . .	61

<b>6</b>	<b>The evolving matrix eigenvalue problem</b>	<b>68</b>
6.1	Introduction . . . . .	68
6.2	The EMEP: computational costs and spectral properties . . . . .	68
6.3	The implicitly restarted Arnoldi method . . . . .	74
<b>7</b>	<b>The adaptive parameter method for the EMEP</b>	<b>83</b>
7.1	Introduction . . . . .	83
7.2	The adaptive parameter method for the EMEP . . . . .	83
7.3	Tracking the EMEP spectrum with the adaptive parameter method . . . .	88
7.4	Numerical results for the adaptive parameter method . . . . .	95
<b>8</b>	<b>Conclusion</b>	<b>101</b>
<b>A</b>	<b>Proof of PhaseLift-<math>l_1</math> gauge duality</b>	<b>102</b>
<b>B</b>	<b>Further justification of residuals established in Chapter 5</b>	<b>103</b>

# Chapter 1

## Introduction

### 1.1 Introduction and contributions

Phase retrieval has a wide range of solution methods, yet few exist for handling noisy observations without imposing additional restrictions such as signal sparsity. One recent noisy phase retrieval model which requires no underlying assumptions is the gauge dual of the PhaseLift model (PLGD)

$$\begin{aligned} \min_y \quad & \lambda_1(\mathcal{A}^*y) \\ \text{(PLGD) s.t.} \quad & \langle b, y \rangle - \epsilon \|y\| \geq 1. \end{aligned} \tag{1.1.1}$$

First introduced and analyzed in [FM16], the PLGD model (1.1.1) is based on the PhaseLift method [CESV13] in which a desired signal of  $n$  elements (e.g., pixels) is lifted into the space of  $n \times n$  positive semidefinite matrices, creating a convex, large-scale recovery problem. The PLGD model (1.1.1) maintains the convergence guarantees of the convex PhaseLift model, yet also allows for efficient first-order methods.

To optimize the PLGD model (1.1.1) we use the first-order method proposed in [FM16, Section 4.4], which we restate in Section 4.2 as Algorithm 3. The authors of [FM16] demonstrate that Algorithm 3 is far more efficient than an alternate method for optimizing the PhaseLift model and returns signals with greater accuracy than wflow, another method for phase retrieval (see Section 2.3.3 for details regarding PhaseLift and wflow). Yet the PLGD model (1.1.1) faces two significant challenges for noisy phase retrieval. Computationally, each evaluation of the objective function  $\lambda_1(\mathcal{A}^*y)$  involves a large-scale eigenvalue problem which may require significant runtime for large signals. Additionally, when the phase retrieval problem has nontrivial noise, first-order methods for the PLGD model (1.1.1) such as Algorithm 3 typically stagnate before converging.

This dissertation offers two main contributions. First, we address the algorithmic stagnation of Algorithm 3 for noisy phase retrieval and establish new termination conditions which indicate stagnation of signal recovery progress. Second, we develop a new strategy

for handling the sequence of eigenvalue problems in Algorithm 3 which decreases the computational cost and overall runtime of Algorithm 3 by 50 – 90% for problems with minimal oversampling.

This dissertation is organized in the following manner. Chapter 2 introduces the phase retrieval problem and provides a survey of phase retrieval methods. We close Chapter 2 by demonstrating that Algorithm 3 is generally more accurate for noisy phase retrieval than the wflow method. Chapter 3 presents the gauge duality theory necessary for analyzing and optimizing the PLGD model (1.1.1). Chapter 4 then develops Algorithm 3, a first-order method for the PLGD model (1.1.1), and demonstrates the effectiveness of this algorithm for noiseless phase retrieval.

Chapter 5 demonstrates that Algorithm 3 typically stagnates prior to converging for noisy phase retrieval problems. We then identify the cause of this stagnation and establish new termination conditions for Algorithm 3. Chapter 6 defines the sequence of eigenvalue problems in Algorithm 3 as the *evolving matrix eigenvalue problem* (EMEP). We see that the EMEP is the computational bottleneck of Algorithm 3. The spectrum of these eigenvalue problems evolves in a predictable way, with the algebraically largest eigenvalues clustering for later EMEP iterates. This clustering causes later EMEP iterates to have more difficult eigenvalue problems. We close Chapter 6 with a review of the *implicitly restarted Arnoldi method* (IRAM), a common large-scale eigenvalue method.

Chapter 7 then proceeds to develop an efficient, adaptive strategy (Algorithm 8) for choosing IRAM parameters to handle the EMEP. This chapter shows that Algorithm 8 effectively tracks the clustering of the algebraically largest eigenvalues from earlier to later EMEP iterates, thus selecting more desirable parameters for the IRAM. We close Chapter 7 with a few experiments demonstrating that Algorithm 8 decreases computational cost of the EMEP, thus making Algorithm 3 more efficient for noisy phase retrieval. Chapter 8 concludes this dissertation with a summary of our contributions and suggestions for future work.

## 1.2 Notation

This dissertation uses the following notation. Additional notation and definitions related to gauge duality are stated in Sections 3.2. The  $(i, j)$  entry of a matrix  $A$  is denoted  $[A]_{i,j}$  or  $A(i, j)$ , and the  $i$ -th component of a vector  $a$  is denoted  $a_i$  or  $[a]_i$ . Vector norms are the standard  $p$ -norms, with  $\|\cdot\| \equiv \|\cdot\|_2$ . Matrix norms for  $A \in \mathbb{C}^{m \times n}$  are Schatten  $p$ -norms, which apply the  $p$ -norm to the vector of singular values, i.e.,

$$\|A\|_p = \left( \sum_{i=1}^{\min\{m,n\}} \sigma_i^p(A) \right)^{1/p}. \quad (1.2.1)$$

The special case of  $p = 2$  gives the Frobenius norm

$$\|A\|_F = \left( \sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A) \right)^{1/2}. \quad (1.2.2)$$

Gauge duality is a duality based on multiplicative relations, and thus Schatten norms are essential to gauge duality and developing the PLGD model (1.1.1). In contrast, the EMEP requires measurements with the vector-induced 2-norm. Thus we define the *matrix norm* as

$$\|A\| = \sup_{\|v\|_2=1} \|Av\|_2 = \sigma_{\max}(A). \quad (1.2.3)$$

The standard basis vector is denoted  $e_i$ , where  $[e_i]_i = 1$  and all other components are zero. Given a vector  $d$  in  $\mathbb{R}^n$  or  $\mathbb{C}^n$  with components  $d_1, d_2, \dots, d_n$ , the *diagonal operator* is defined as

$$\text{Diag}(d) = \text{Diag}(d_1, d_2, \dots, d_n)_{ij} = \begin{cases} d_i & \text{if } i = j \\ 0 & \text{else.} \end{cases} \quad (1.2.4)$$

Additionally, if  $A$  is a matrix in  $\mathbb{R}^{n \times n}$  or  $\mathbb{C}^{n \times n}$  then the *diagonal operator* is defined as

$$\text{diag}(A) = \begin{bmatrix} A_{1,1} \\ A_{2,2} \\ \vdots \\ A_{n,n} \end{bmatrix}. \quad (1.2.5)$$

Given  $\mathcal{S}$ , a subset of a finite-dimensional Euclidean space  $\mathcal{X}$ , the *indicator* function of  $\mathcal{S}$  is defined as

$$\delta_{\mathcal{S}}(x) = \begin{cases} 0 & x \in \mathcal{S} \\ +\infty & x \notin \mathcal{S}. \end{cases} \quad (1.2.6)$$

It is easily seen that if  $\mathcal{S}$  is convex, then  $\delta_{\mathcal{S}}$  will be convex. Thus the indicator function is useful for tasks like embedding a domain constraint of an optimization model into the objective function and is used frequently in Chapter 3 for proving gauge duality results.

If  $\mathcal{C}$  is a convex subset of a finite-dimensional Euclidean space, then the *normal cone* of  $\mathcal{C}$  at  $y_0 \in \mathcal{C}$  is defined as

$$N_{\mathcal{C}}(y_0) = \{g \in \mathcal{X} \mid \langle g, y - y_0 \rangle \leq 0 \quad \forall y \in \mathcal{C}\}. \quad (1.2.7)$$

By convention, if  $y_0$  is not in  $\mathcal{C}$ , then  $N_{\mathcal{C}}(y_0)$  is the empty set.

Given a subspace  $S$  of  $\mathbb{R}^n$  or  $\mathbb{C}^n$ , the *orthogonal complement* of  $S$  is defined as

$$S^\perp = \{v \mid \langle v, w \rangle = 0 \text{ for all } w \in S\}. \quad (1.2.8)$$

Given a symmetric (or Hermitian) matrix  $A$  in  $\mathbb{R}^{n \times n}$  (or  $\mathbb{C}^{n \times n}$ ), its eigenvalues are ordered

$$\lambda_1(A) \geq \lambda_2(A) \geq \dots \geq \lambda_n(A), \quad (1.2.9)$$

where  $\lambda_1(A)$  or simply  $\lambda_1$  is the algebraically largest eigenvalue of  $A$ , and  $\lambda_n(A)$  or  $\lambda_n$  is the smallest algebraic eigenvalue. The *spectrum* of  $A$  is the set of all of its eigenvalues  $\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_n\}$ . If  $S$  is a subspace of  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) then  $(\theta, u)$  is a *Ritz pair for  $A$  with respect to  $S$*  if

$$\langle w, (Au - \theta u) \rangle = 0 \quad \forall w \in S. \quad (1.2.10)$$

Likewise,  $\theta$  is a *Ritz value* and  $u$  the corresponding *Ritz vector for  $A$  with respect to  $S$* .

For a pair of matrices  $A, B \in \mathbb{C}^{n \times n}$ , their inner product is that induced by the trace

$$\langle A, B \rangle := \text{tr}(A^* B) = \sum_{i=1}^n \sigma_i(A^* B). \quad (1.2.11)$$

Given  $\mathcal{C}$ , a convex subset of a finite-dimensional Euclidean space  $\mathcal{X}$ , we define projection of  $x \in \mathcal{X}$  onto  $\mathcal{C}$  as  $\Pi_{\mathcal{C}}(x)$ . Since the PLGD (1.1.1) objective function  $f(y) = \lambda_1(\mathcal{A}^* y)$  is not generally differentiable, we consider the subdifferential of  $f$ . Given a convex function  $f : \mathcal{U} \rightarrow \mathbb{R}$  defined on an open, convex subset  $\mathcal{U}$  of a finite-dimensional Euclidean space  $\mathcal{X}$ , the *subdifferential* of  $f$  at  $y_0$  is defined as

$$\partial f(y_0) = \{g \in \mathcal{X} \mid f(y) \geq f(y_0) + \langle g, y - y_0 \rangle \quad \forall y \in \mathcal{U}\}, \quad (1.2.12)$$

and each element of  $\partial f(y_0)$  is a *subgradient* of  $f$ .

Given a linear operator  $\mathcal{A} : \mathcal{X} \rightarrow \mathcal{Y}$  over finite-dimensional Euclidean spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , its *adjoint*  $\mathcal{A}^* : \mathcal{Y} \rightarrow \mathcal{X}$  is defined as the operator which satisfies

$$\langle \mathcal{A}x, y \rangle = \langle x, \mathcal{A}^* y \rangle \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y}. \quad (1.2.13)$$

Since  $\mathcal{X}$  and  $\mathcal{Y}$  are finite and  $\mathcal{A}$  is linear,  $\mathcal{A}$  is also continuous. Thus the Riesz representation theorem guarantees that there will exist a unique linear operator  $\mathcal{A}^*$  [RS80, Section 6.2]. In this dissertation, we will be concerned specifically with linear operators  $\mathcal{A} : \mathcal{H}^n \rightarrow \mathbb{R}^m$ , where  $\mathcal{H}^n$  is the set of  $n \times n$  Hermitian matrices. It is easily shown that all such linear operators  $\mathcal{A}$  will have the form

$$\mathcal{A}(X) = \begin{bmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_m, X \rangle \end{bmatrix}, \quad (1.2.14)$$

where each  $A_i$  is some matrix in  $\mathcal{H}^n$ . In this case, the adjoint of  $\mathcal{A}$  is given by

$$\langle \mathcal{A}(X), y \rangle = \sum_{i=1}^m \langle A_i, X \rangle y_i = \sum_{i=1}^m \langle y_i A_i, X \rangle = \langle X, \sum_{i=1}^m y_i A_i \rangle = \langle X, \mathcal{A}^* y \rangle, \quad (1.2.15)$$

where  $X \in \mathcal{H}^n$  and  $y \in \mathbb{R}^m$ . Thus we have  $\mathcal{A}^*y = \sum_{i=1}^m y_i A_i$ .

The *Gaussian distribution* (or *normal distribution*)  $\mathcal{N}(\mu, \sigma^2)$  is the distribution defined by the probability density function

$$f(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad (1.2.16)$$

where  $\mu$  is the mean and  $\sigma^2$  the variance of the distribution. A real vector has Gaussian distribution  $\nu \sim \mathcal{N}(\mu, \sigma^2)$  if all its elements have Gaussian distribution. Unless otherwise specified, the Gaussian distribution refers specifically to the *standard Gaussian distribution*, where  $\mu = 0$  and  $\sigma^2 = 1$ . The *complex standard Gaussian distribution* is defined by the probability density function

$$f(z) = \frac{1}{\pi} e^{-|z|^2}. \quad (1.2.17)$$



## Chapter 2

# Noisy phase retrieval: a brief survey of applications and methods

### 2.1 Mathematical model

Phase retrieval is the problem of recovering a signal from magnitude-only observations with little or no knowledge of the signal phase. Let  $\bar{x}$  be a one-dimensional signal in  $\mathbb{R}^n$  or  $\mathbb{C}^n$  which has been observed with sensing (or sampling) vectors  $a_i \in \mathbb{C}^m$ , resulting in squared measurements  $|\langle a_i, \bar{x} \rangle|^2 = \bar{b}_i \in \mathbb{R}$  for  $i = 1, \dots, m$ . Also assume the *true observation* vector  $\bar{b} = [\bar{b}_1, \dots, \bar{b}_m]^T \in \mathbb{R}^m$  has been contaminated by possibly nontrivial noise  $\eta = [\eta_1, \dots, \eta_m]^T \in \mathbb{R}^m$ , giving the observation vector  $b = \bar{b} + \eta \in \mathbb{R}^m$ . Then we define the *phase retrieval problem* as

$$\begin{aligned} &\text{find } x \\ &\text{s.t. } |\langle a_i, x \rangle|^2 = b_i = \bar{b}_i + \eta_i \quad 1 \leq i \leq m. \end{aligned} \tag{2.1.1}$$

If  $\eta = 0$ , then (2.1.1) is the noiseless phase retrieval problem (from [Fie82] and [CLS15] among many others). In this dissertation however, we are primarily concerned with nontrivial noise and will refer to (2.1.1) with  $\|\eta\| > 0$  as *noisy phase retrieval*. Additionally, we are concerned with noise  $\eta$  which has a random Gaussian distribution, as discussed in [CESV13] and [FM16].

Each sensing vector  $a_i \in \mathbb{C}^n$  is typically the conjugate of the  $i$ th row of the  $n$ -dimensional discrete Fourier transformation (DFT) matrix  $F$  [BB86, Chapter 11]. This gives the constraint  $|Fx|^2 = b$  (where the square operator is applied element-wise).

Often the number of observations  $m = nL$  is oversampled by a factor of  $L$  to promote uniqueness of a signal solution and convergence of a given algorithm. The domain of the constraint in (2.1.1) is a high-dimensional torus, and thus phase retrieval is inherently nonconvex. When deciding how to handle a particular phase retrieval problem, this nonconvexity presents a key fork in the road in terms of choosing an appropriate algorithm.

In this chapter, we begin by discussing a few typical applications of phase retrieval and briefly discuss the experimental models used for numerical testing (Chapter 7). We then review several methods for noisy phase retrieval. These methods are split into three classes: alternating projection methods (Section 2.3.1), structured optimization methods (Section 2.3.2) which rely on additional assumptions like sparsity of the signal  $x$ , and unstructured optimization methods (Section 2.3.3) which only require an observation vector  $b$  and noise ratio  $\epsilon_{\text{rel}} = \|\eta\|/\|b\|$ . This final class of methods contains the PhaseLift gauge dual (1.1.1) which is central to this dissertation, and therefore we provide greater detail for theoretical guarantees and numerical performance behavior.

## 2.2 Applications and experimental models

Phase retrieval has a broad range of applications across the sciences, many of which fall into the general category of coherent diffraction imaging (CDI) [MCKS99]. This section provides a brief overview of CDI and closes with a description of the phase retrieval experiment models used in [CESV13], [CLS15], [FM16], and this dissertation.

More specific applications of phase retrieval can be found in astronomy [FD87], diffraction and array imaging [BDP<sup>+</sup>07] [CMP10], microscopy [MISE08], optics [Wal63], and x-ray crystallography [Har93], [Mil90] (for a recent benchmark set of crystallography problems, see [ELB17]). For a comprehensive introduction to optical phase retrieval and an overview of recent theory and methods, see the survey [SEC<sup>+</sup>15].

CDI is a method for reconstructing 2- or 3-dimensional nano-structures (e.g., nanotubes, nanocrystals, proteins). Highly coherent waves (e.g., x-rays, electrons, photons) are projected at a given object. The resulting diffraction creates a pattern of intensities which are measured with a detector, resulting in magnitude-only measurements. Figure 2.1 below depicts the CDI observation process.

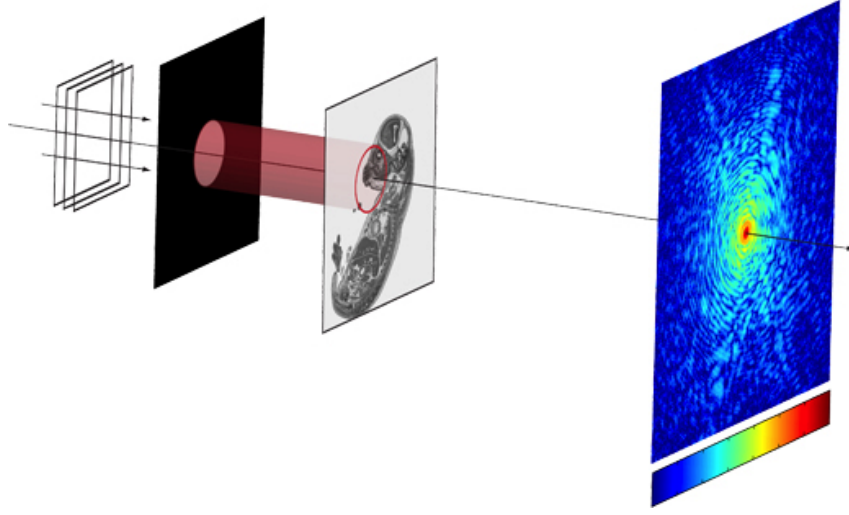


Figure 2.1: Depiction of coherent diffractive imaging. Coherent waves (left) are projected at an image (center) which causes a diffraction pattern that is measured by a detector (right). Image from [GS18b].

Because CDI does not involve optical lenses, there is no optical aberration (blurring or distorting). The resolution depends is instead dependent on the limits of diffraction and dose. Many efficient methods exist for handling low-noise phase retrieval (see Section 2.3.1 for an overview of a few common methods). However, due to the nonconvexity of the constraints in (2.1.1), these low-noise methods are often likely to diverge or converge to a suboptimal local minimum if there is modest noise or an insufficient number of observations. Thus accurate CDI typically requires minimal noise and multiple observations to recover a high-resolution solution.

When developing an experimental model, there are many methods for increasing the number of observations of a signal. Some options include rotating the position of the object, using a spatial light modulator to defocus the observations, and inserting phase plates, or *masks*, in line with the waves (see the survey [DMM<sup>+</sup>11] for a discussion of these methods). Our experimental models use the same masking method described in [CESV13, Section 2], [CLS15, Sections 4.2, 4.3], and [FM16, Section 5.1] (see Section 2.3.3 for an overview of these papers). This masking method involves placing a phase plate with a known structure oriented normal to the projected waves. The phase plate can be placed on either side of the object; in our experiments, the plate lies between the object and the detector. The mask is then shifted and multiple observations are collected.

Mathematically, the application of a phase plate to the phase problem 2.1.1 is equivalent to replacing the sensing operator  $|Fx|^2$  with  $|FC_jx|^2$ , where the matrices  $C_j \in \mathbb{C}^{n \times n}$  are diagonal with standard Gaussian distribution entries  $C_j(i, i) \sim \mathcal{N}(0, 1)$ , representing the

diffraction patterns of the shifted phase plate. This gives the observation constraint

$$\begin{vmatrix} |FC_1x|^2 \\ \vdots \\ |FC_Lx|^2 \end{vmatrix} = b. \quad (2.2.1)$$

In certain cases only a limited number of observations can be collected. For instance, in x-ray imaging, overexposure of the incident waves to the object (living tissue) can be dangerous. Thus the number of observations  $L$  may be relatively small compare to the signal size  $n$ , making signal recovery difficult for nonconvex methods.

To analyze the effectiveness of various phase retrieval methods, we follow the method for generating experimental models established in [CLS15] and extended in [FM16]. Figure 2.2 below depicts a test image of size  $128 \times 128$  and a few particular iterates returned by the method considered in this dissertation (Algorithm 3 in Section 4.2).



Figure 2.2: Results from a first-order method (Algorithm 3 in Section 4.2) applied to a test image with oversampling rate  $L = 8$  and noise ratio 0.30.

## 2.3 Survey of noisy phase retrieval methods

Due to the variety of phase retrieval applications and the difficulty of solving the phase retrieval problem (2.1.1), a wide range of phase retrieval methods have been developed for both noisy and noiseless phase retrieval. This section provides a review of methods for noisy phase retrieval (2.1.1 with  $\|\eta\| > 0$ ), which we separate into three classes: alternating projection methods (Section 2.3.1), structured optimization methods (Section 2.3.2), and unstructured optimization methods (Section 2.3.3). These three classes were chosen to mirror the historical progression of methods (as the alternating projection methods preceded the others by a few decades) and emphasize the uniqueness of the three unstructured methods discussed in Section 2.3.3. For a recent survey of noiseless phase retrieval methods (2.1.1 with  $\eta = 0$ ), see [JEH15].

### 2.3.1 Alternating projection methods

We begin by discussing alternating projection methods for phase retrieval. These methods were established in the 1970s and 1980s as the first strategy for solving (2.1.1) and rely on prior information about the signal, such as support constraints or positivity. Originally referred to as *error reduction algorithms*, these methods were later identified as nonconvex alternating projection methods [LS84]. Each of these methods relies on projecting an iterate between the object (or time) domain and the frequency domain.

We begin with two common alternating projection methods: the Gershberg-Saxton (GS) algorithm [GS72] and the hybrid input-output (HIO) algorithm [Fie82], a GS variant still commonly used in practice. Although these methods are usually not effective for noisy phase retrieval and typically require additional prior information, the basic alternating projection algorithm will be useful for our discussion of the wflow algorithm [CLS15] (Section 2.3.3), which can also be viewed as an alternating projection method. We then discuss recent alternating projection methods for noisy phase retrieval (oversampling smoothness (OSS) [RXC<sup>+</sup>13], error reduction (ER-) HIO and noise robust (NR-) HIO [MWL<sup>+</sup>12]).

Developed in 1972, the GS algorithm was the first alternating projection method for phase retrieval and serves as an algorithmic foundation which later alternating projection methods would modify to improve the likelihood of convergence.

---

#### Algorithm 1 Gershberg-Saxton (GS) algorithm

---

**Input:** Frequency measurement vector  $b \in \mathbb{R}_+^m$ , signal measurement vector  $c \in \mathbb{R}_+^m$ .

**Output:** Approximate solution signal  $x$ .

- 1: *Initialize:* Choose a random signal  $x_0$ , compute DFT  $y_0 = Fx_0$ , set  $res = |y_0|^2 - b$ ,  $k = 0$ .
  - 2: **while**  $\|res\| > \text{tol}$  **do**
  - 3:   Compute DFT of  $x_k$  and residual:  $y_{k+1} = Fx_k$ ,  $res = |y_{k+1}|^2 - b$ .
  - 4:   Impose frequency magnitude constraints:  $[y_{k+1}]_i = \frac{[y_{k+1}]_i}{|[y_{k+1}]_i|} \sqrt{b_i}$  for all  $i = 1, \dots, m$ .
  - 5:   Compute inverse DFT of  $y_{k+1}$ :  $x_{k+1} = F^{-1}y_{k+1}$ .
  - 6:   Impose signal magnitude constraints:  $[x_{k+1}]_i = \frac{[x_{k+1}]_i}{|[x_{k+1}]_i|} \sqrt{c_i}$  for all  $i = 1, \dots, m$ .
  - 7:    $k = k + 1$ .
  - 8: **end while**
  - 9: return  $x \leftarrow x_k$ .
- 

During an iteration of the GS algorithm, knowledge of the object and frequency magnitudes, as well as any additional information, is applied when the iterate reaches that respective domain (steps 6 and 4, respectively). While the GS algorithm is notable as the first algorithmic phase retrieval method, this algorithm is also very likely to converge to a local rather than global minima [JEH15].

In [Fie82], Fienup interpreted the GS algorithm as a nonlinear feedback control system

(Figure 2.3), where the *System* component corresponds to the map  $\tilde{P}_{\text{freq}} = F^{-1}P_{\text{freq}}F$  (combining steps 3-5 of the GS algorithm, where  $P_{\text{freq}}$  is projection onto the frequency constraints in step 4). Since the GS algorithm interprets phase retrieval as an error-reduction problem, the system output  $\tilde{P}_{\text{freq}}(x_k)$  is viewed as the current candidate for the desired signal  $\bar{x}$ , and the object domain constraints are imposed to arrive at a new input  $x_{k+1}$  which is considered the current approximation to the solution signal. By viewing phase retrieval as a nonlinear feedback problem, the input  $x_k$  is no longer treated as a signal approximation, and instead serves as feedback information for the system. Thus  $x_k$  need not satisfy the object domain constraints.

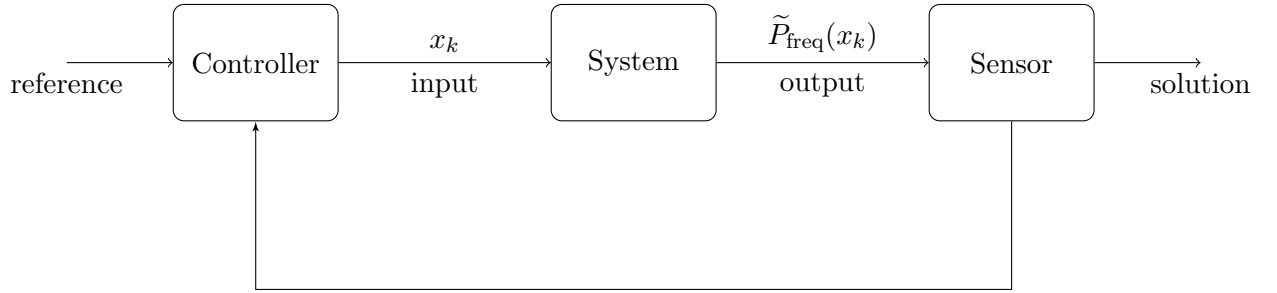


Figure 2.3: Phase retrieval as a nonlinear feedback control system

Fienup observed that a small change in the input will result in an output that is approximately a constant  $\alpha$  times the change in the input, that is

$$\tilde{P}_{\text{freq}}(x + \Delta x) - \tilde{P}_{\text{freq}}(x) \approx \alpha \Delta x. \quad (2.3.1)$$

Thus to force a change of  $\Delta x$  in the output  $\tilde{P}_{\text{freq}}(x)$ , the input would logically be changed by  $\beta \Delta x$ , where  $\beta = 1/\alpha$  from (2.3.1). This observation led Fienup to identify three new potential strategies for selecting an update (the *Controller* in Figure 2.3). The most successful of these methods is the HIO algorithm, which has the following index-wise update replacing step 6 in the GS algorithm:

$$[x_{k+1}]_i = \begin{cases} [\tilde{P}_{\text{freq}}(x_k)]_i & i \notin \mathcal{V} \\ [x_k]_i - \beta [\tilde{P}_{\text{freq}}(x_k)]_i & i \in \mathcal{V}, \end{cases} \quad (2.3.2)$$

where  $\mathcal{V}$  is the set of indices in which the update violates the object domain constraints. This simple corrective step HIO applies to the object domain makes the resulting algorithm much more likely than the GS algorithm to recover a signal from a low-noise observation [JEH15]. Nevertheless, HIO has a tendency to converge to local minima and is not robust to noise. In practice, the user will often select a large set of random initial iterates to initialize HIO and take the best resulting signal.

Recent developments in alternating projection methods for phase retrieval have focused on handling noisy observations. In [RXC<sup>+</sup>13], the authors modify the GS algorithm by

taking the HIO update (2.3.2) and performing an Gaussian smoothing step in the frequency domain on the indices (pixels) violating the support constraint. Their OSS algorithm takes the update  $x_{k+1}$  from HIO step 6 and applies the additional update:

$$[x'_{k+1}]_i = \begin{cases} [x_{k+1}]_i & i \notin \mathcal{V} \\ [F^{-1} [F(x_{k+1}) \cdot \bar{w}(j, \alpha)]]_i & i \in \mathcal{V}, \end{cases} \quad (2.3.3)$$

where  $j = 1, \dots, N$  is the Fourier index,  $\cdot *$  is pointwise multiplication, and  $\bar{w}(j, \alpha)$  is a normalized Gaussian function

$$\bar{w}(j, \alpha) = e^{-\frac{1}{2}(\frac{j}{\alpha})^2}.$$

As  $\alpha \rightarrow \infty$ , the original HIO algorithm is recovered. The authors select  $\alpha$  heuristically, with early  $\alpha = \mathcal{O}(N)$  and later  $\alpha = \mathcal{O}(1/N)$ , which causes OSS to behave similar to HIO during early iterations while damping high frequency violations in later iterations. Their experiments [RXC<sup>+</sup>13, Section 3] apply Poisson noise with noise levels ranging from 0.05 to 0.25 (and measured in the 1-norm, but approximately equivalent to our definition  $\epsilon_{\text{rel}} = \|\eta\|/\|b\|$ ). The results indicate the accuracy and consistency of OSS is superior to HIO and two HIO variants: ER-HIO and NR-HIO of [MWL<sup>+</sup>12].

While these HIO-type methods are computationally efficient and still commonly used in practice [JEH15], they also requires special parameter tuning (e.g.,  $\beta$  in HIO and  $\alpha, \beta$  in OSS) and prior information about the desired signal. Additionally, none of these methods are wholly robust to noise or guaranteed to converge, and require large batches of random initializations in practice to generate an adequate solution. To overcome these deficiencies, recent phase retrieval methods have largely avoided the alternating projection framework, instead casting (2.1.1) as a structured optimization problem.

### 2.3.2 Structured optimization methods

Next we discuss structured optimization methods for noisy phase retrieval. In the past decade, a wide range of methods have been crafted to take advantage of the structures of particular phase retrieval problems. The survey in this section highlights typical methods, their applications, and theoretical guarantees for exact recovery or error bounds.

These methods are grouped based on the structural property they seek to exploit. *Compressive phase retrieval* methods assume sparsity in the signal  $x$  ([SESS11], the GrEedy Sparse PhAse Retrieval (GESPAR) method [SBE14], and thresholded wflow [CLM<sup>+</sup>16]). *Robust phase retrieval* methods assume sparsity in either the observation  $b$  ([Kat17]) or the noise  $\eta$  ([JSL17]). *Supervised phase retrieval* methods require an approximate solution signal  $\hat{x}$  for initialization ([GS18a] and [BR16]).

Many methods have been considered for handling compressive phase retrieval. In [SESS11], the authors consider sparse signals with noisy observations and proceed by lifting the signal and sensing vectors (for details on lifting, see Section 2.3.3). They construct a

rank minimization problem similar to the PhaseLift rank minimization model, but with an additional mixed 1, 2-norm constraint

$$\sum_{i=1}^n \left( \sum_{j=1}^n X_{i,j}^2 \right)^{1/2} \leq \zeta$$

which promotes row-sparsity of the lifted solution matrix  $X$ . The resulting algorithm also includes a thresholding step on the spectrum of  $X$  to enforce low rankness in the solution. The authors provide comparative numerical results indicating the addition of this constraint/thresholding strategy in their optimization method decreases reconstruction error as the noise ratio increases. However, the overall model is nonconvex, the iterations are expensive (each requiring the inversion of a lifted matrix), and no theory is provided to guarantee convergence or signal recovery quality.

The authors of [EM12] prove signal uniqueness and stable recovery (a property stronger than invertibility) results for a variety of sparsity assumptions. If the signal  $x$  is  $k$ -sparse and observation  $b$  is noiseless,  $\mathcal{O}(k \log(n/k))$  observations are necessary for signal uniqueness. If  $x$  is  $k$ -sparse and  $b$  is noisy, then  $\mathcal{O}(k \log(n/k) \log(k))$  observations are necessary to guarantee stable recovery (which gives  $\mathcal{O}(n \log(n))$  observations if  $x$  dense).

The GESPAR method of [SBE14] again assumes  $x$  is  $k$ -sparse and  $b$  is noisy, and constructs an algorithm which maintains an active set  $S$  of indices which converges to the appropriate  $k$  sized set of active indices in the solution signal. This local search method does not require matrix lifting, making it efficient for large-scale phase retrieval. Numerical results [SBE14, Section 5] indicate that as sparsity increases, GESPAR has a higher recovery probability than PhaseLift and a sparse variant of HIO.

A thresholded version of the wflow algorithm (see Section 2.3.3 for wflow) is considered in [CLM<sup>+</sup>16]. Under the assumption that  $x$  is sparse and  $b$  is noisy, the authors develop a thresholded gradient descent algorithm by adding  $\tau \|x\|_1$  to the wflow objective function

$$\min_x \frac{1}{2m} \sum_{i=1}^m (|a_i^* x|^2 - b_i)^2 + \tau \|x\|_1. \quad (2.3.4)$$

They show that this phase retrieval analog to the basis pursuit denoising model [CDS01] has the minimax optimal rate of convergence,  $\mathcal{O}(k \log(n)/m)$ .

Another recent method for compressive phase retrieval considers unstructured noise  $\eta$  in the observation, and interprets the recovery of a  $k$ -sparse signal  $x$  as a covariance maximization problem between the observations  $b_i$  and the sensing values  $|a_i^* x|^2$  over the appropriate  $k$ -dimensional subspace [ZYLX17]. The authors show that only  $\mathcal{O}(k)$  measurements are required for their algorithm to converge within  $\epsilon$  error in a runtime of  $\mathcal{O}(nk \log(1/\epsilon))$ .

A few recent papers have considered sparsity in the observation space rather than the signal space, which can be considered robust phase retrieval. In one case, the paper [JSL17] assumes the noise  $\eta$  is sparse. Whereas the minimization of the 2-norm is optimal for fitting against Gaussian noise, the 1-norm is optimal for fitting against outliers. Thus the



authors suggests minimizing the objective function  $\|\mathcal{A}(xx^*) - b\|_1$  and develop an alternating direction method of multipliers (ADMM) algorithm. Their work provides numerical results demonstrating this method achieves better accuracy than the wflow algorithm when recovering from an observation with Gaussian noise and 10% outliers.

In contrast, the authors of [Kat17] consider phase retrieval when the observation vector  $b$  is itself sparse and noisy. Their algorithm is infact a HIO-type algorithm, with noise suppression in the frequency domain and phase/amplitude filtering in the object domain. Numerical results show this algorithm achieves better accuracy under the given sparsity assumptions than wflow and a filtered version of the GS algorithm.

Rather than assuming sparsity of the signal or observations, the authors of [GS18a] and [BR16] assume the existence of an approximation  $\hat{x}$  to the solution signal  $\bar{x}$ . In this supervised phase retrieval paradigm, the authors define the PhaseMax problem

$$\begin{aligned} \max_x \quad & \operatorname{Re}\langle \hat{x}, x \rangle \\ \text{s.t.} \quad & |\langle a_i, x \rangle|^2 \leq b_i, \quad i = 1, \dots, m \end{aligned} \tag{2.3.5}$$

which has the benefits of being convex and without requiring signal lifting. The dual to this problem is the widely-studied basis pursuit problem, which has several efficient solution techniques as discussed in (2.3.10) below. The authors of [GS18a] prove a relationship between the angle between  $\bar{x}$  and  $\hat{x}$  and the probability of exact recovery, and demonstrate numerically that this probability quickly approaches 1 as oversampling increases.

### 2.3.3 Unstructured optimization methods

The final class of methods we examine are unstructured optimization methods. Unlike the methods discussed in Section 2.3.2, these methods place no assumptions on the signal  $x$ , the sensing vectors  $a_i$ , or the observation  $b$  other than the knowledge of the noise ratio  $\epsilon_{\text{rel}} = \|\eta\|/\|b\|$ .

This section begins with an explanation of matrix lifting for the signal and sensing vectors. Next we discuss the PhaseLift method [CESV13] and the wflow algorithm [CLS15]. Theoretical and numerical results are included to highlight the effectiveness and robustness of these methods (for complete theoretical results, see [CLS15], [SQW16] for wflow and [CL14], [CSV13] for PhaseLift). This discussion of the PhaseLift and wflow methods leads to the PhaseLift gauge dual method [FM16], another unstructured optimization method which is the subject of Chapter 3.

The PhaseLift model was first introduced in [CESV13], with additional theoretical results established in [CSV13]. This model is based on the concept of matrix lifting. First we define the linear operator  $\mathcal{A}$  which allows us to lift the nonlinear phase retrieval observation constraint (2.1.1) into a higher-dimensional linear constraint. If we lift the  $n$ -dimensional signal  $x$  and sensing vectors  $a_i$  into rank-one Hermitian matrices  $X = xx^*$ ,  $A_i = a_i a_i^* \in \mathcal{H}^n$ , then the sensing operator  $\mathcal{A}$  is defined coordinate-wise to satisfy  $[\mathcal{A}(X)]_i = \langle A_i, X \rangle = \operatorname{tr}(a_i a_i^* x x^*) = |\langle a_i, x \rangle|^2$ . Thus  $\mathcal{A} : \mathcal{H}^n \rightarrow \mathbb{R}^m$  and its adjoint

$\mathcal{A}^* : \mathbb{R}^m \rightarrow \mathcal{H}^n$  are defined as

$$\mathcal{A}(xx^*) := \begin{bmatrix} \langle a_1 a_1^*, xx^* \rangle \\ \vdots \\ \langle a_m a_m^*, xx^* \rangle \end{bmatrix} \quad \text{and} \quad \mathcal{A}^* y := \sum_{i=1}^m y_i a_i a_i^*. \quad (2.3.6)$$

In particular, if the experimental model is the one discussed in Section 2.2 then we have

$$\mathcal{A}(xx^*) = \begin{bmatrix} FC_1 x \\ \vdots \\ FC_L x \end{bmatrix}^2 = \text{diag} \left[ \begin{pmatrix} FC_1 \\ \vdots \\ FC_L \end{pmatrix} (xx^*) \begin{pmatrix} FC_1 \\ \vdots \\ FC_L \end{pmatrix}^* \right], \quad (2.3.7)$$

$$\mathcal{A}^* y = \sum_{j=1}^L [FC_j]^* \text{Diag}(y_j) FC_j.$$

The lifted rank-one matrix  $X = xx^* \in \mathcal{H}^n$  and sensing operator  $\mathcal{A}$  are then used to lift the nonlinear phase retrieval constraints  $|\langle a_i, x \rangle|^2 = b_i$  for  $i = 1, \dots, m$  into the linear constraint  $\mathcal{A}(X) = b$ . Thus (2.1.1) is equivalent to the left-most problem in sequence of problems

$$\begin{array}{llll} \text{find } x & & \text{find } X & \\ \text{s.t. } \mathcal{A}(xx^*) = b & \iff & \text{s.t. } \begin{array}{l} \mathcal{A}(X) = b \\ X \succeq 0 \\ \text{rank}(X) = 1 \end{array} & \iff \begin{array}{ll} \min_{X \in \mathcal{H}^n} & \text{rank}(X) \\ \text{s.t.} & \mathcal{A}(X) = b \\ & X \succeq 0. \end{array} \end{array} \quad (2.3.8)$$

If a rank-one solution  $X = xx^*$  exists then this matrix satisfies the lifted model in the middle of (2.3.8), and the left implication holds. Likewise, the rank minimization model at right in (2.3.8) will have a rank-one solution, and the right implication holds.

The resulting rank minimization problem at right of (2.3.8) is NP hard, as it includes the cardinality minimization problem as a special case (see [Nat95], [RFP10]). Thus the next step is to relax the discrete, nonconvex objective function  $\text{rank}(X)$ . To generalize the model for noisy phase retrieval, a norm bound is also applied to the sensing constraint. This gives the semidefinite program which defines the PhaseLift model [CESV13], [CSV13]

$$\begin{array}{ll} \min_{X \in \mathcal{H}^n} & \|X\|_1 = \sum_{i=1}^n \sigma_i(X) \\ \text{(PhaseLift)} \quad \text{s.t.} & \|\mathcal{A}(X) - b\|_2 \leq \epsilon \\ & X \succeq 0. \end{array} \quad (2.3.9)$$

Here the objective function  $\|X\|_1$  refers to the Schatten  $p$ -norm (which is expressed as  $\sum_{i=1}^n \sigma_i(X)$ ,  $\text{tr}(X)$ , or  $\langle I, X \rangle$  in various literature). Also, the term  $\epsilon = \|\eta\|$  measures the total noise, whereas  $\epsilon_{\text{rel}} = \|\eta\|/\|b\|$  discussed in Section 2.2 measures the noise ratio or

relative noise. We choose the Schatten  $p$ -norm to highlight the fact that (2.3.9) is the semidefinite analog to the celebrated basis pursuit denoising problem [CDS01], [CRT06],

$$\begin{aligned} \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|Ax - b\|_2 \leq \epsilon, \end{aligned} \tag{2.3.10}$$

where minimizing  $\|x\|_1$  serves as a convex relaxation to minimizing the discrete, nonconvex function  $\|x\|_0 = \text{nnz}(x)$ .

The next two theorems quantify the relaxation of the PhaseLift model (2.3.9) from (2.1.1), establishing exact and approximate recovery guarantees for the PhaseLift model [CL14], [CSV13]. Theorem 2.3.1 applies to the noiseless case (where  $\epsilon = 0$  in the PhaseLift model), establishing a probability of equivalence between (2.1.1) and PhaseLift (2.3.9).

**Theorem 2.3.1.** *Consider an arbitrary signal  $\bar{x}$  in  $\mathbb{R}^n$  or  $\mathbb{C}^n$  and assume the sensing vectors  $a_i$  are independent, uniformly distributed on the unit sphere of  $\mathbb{R}^n$  or  $\mathbb{C}^n$ . Suppose that the number of measurements obeys  $m \geq c_0 n$ , where  $c_0$  is a sufficiently large constant. Then in both the real and complex cases, the solution to (2.3.9) is exact with high probability in the sense that the noiseless PhaseLift problem (2.3.9 with  $\epsilon = 0$ ) has a unique solution obeying*

$$X = \bar{x}\bar{x}^*.$$

*This holds with probability at least  $1 - \mathcal{O}(e^{-\gamma m})$ , where  $\gamma$  is a positive absolute constant.*

*Proof.* See [CL14, Section 2]. □

The assumptions of Theorem 2.3.1 are met by the experimental models used in [CESV13], [CSV13], [FM16], and this dissertation (see Section 2.2, where the masks  $C_i \in \mathbb{C}^{n \times n}$  from (2.2.1) have diagonal entries chosen with random Gaussian distribution). Thus, if  $L$  is the oversampling rate (i.e.,  $m = nL$ ) then Theorem 2.3.1 shows that only  $L = \mathcal{O}(1)$  masks are required to guarantee exact signal recover (up to global phase) with high probability.

Theorem 2.3.2 applies to the noisy phase retrieval case ( $\epsilon > 0$ ), where there is no guarantee that the PhaseLift (2.3.9) solution matrix  $X$  is low rank. Thus the solution signal  $x$  is set to an appropriate rescaling of the eigenvector  $v_1$  corresponding to the algebraically largest eigenvalue  $\lambda_1 = \lambda_1(X)$ , giving

$$x = \sqrt{\lambda_1} v_1. \tag{2.3.11}$$

**Theorem 2.3.2.** *Consider an arbitrary signal  $\bar{x}$  in  $\mathbb{R}^n$  or  $\mathbb{C}^n$  and assume the sensing vectors  $a_i$  are independent, uniformly distributed on the unit sphere of  $\mathbb{R}^n$  or  $\mathbb{C}^n$ . And suppose that the number of measurements obeys  $m \geq C_0 n \log n$ , where  $C_0$  is a sufficiently large constant. Then the PhaseLift (2.3.9) solution  $X$  obeys*

$$\|X - \bar{x}\bar{x}^*\|_F \leq C_0 \epsilon, \tag{2.3.12}$$

for some positive constant  $C_0$ . Additionally, we have

$$\|x - e^{i\theta}\bar{x}\|_2 \leq C_0 \min(\|\bar{x}\|_2, \epsilon/\|\bar{x}\|_2), \quad (2.3.13)$$

where  $x$  is defined as in (2.3.11) and we have some  $\theta \in [0, 2\pi]$ . Both these estimates hold with probability at least  $1 - \mathcal{O}(e^{-\gamma \frac{m}{n}})$ , where  $\gamma$  is a positive absolute constant.

*Proof.* See [CSV13, Section 6]. □

The strength of the PhaseLift model lies in its convexity and generality (no signal or observation assumptions are necessary). Yet the weakness is the difficulty in optimizing the objective function  $\|X\|_1$ , as its evaluation requires a partial SVD.

Like the PhaseLift method, the recent Wirtinger flow (wflow) algorithm [CLS15] also seeks to solve the phase retrieval problem (2.1.1) without any assumptions on the structure or sparsity of the signal or observation. However, unlike the PhaseLift model, the wflow model does not involve matrix lifting and instead frames (2.1.1) as the least-squares problem

$$\min_x \frac{1}{2m} \sum_{i=1}^m (|a_i^* x|^2 - b_i)^2 = \frac{1}{2m} \|\mathcal{A}(xx^*) - b\|^2. \quad (2.3.14)$$

While this model is nonconvex, the authors develop an efficient gradient descent-like method which has a provable guarantee on exact recovery when initialized appropriately. Initialization of the wflow algorithm involves a rescaling of the algebraically largest eigenvector  $v_1$  of the sum of the outer products of the measurement vectors, scaled with the observation magnitudes

$$\lambda_1 = \lambda_1(\mathcal{A}^*b), \quad \mathcal{A}^*b := \sum_{i=1}^m b_i a_i a_i^*. \quad (2.3.15)$$

The authors motivate this choice of initialization by noting that if the sensing vectors  $a_i$  are i.i.d. with standard normal distribution, and  $\|\bar{x}\| = 1$ , then the expected value of the outer product sum  $\mathcal{A}^*b$  will be

$$\mathbb{E} \left[ \frac{1}{m} \mathcal{A}^*b \right] = I + 2\bar{x}\bar{x}^*. \quad (2.3.16)$$

Thus for large  $m$ , (2.3.15) has a high probability of returning a vector  $v_1$  closely aligned with the solution  $\bar{x}$ . Intuitively, this initialization can also be seen as maximizing  $\langle \mathcal{A}(vv^*), b \rangle = \langle vv^*, \mathcal{A}^*b \rangle = \langle v, [\mathcal{A}^*b]v \rangle$ , and hence the eigenvalue problem (2.3.15) will encourage  $\mathcal{A}(vv^*)$  to be collinear with  $b$ .

With this initialization, the authors recommend a gradient descent-like method with a preset stepsize strategy. Let  $f(x) = \frac{1}{2} \|\mathcal{A}(xx^*) - b\|^2$ , and the residual  $r = \mathcal{A}(xx^*) - b$ . The mapping  $f : x \rightarrow \frac{1}{2} \|\mathcal{A}(xx^*) - b\|^2$  from  $\mathbb{C}^n$  to  $\mathbb{R}$  is not holomorphic, and thus not

complex-differentiable. As a result, the authors appeal to Wirtinger derivatives [CLS15, Section 6] for the descent direction

$$\begin{aligned}\nabla f(x) &= [\mathcal{A}^*(\mathcal{A}(xx^*) - b)]x \\ &= [\mathcal{A}^*r]x \\ &= \left[ \sum_{j=1}^L C_j^* F^* \text{Diag}(r_j) F C_j \right] x.\end{aligned}\tag{2.3.17}$$

The stepsize is then chosen using the heuristics

$$\mu_k = \min\{1 - e^{-k/k_0}, \mu_{\max}\} \quad k_0 = 330, \mu_{\max} = 0.4,\tag{2.3.18}$$

where  $\mu_k$  starts small and increases [CLS15, Section 2]. This descent direction and stepsize choice lead to Algorithm 2.

---

**Algorithm 2** wflow algorithm

---

**Input:** Sensing operator  $\mathcal{A}$  (2.3.6) with sensing vectors  $a_i$ , frequency measurement vector  $b \in \mathbb{R}_+^m$ .

**Output:** Approximate solution signal  $x$ .

- 1: *Initialize:* Compute the algebraically largest eigenpair  $(\lambda_1, v_1)$  of  $\mathcal{A}^*b$ , set  $x_0 = \alpha v_1$  where  $\alpha^2 = n \sum_{i=1}^m b_i / \sum_{i=1}^m \|a_i\|^2$ , set  $k = 0$ .
  - 2: **while**  $\|res\| > \text{tol}$  **do**
  - 3:     Compute Wirtinger derivative:  $\nabla f(x_k)$  from (2.3.17).
  - 4:     Compute stepsize:  $\mu_{k+1}$  from (2.3.18) and set  $\alpha_k = \frac{\mu_{k+1}}{\|x_0\|^{2m}}$ .
  - 5:     Compute signal update:  $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$ .
  - 6:      $k = k + 1$ .
  - 7: **end while**
  - 8: return  $x \leftarrow x_k$ .
- 

In [CLS15, Section 2.3], it is observed that the wflow algorithm can be interpreted as a stochastic gradient scheme, where  $\nabla f(x)$  is an unbiased estimate of an ideal gradient  $\nabla F(x)$ , with

$$F(x) = x^* (I - \bar{x}x^*) x - \frac{3}{4} (\|x\|^2 - 1)^2.\tag{2.3.19}$$

Interestingly, we observe that this algorithm can also be interpreted as an alternating projection HIO-type method, where computation of the Wirtinger derivative (2.3.17) corresponds to steps 3-5 of the GS algorithm and the signal update  $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$  corresponds to step 6 of the GS algorithm. More precisely, in the last line of (2.3.17) we see  $x_k$  is first mapped to the frequency domain (Algorithm 1, step 3). This vector is then damped coordinate-wise by the corresponding residual values  $\text{Diag}(r_j)$  for  $j = 1, \dots, L$

(Algorithm 1, step 4). Finally, the observation is mapped back to the object domain (Algorithm 1, step 5) and damped by  $\alpha_k$  to generate the signal update  $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$  (similar to the HIO update (2.3.2) which replaces Algorithm 1, step 6).

Yet unlike other alternating projection methods, the wflow algorithm does not require additional frequency domain information or heuristics. Wflow replaces the frequency domain damping in Algorithm 1, step 4 with a damping which uses the observation residual  $r = \mathcal{A}(xx^*) - b$  as an exact, real-time measurement of the frequency domain violations of  $x$ . In this view,  $\text{Diag}(r_j)$  in (2.3.17) provides a heuristic-free frequency domain damping. Consequentially, the wflow algorithm merges the computational simplicity of an alternating projection method (Section 2.3.1) with the broad utility of an unstructured optimization method. Theorem 2.3.3 provides an exact recovery probability for Algorithm 2 applied to noiseless phase retrieval problems.

**Theorem 2.3.3.** *Consider an arbitrary signal  $\bar{x}$  in  $\mathbb{C}^n$  and assume the sensing vectors  $a_i$  are independent, uniformly distributed on the unit sphere of  $\mathbb{C}^n$ . Suppose that the number of measurements obeys  $m \geq c_0 n \log(n)$ , where  $c_0$  is a sufficiently large constant. Then the wflow initial estimate  $x_0$  normalized to have squared Euclidean norm equal to  $(\sum_{i=1}^m b_i)/m$  has*

$$\text{dist}(x_0, \bar{x}) := \min_{\theta \in [0, 2\pi]} \|x_0 - e^{i\theta} \bar{x}\| \leq \frac{1}{8} \|x\| \quad (2.3.20)$$

with probability at least  $1 - 10e^{-\gamma n} - 8/n^2$ , where  $\gamma$  is a positive absolute constant. Additionally, assume the stepsize is a bounded constant  $\mu_k = \mu \leq c_1/n$  for some fixed numerical constant  $c_1$ . Then there is an event of probability at least  $1 - 13e^{-\gamma n} - me^{-1.5m} - 8/n^2$ , such that on this event, starting from any initial solution  $x_0$  obeying (2.3.20), we have

$$\text{dist}(x_k, \bar{x}) \leq \frac{1}{8} \|\bar{x}\| \left(1 - \frac{\mu}{4}\right)^{k/2}. \quad (2.3.21)$$

*Proof.* See [CLS15, Section 7]. □

Computationally, the wflow algorithm is very efficient at recovering the signal of an oversampled noiseless observation. Numerical experiments in [CLS15, Section 4] demonstrate the effectiveness of this algorithm for 1-D random signals, 2-D natural images, and 3-D molecular structures.

More generally, it was proved in [SQW16] that if  $a_i$  are independent, uniformly distributed on the unit sphere of  $\mathbb{C}^n$  and there are  $m \geq c_0 n \log^3(n)$  measurements, then with probability at least  $1 - c_0/m$  the function  $f(x) = \frac{1}{2} \|\mathcal{A}(xx^*) - b\|^2$  has the properties:

- $f$  has no spurious local minima,
- all global minima are equal to  $\bar{x}$  up to some phase constant:  $x = e^{i\theta} \bar{x}$ ,
- $f$  has negative directional curvature at all saddle points.

Thus when solving (2.3.14), methods such as wflow do not require specialized initialization, and are likewise guaranteed to find a global minimum. Nevertheless, wflow can diverge when significant noise exists and an insufficient number of observations are present (see Section 5.3).

As we have seen in this section, many phase retrieval methods have been developed to utilize structural properties of the given phase retrieval problem (2.1.1). Yet if the problem is unstructured (only the observation and noise ratio are available), then only a few methods are available. We will now examine these methods to establish why the PLGD model is best suited for handling large-scale, noisy phase retrieval.

## 2.4 Justification for optimizing the PLGD model

In this section we justify using the PLGD model (1.1.1) to solve large-scale, unstructured noisy phase retrieval problems (2.1.1). First, we examine several *PhaseLift-type* models, which are either duals or modifications of the PhaseLift model (2.3.9) discussed in Section 2.3.3. Among these PhaseLift-type models, we explain why the PLGD model is the best suited for developing a first-order optimization method. Next, we demonstrate experimentally that the first-order PLGD method used in this dissertation is more accurate and robust than the wflow algorithm (Algorithm 2) for noisy phase retrieval with minimal oversampling.

Note that the first-order PLGD method discussed in this section is presented fully in Section 4.2, where we also review the steps necessary to develop a generic first-order method. For now we simply assume that any first-order method for minimizing an objective function  $f(x)$  over some convex set  $\mathcal{C}$  will require several evaluations of  $f(x)$  and its gradient  $\nabla f(x)$  or subdifferential (1.2.12)  $\partial f(x)$ , and several projections onto  $\mathcal{C}$ . Also note that the method used in this section for creating noisy phase retrieval problems is presented in Section 5.2.

We begin by examining several PhaseLift-type models and discussing the computational costs involved in a typical first-order method for each model. As discussed in equation (2.3.8), any PhaseLift-type model will have as a variable a lifted matrix  $X$  of size  $n \times n$ , where  $n$  is the dimension of the desired signal  $\bar{x}$ . Since we are focused on large-scale phase retrieval problems, we seek to avoid repeated partial singular value decompositions (SVDs) of  $n \times n$  matrices. For instance, the PhaseLift model (2.3.9) has the objective function  $f(X) = \|X\|_1$  and constraints  $\|\mathcal{A}(X) - b\|_2 \leq \epsilon$  and  $X \succeq 0$ . The evaluation of  $f$  and its subdifferential require a partial SVD, and projection onto  $X \succeq 0$  requires an additional partial SVD. Thus the PhaseLift model (2.3.9) is not well suited for first-order methods.

In contrast with the PhaseLift model (2.3.9), the PLGD model

$$\begin{aligned} \min_y \quad & \lambda_1(\mathcal{A}^*y) \\ \text{(PLGD) s.t.} \quad & \langle b, y \rangle - \epsilon \|y\| \geq 1 \end{aligned} \tag{2.4.1}$$

has an objective function  $f(y) = \lambda_1(\mathcal{A}^*y)$  whose evaluation and gradient require a standard eigenvalue problem, and projection onto  $\mathcal{C} = \{y \mid \langle b, y \rangle - \epsilon \|y\| \geq 1\}$  is an  $\mathcal{O}(n)$  operation (see Section 4.2 for details). Likewise, recovery of the desired approximate signal  $x$  from the variable  $y$  involves a wflow-like problem which is very efficient for large-scale problems. Thus the PLGD model (2.4.1) is well suited for developing first-order methods.

Another PhaseLift-type model is the PhaseLift Lagrange dual

$$\begin{aligned} \max_y \quad & \langle b, y \rangle - \epsilon \|y\| \\ \text{s.t.} \quad & I \succeq \mathcal{A}^*y \end{aligned} \tag{PLD} \tag{2.4.2}$$

(see [BV04, Chapter 5] for derivation). The PLD model (2.4.2) has a simple objective function to evaluate. Yet the PLD constraint is a complicated linear matrix inequality and projection onto the feasible set  $\{y \mid I \succeq \mathcal{A}^*y\}$  is a separate eigenvalue optimization problem similar to PLGD model (2.4.1). Thus the PLGD model is better suited for first-order methods than the PLD model (2.4.2).

Another method for dualizing the PhaseLift model (2.3.9) is considered in [CSV13] where the authors demonstrate the efficacy of the PhaseLift model by applying a Lagrange multiplier to the constraint  $\|\mathcal{A}(X) - b\| \leq \epsilon$  in (2.3.9), giving the model

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathcal{A}(X) - b\|^2 + \tau \|X\|_1 \\ \text{s.t.} \quad & X \succeq 0. \end{aligned} \tag{2.4.3}$$

To optimize (2.4.3), the authors use the Templates for First-Order Conic Solvers (TFOCS) package [BCG11], which dualizes a given model, applies a smoothing, and then solves the smooth dual with a first-order method. For the appropriate value  $\tau(\epsilon)$ , the model (2.4.3) is equivalent to the PhaseLift model (2.3.9) [Roc70, Section 28]. Thus (2.3.9) is solved by maximizing  $\tau(\epsilon)$  with a bisection method which requires solving a sequence of models (2.4.3) using TFOCS.

This PhaseLift bisection method successfully demonstrates that the accuracy of the PhaseLift optimal signal improves with oversampling and decreases gradually as the error  $\epsilon$  increases [CSV13, Section 7]. However, this method is computationally expensive, as the function  $\|X\|_1$  in the objective requires a partial SVD and the constraint  $\|\mathcal{A}(xx^*) - b\| \leq \epsilon$  must be embedded into the objective. Additionally, this method tends to fail without sufficient oversampling and is unable to achieve a high level of accuracy for noiseless models [FM16, Section 5, Table 1].

Another potential PhaseLift-type model we may consider optimizing is the *PhaseLift- $l_1$  model*

$$\begin{aligned} \text{(PLP-}l_1) \quad \min \quad & \|\mathcal{A}(X) - b\|_1 \\ \text{s.t.} \quad & X \succeq 0 \end{aligned} \tag{2.4.4}$$

discussed in [CL14]. This paper showed that the number of observations required in Theorem 2.3.2 to guarantee the solution signal error bounds (2.3.12), (2.3.13) are reduced to  $\mathcal{O}(n)$  if we instead consider the PhaseLift- $l_1$  model (2.4.4). Thus we investigate whether



(2.4.4) or its duals could be used to develop a computationally efficient large-scale first-order method.

A first-order method for the PhaseLift- $l_1$  model (2.4.4) will require projection onto the positive semidefinite constraint  $X \succeq 0$ , again requiring a partial SVD which is prohibitive for large-scale problems. We may also consider the gauge dual and Lagrange dual

$$\begin{aligned}
& \min_y \|y\|_\infty & \max_y \langle b, y \rangle \\
\text{(PLGD-}l_1\text{)} \quad \text{s.t.} \quad & -\mathcal{A}^*y \succeq 0 & \text{(PLD-}l_1\text{)} \quad \text{s.t.} \quad & -\mathcal{A}^*y \succeq 0 \\
& \langle b, y \rangle \geq 1 & & \|y\|_\infty \leq 1
\end{aligned} \tag{2.4.5}$$

to the PhaseLift- $l_1$  model (2.4.4) (see Appendix A for the derivation of PLGD- $l_1$  and [BV04, Chapter 5] for PLD- $l_1$ ). Yet like the PhaseLift Lagrange dual (2.4.2), both of these models (2.4.5) have a linear matrix inequality in the constraint. As a result, projection onto this constraint  $\{y \mid -\mathcal{A}^*y \succeq 0\}$  will again require a separate eigenvalue optimization problem similar to PLGD model (2.4.1), making both models (2.4.5) computationally prohibitive to optimize. Thus, among the PhaseLift-type models discussed above, the PLGD model (2.4.1) is the best suited for developing a first-order method.

Next, we demonstrate that the first-order PLGD method discussed in this dissertation (Algorithm 3 of Section 4.2) is typically more accurate and robust than the wflow algorithm (Algorithm 2) for noisy phase retrieval problems (2.1.1) with minimal oversampling. To compare these two algorithms, we generate a set of noisy phase retrieval problems using the method described in Section 5.2 (the phase retrieval problem with Gaussian noise (5.2.3)). Successful signal recovery occurs when the approximate observation  $\mathcal{A}(xx^*)$  as defined in (2.3.7) closely matches the true observation  $\bar{b} = \mathcal{A}(\bar{x}\bar{x}^*)$  rather than the noisy observation  $b$ . Thus we use the primal true relative error

$$\frac{\|\mathcal{A}(xx^*) - \bar{b}\|}{\|\bar{b}\|} \leq \tau \epsilon_{\text{rel}} \tag{2.4.6}$$

to measure the accuracy of these algorithms. A signal is considered successfully recovered if it satisfies the inequality (2.4.6), where  $\tau = 1$  indicates accuracy within the expected error, and  $\tau < 1$  indicates a higher level of accuracy.

The ability of Algorithm 3 to denoise the noisy observation  $b = \bar{b} + \eta$  is a consequence of the property that PLGD variable  $y_k$  approximates the noise term  $\eta$  with increasing accuracy as Algorithm 3 converges (see Section 3.4 for details regarding the primal refinement method). Thus we also measure the angle between  $\eta$  and the final dual variable  $y$  returned by Algorithm 3

$$\cos \angle(\eta, y) = \frac{\eta^* y}{\|\eta\| \|y\|}. \tag{2.4.7}$$

Table 2.1 displays the results of the first-order PLGD method (Algorithm 3 of Section 4.2) and wflow algorithm (Algorithm 2).

		Algorithm 3				wflow		
$L$	$\epsilon_{\text{rel}}$	$\cos \angle(\eta, y_{100})$	xErr	% success		xErr	% success	
				$\tau = 1.0$	$\tau = 0.8$		$\tau = 1.0$	$\tau = 0.8$
4	0.050	1.12 <sub>-1</sub>	1.50 <sub>-1</sub>	0.86	0.00	3.92 <sub>-1</sub>	0.01	0.01
4	0.150	3.57 <sub>-1</sub>	6.21 <sub>-1</sub>	0.07	0.00	5.58 <sub>-1</sub>	0.00	0.00
4	0.300	5.65 <sub>-1</sub>	1.17 <sub>0</sub>	0.30	0.00	1.00 <sub>0</sub>	0.04	0.00
6	0.050	1.89 <sub>-1</sub>	6.92 <sub>-2</sub>	1.00	0.96	1.25 <sub>-1</sub>	0.64	0.64
6	0.150	4.27 <sub>-1</sub>	2.58 <sub>-1</sub>	1.00	0.93	2.40 <sub>-1</sub>	0.49	0.49
6	0.300	6.12 <sub>-1</sub>	6.72 <sub>-1</sub>	1.00	0.20	4.21 <sub>-1</sub>	0.64	0.32
8	0.050	4.00 <sub>-1</sub>	4.61 <sub>-2</sub>	1.00	1.00	4.53 <sub>-2</sub>	1.00	1.00
8	0.150	5.32 <sub>-1</sub>	1.55 <sub>-1</sub>	1.00	1.00	1.42 <sub>-1</sub>	0.98	0.98
8	0.300	6.68 <sub>-1</sub>	4.01 <sub>-1</sub>	1.00	1.00	2.97 <sub>-1</sub>	0.98	0.94

Table 2.1: Rate of successful signal recovery and mean residual values for sets of 100 noisy phase retrieval problems with random Gaussian signals of size  $n = 128$  with oversampling rate  $L$  and relative error  $\epsilon_{\text{rel}}$ . The term  $xErr$  is signal relative error  $\|xx^* - \bar{x}\bar{x}^*\|_F / \|\bar{x}\bar{x}^*\|_F$ . Recovery is determined successful if the inequality (2.4.6) is satisfied for a given  $\tau$ . The first-order PLGD method (Algorithm 3) is set to terminate at 100 iterations. Numbers  $n_{-k}$  are shorthand for  $n \times 10^{-k}$ .

As we see in Table 2.1, Algorithm 3 generally has a greater likelihood of successful recovery than wflow when observations have a lower rate of oversampling, regardless of the noise level. Additionally, if models have greater noise and greater oversampling, this increases the accuracy of the PLGD variable  $y$  approximating the noise term  $\eta$ . Figure 2.4 depicts the ability of Algorithm 3 to recover a much larger signal with moderate noise and minimal oversampling.

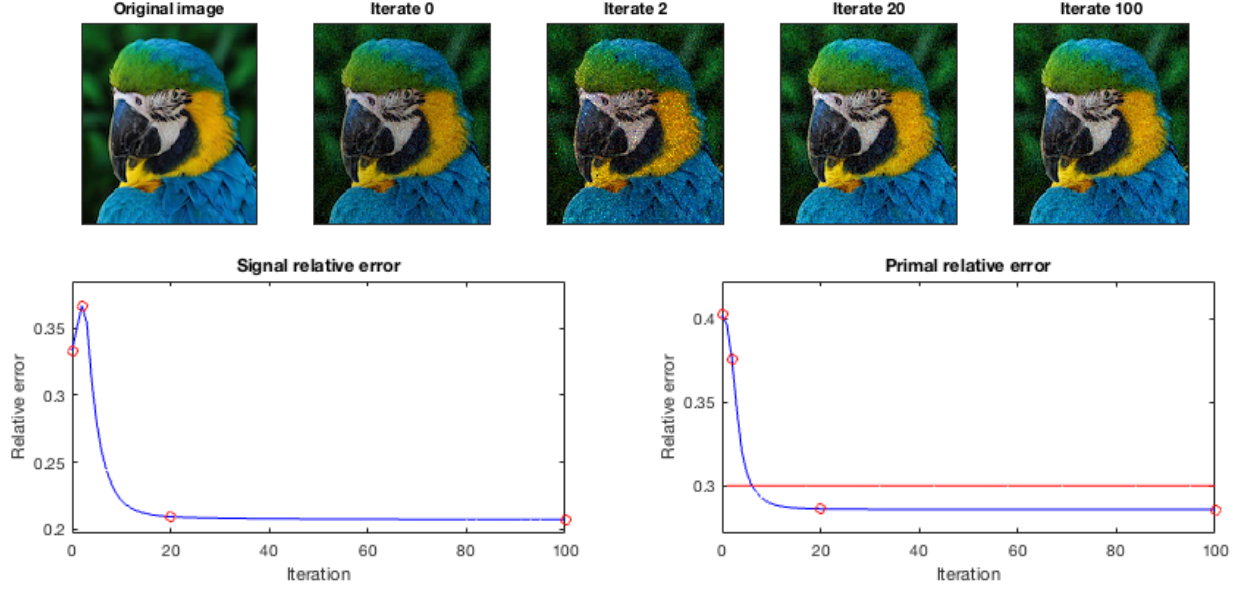


Figure 2.4: Results from the first-order PLGD method (Algorithm 3) applied to a test image separated into its three RGB channels. Left: signal relative error  $\|xx^* - \bar{x}\bar{x}^*\|_F / \|\bar{x}\bar{x}^*\|_F$ . Right: primal relative error  $\|\mathcal{A}(xx^*) - b\| / \|b\|$  with red line indicating primal feasibility. Red circles denote the iterate signals pictured above. Original signal is  $128 \times 128$  pixels, with an oversampling of  $L = 8$  and noise ratio  $\epsilon_{\text{rel}} = 0.30$ . Measurements use the mean of the three color channel values.

Figure 2.4 demonstrates the tendency of Algorithm 3 to make significant progress during early iterates. When the same set of models in Figure 2.4 were solved with wflow, the red channel converged to an infeasible solution (with primal residual 0.300068), the green channel converged to a feasible solution (with primal residual 0.288473), and the blue channel diverged.

Table 2.1 and Figure 2.4 highlight the effectiveness of the first-order PLGD method (Algorithm 3) for noisy phase retrieval problems. This method is typically more robust than wflow (Algorithm 2) and tends to recover signals successfully when there is minimal oversampling. Among the phase retrieval models and methods considered in this chapter, the PLGD model (2.4.1) is the best suited for large-scale, unstructured noisy phase retrieval problems. In the next two chapters we present the theory behind the PLGD model and develop a first-order method (Algorithm 3) for optimizing this model.

## Chapter 3

# Gauge duality theory and the PhaseLift gauge dual model

### 3.1 Introduction

This chapter examines the PLGD model for solving the phase retrieval problem (2.1.1) and presents the gauge duality theory necessary for examining this model. The PLGD model is the constrained eigenvalue optimization problem

$$\begin{aligned} \min_y \quad & \lambda_1(\mathcal{A}^*y) \\ \text{s.t.} \quad & \langle b, y \rangle - \epsilon \|y\|_2 \geq 1. \end{aligned} \tag{3.1.1}$$

We begin in Section 3.2 by establishing the PhaseLift primal and gauge dual (PLP, PLGD) pair of models along with some basic definitions. We then contrast gauge duality with Lagrange duality and provide a brief history of gauge duality theory. Section 3.3 develops a sequence of propositions which are used to prove the gauge duality theorem (Theorem 3.3.1) along with weak duality, strong duality, and optimality conditions for an inequality constrained gauge dual pair. Finally, Section 3.4 uses these results to establish key properties of the primal and gauge dual models which are necessary for developing an efficient method of optimizing the PLGD model and recovering the signal  $x$  from the dual variable  $y$ .

All of the gauge properties and gauge duality results in this chapter were previously established in [Roc70], [Fre87], [FMP14], and [FM16]. Gauge functions were first analyzed in [Roc70]. Gauge duality was then introduced in [Fre87], where the author focused on quadratic programming applications. In [FMP14], the authors developed a broad set of antipolar calculus results for the analysis of gauge duality. These results were then applied to the PLP-PLGD pair to develop a first-order method in [FM16].

Our contribution is to provide a self-contained, comprehensive treatment of gauge duality theory focused solely on establishing the duality theorem and optimality conditions

for the PLP-PLGD pair. Prior to this treatment, the results in this chapter were spread throughout the original texts ([Roc70], [FMP14], and [FM16]), occasionally with differing notation or style. In this chapter, we present a single notation and style (Section 4.1) which we use to develop gauge duality theory (Section 3.3) and then apply this theory to the PLP-PLGD pair (Section 3.4). The results of this chapter provide the theoretical foundation for developing a first-order optimization method for the PLGD model in Chapter 4.

## 3.2 Definitions and primal-dual pairs

The PhaseLift primal semidefinite program (restated from (2.3.9) in Section 2.3.3) and its gauge dual are

$$\begin{aligned}
 \min_X \quad & \|X\|_1 = \sum_{i=1}^n \sigma_i(X) & \min_y \quad & \lambda_1(\mathcal{A}^*y) \\
 \text{(PLP)} \quad \text{s.t.} \quad & \|\mathcal{A}(X) - b\| \leq \epsilon & \text{(PLGD)} \quad \text{s.t.} \quad & \langle b, y \rangle - \epsilon \|y\| \geq 1. \\
 & X \succeq 0 & & 
 \end{aligned} \tag{3.2.1}$$

If  $\epsilon > 0$  then we refer to (3.2.1) as the *noisy PLP-PLGD pair*, where the observation vector  $b = \bar{b} + \eta$  is the sum of the true observation  $\bar{b}$  with some nontrivial noise  $\eta$  as described in the phase retrieval problem (2.1.1). This section places the PLP-PLGD pair (3.2.1) in the context of gauge duality theory. We begin by discussing definitions and basic properties which are relevant to gauge duality theory. We then take a moment to highlight some parallels between gauge duality and Lagrange duality before closing with a summary of major developments in gauge duality theory.

Following the convention of these major theoretical advances in gauge duality ([Fre87], [FM16], [ABD<sup>+</sup>17]), all notation and definitions in this section are based on [Roc70]. All unproven results in this section are accompanied by a reference to the appropriate section of [Roc70]. Along the way, we present three additional, more general primal-gauge dual models which have PLP-PLGD (3.2.1) as a special case. Two of these models are then used in Section 3.3 to develop the theory which establishes the gauge duality of the pair (3.2.1).

The PLP-PLGD pair (3.2.1) are examples of a more general primal-gauge dual pair which we define as the *nonlinearly-constrained* pair

$$\begin{aligned}
 \min_{x \in \mathcal{X}} \quad & \kappa(x) & \min_{z \in \mathcal{X}} \quad & \kappa^\circ(z) \\
 \text{(P-nonline)} \quad \text{s.t.} \quad & x \in \mathcal{C} & \text{(GD-nonline)} \quad \text{s.t.} \quad & z \in \mathcal{C}'.
 \end{aligned} \tag{3.2.2}$$

Here,  $\mathcal{C}$  and  $\mathcal{C}'$  are subsets of  $\mathcal{X}$ , a finite-dimensional Euclidean space. The set  $\mathcal{C}'$  is the *antipolar* of  $\mathcal{C}$ , defined as

$$\mathcal{C}' = \{z \mid \langle x, z \rangle \geq 1 \ \forall x \in \mathcal{C}\}. \tag{3.2.3}$$

This is in contrast to the *polar* of  $\mathcal{C}$ , which is defined as

$$\mathcal{C}^\circ = \{z \mid \langle x, z \rangle \leq 1 \ \forall x \in \mathcal{C}\}. \quad (3.2.4)$$

Although we refer to the primal-gauge dual pair (3.2.2) as nonlinear, we are primarily concerned with closed, convex sets  $\mathcal{C}$  which do not contain the origin (e.g., the PLP (3.2.1) constraint set). Thus gauge dual models with linear constraints which do not include the origin are a subset of the models described by (3.2.2).

The functions  $\kappa, \kappa^\circ : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$  are *gauge* functions, meaning they are convex, nonnegative, positively homogeneous ( $\kappa(\alpha x) = \alpha \kappa(x)$  for all  $\alpha > 0$ ), and vanish at the origin. The gauge function generalizes the notion of norm, allowing for flexibility in modeling the phase retrieval problem. In the PLP model for instance, the Schatten 1-norm  $\kappa(X) := \|X\|_1$  and vector 2-norm  $\rho(y) := \|y\|$  are both gauges.

Given a gauge function  $\kappa$ , the *polar* of this function is the function  $\kappa^\circ$  that most tightly satisfies the inequality

$$\langle x, z \rangle \leq \kappa(x) \kappa^\circ(z) \quad \forall x \in \text{dom } \kappa, \ \forall z \in \text{dom } \kappa^\circ. \quad (3.2.5)$$

Equivalently, the polar may be defined as [Roc70, Section 15]

$$\kappa^\circ(z) = \inf \{ \mu > 0 \mid \langle x, z \rangle \leq \mu \kappa(x) \ \forall x \}. \quad (3.2.6)$$

Note that the polar is a generalization of the *dual norm*, which is defined as

$$\|z\|_* = \sup_x \{ \text{Re} \langle x, z \rangle \mid \|x\| \leq 1 \}. \quad (3.2.7)$$

The *preimage* of  $A$  over the set  $\mathcal{S} \subseteq \mathcal{Y}$  is defined as

$$A^{-1}\mathcal{S} = \{x \in \mathcal{X} \mid Ax \in \mathcal{S}\}. \quad (3.2.8)$$

The *closure* of a set  $\mathcal{S} \subseteq \mathcal{X}$  is denoted  $\text{cl}(\mathcal{S})$ . The *affine hull* of  $\mathcal{S}$  is the set of all affine combinations of elements of  $\mathcal{S}$ , or

$$\text{aff}(\mathcal{S}) = \left\{ \sum_{i=1}^k \alpha_i x_i \mid k > 0, \ x_i \in \mathcal{S}, \ \alpha_i \in \mathbb{R}, \ \sum_{i=1}^k \alpha_i = 1 \right\}. \quad (3.2.9)$$

The *relative interior* of  $\mathcal{S}$ , denoted  $\text{ri}(\mathcal{S})$ , is the interior within the affine hull of  $\mathcal{S}$ , i.e.,

$$\text{ri}(\mathcal{S}) = \{x \in \mathcal{S} \mid \exists \epsilon > 0, \ B_\epsilon(x) \cap \text{aff}(\mathcal{S}) \subseteq \mathcal{S}\}, \quad (3.2.10)$$

where  $B_\epsilon(x)$  is a ball of radius  $\epsilon$  centered at  $x$ . The *support* function  $\sigma_{\mathcal{C}}$  of a nonempty convex set  $\mathcal{C}$  is defined as

$$\sigma_{\mathcal{C}}(z) = \sup_{x \in \mathcal{C}} \langle x, z \rangle. \quad (3.2.11)$$

And the *Minkowski* function  $\gamma_{\mathcal{C}}$  is defined as

$$\gamma_{\mathcal{C}}(x) = \inf \{ \lambda \geq 0 \mid x \in \lambda \mathcal{C} \}. \quad (3.2.12)$$

If there is no  $\lambda$  such that  $\lambda x \in \mathcal{C}$ , then  $\gamma_{\mathcal{C}}(x) = +\infty$ . Note that any gauge  $\kappa$  is a Minkowski function  $\gamma_{\mathcal{C}}$  for  $\mathcal{C} = \{x \mid \kappa(x) \leq 1\}$  [Roc70, Section 15]. Given a function  $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$ , the *epigraph* of  $f$  is defined as

$$\text{epi}(f) = \{(x, \tau) \mid f(x) \leq \tau\}. \quad (3.2.13)$$

Note that  $f$  is convex if and only if  $\text{epi}(f)$  is convex [Roc70, Section 7]. Thus the function  $f$  is said to be *closed* if  $\text{epi}(f)$  is closed. Additionally,  $f$  is closed if and only if it is lower-semicontinuous (that is,  $\liminf_{x \rightarrow x_0} f(x) \geq f(x_0)$  for all  $x_0$  in  $\text{dom}(f)$ ) [Roc70, Section 7]. Also,  $f$  is *proper* if the domain of  $f$ ,  $\text{dom}(f) = \{x \mid f(x) < +\infty\}$  is nonempty.

If  $\kappa$  is a closed gauge, then its polar may also be expressed as [Roc70, Section 15]

$$\kappa^\circ(z) = \sup_x \{ \langle x, z \rangle \mid \kappa(x) \leq 1 \}, \quad (3.2.14)$$

highlighting the fact that the polar function is a generalization of the dual norm (3.2.7). Additionally, if  $\kappa$  is also positive everywhere except at the origin then its polar may also be defined as [Roc70, Section 15]

$$\kappa^\circ(z) = \sup_x \left\{ \frac{\langle x, z \rangle}{\kappa(x)} \right\}. \quad (3.2.15)$$

Given the gauge duality notation discussed above, we take a moment to contrast gauge duality with the much more common Lagrange duality. Whereas gauge duality involves multiplicative duality transformations, Lagrange duality is additive in nature. The reader may see [BV04, Chapter 5] for a comprehensive introduction to Lagrange duality, and [Roc70, Section 28] or [BTN01, Chapter 2] for a treatment of Lagrange duality theory.

The nonlinear pair (3.2.2) demonstrates that the gauge dual model GD-nonline can be described simply using the polar function  $\kappa^\circ$  and the antipolar set  $\mathcal{C}'$ . Similarly, the Lagrange dual model can be described using the appropriate function and set transformations. Given a function  $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{\pm\infty\}$ , the *convex conjugate*  $f^*$  is defined as the function that most tightly satisfies the inequality

$$\langle x, z \rangle \leq f(x) + f^*(z) \quad \forall x \in \text{dom } f, \forall z \in \text{dom } f^*. \quad (3.2.16)$$

Equivalently, the convex conjugate may be defined as [Roc70, Section 12]

$$f^*(z) = \sup_x \{ \langle x, z \rangle - f(x) \}. \quad (3.2.17)$$

We see that (3.2.16) and (3.2.17) are the additive analogs of the polar function definitions (3.2.5) and (3.2.15), respectively. Additionally, for a set  $\mathcal{S} \subseteq \mathcal{X}$ , the *dual cone* is defined as

$$\mathcal{S}^* = \{z \mid \langle x, z \rangle \leq 0 \forall x \in \mathcal{S}\}. \quad (3.2.18)$$

Given the convex conjugate  $\kappa^*$  and the dual cone  $\mathcal{C}^*$ , the Lagrange dual D-nonline and gauge dual GD-nonline both have simple forms

$$\begin{array}{ll} \max_z & -\kappa^*(z) \\ \text{(D-nonline)} \quad \text{s.t.} & z \in \mathcal{C}^* \end{array} \qquad \begin{array}{ll} \min_{z \in \mathcal{X}} & \kappa^\circ(z) \\ \text{(GD-nonline)} \quad \text{s.t.} & z \in \mathcal{C}'. \end{array} \quad (3.2.19)$$

Note that Lagrange duality applies for any function  $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{\pm\infty\}$  and set  $\mathcal{S} \subseteq \mathcal{X}$ . Thus Lagrange duality has been studied thoroughly and applied extensively (e.g., see [BV04, Chapter 5]). In contrast, gauge duality places specific restrictions on the objective function and constraint set. As a result, the development of gauge duality theory has a brief, sporadic history.

In 1970, Rockafellar thoroughly analyzed gauge functions and their polars [Roc70, Part III]. The concept of gauge duality was then introduced by Freund in 1987 [Fre87]. This seminal work focuses primarily on the *linearly-constrained* primal and gauge dual pair

$$\begin{array}{ll} \min_{x \in \mathcal{X}} & \kappa(x) \\ \text{(P-lin)} \quad \text{s.t.} & Ax = b \end{array} \qquad \begin{array}{ll} \min_{y \in \mathcal{Y}} & \kappa^\circ(A^*y) \\ \text{(GD-lin)} \quad \text{s.t.} & \langle b, y \rangle = 1. \end{array} \quad (3.2.20)$$

Here  $A : \mathcal{X} \rightarrow \mathcal{Y}$  is a linear operator over finite-dimensional Euclidean spaces. As with Lagrange duality, if the primal constraint is replaced with  $Ax \geq b$  then the gauge dual also includes the constraint  $y \geq 0$ . Freund develops strong duality and optimality conditions for the pair (3.2.20) based on polarity relationships for sets and gauge functions. His work also establishes these conditions for the nonlinear pair (3.2.2). In this case, his work requires  $\mathcal{X}$  and  $\mathcal{Y}$  to be *ray-like* (meaning that for all  $x, y \in \mathcal{X}$  we have  $x + \alpha y \in \mathcal{X}$  for all  $\alpha \geq 0$ ). The addition of the ray-like property to  $\mathcal{X}$  (along with closed, convex, and containing the origin) guarantees  $\mathcal{X}'' = \mathcal{X}$ .

Gauge duality theory was revisited in [FMP14], where the authors consider the *inequality-constrained* primal and gauge dual pair

$$\begin{array}{ll} \min_{x \in \mathcal{X}} & \kappa(x) \\ \text{(P-ineq)} \quad \text{s.t.} & \rho(Ax - b) \leq \epsilon \end{array} \qquad \begin{array}{ll} \min_{y \in \mathcal{Y}} & \kappa^\circ(A^*y) \\ \text{(GD-ineq)} \quad \text{s.t.} & \langle b, y \rangle - \epsilon \rho^\circ(y) \geq 1. \end{array} \quad (3.2.21)$$

As we will see in Section 3.4, the PLP-PLGD pair (3.2.1) pair is recovered when we set  $\kappa(X) = \|X\|_1 + \delta_{(\cdot) \geq 0}(X)$  and  $\rho(y) = \|y\|$ . The authors of [FMP14] develop an antipolar calculus to determine the antipolar for sets like  $\{y \mid \rho(Ax - b) \leq \epsilon\}$  and use this calculus rather than polarity relations to derive the gauge dual pair (3.2.21) and establish conditions such as those for strong duality. Additionally, this antipolar calculus allows the authors to drop the requirement that the sets  $\mathcal{X}$  and  $\mathcal{Y}$  are ray-like.

In contrast to the antipolar calculus framework, the authors of [ABD<sup>+</sup>17] develop gauge duality through a perturbation framework. This even broader setting frames gauge duality as a product of Fenchel-Rockafellar duality, allowing the consideration of gauge duality for general nonnegative convex functions.



### 3.3 Theory for linearly and nonlinearly constrained models

In this section we develop gauge duality theory for the inequality-constrained primal and gauge dual pair (3.2.21). We begin by establishing the propositions necessary to show the gauge duality of (3.2.21) is a consequence of the duality of both the nonlinear (3.2.2) and linear (3.2.20) models. This two-track proof strategy as depicted in Figure 3.1 leads to Theorem 3.3.1, (i) and (ii), respectively.

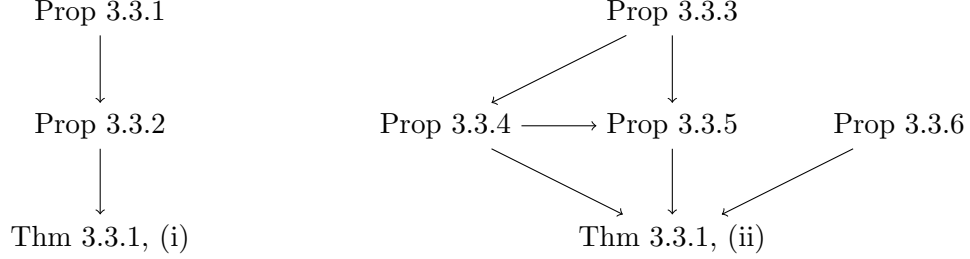


Figure 3.1: Dependency chart for proofs of Theorem 3.3.1, (i) and (ii).

Next we establish weak duality, strong duality, and optimality conditions for (3.2.21). In Section 3.4 we return to the PhaseLift model (3.2.1), where these optimality conditions provide a method for recovery of the primal signal  $x$  from a dual solution  $y$ .

The results in this section are based on [Roc70], [Fre87], and especially [FMP14], and rely on the analysis of polarity relations rather than perturbation analysis as discussed in [ABD<sup>+</sup>17]. In particular, the two-track proof strategy depicted in Figure 3.1 was first developed in [FMP14] as part of a larger treatment of antipolar calculus, and refers to [Roc70] for more elementary results (e.g., Propositions 3.3.3, 3.3.4, and 3.3.6). This section provides a self-contained development of gauge duality for (3.2.21).

The next two propositions are used in Theorem 3.3.1, (i) to show the inequality pair (3.2.21) is an instance of the nonlinear pair (3.2.2). Since (3.2.21) allows for the constraint set  $\mathcal{C}$  to be any closed, convex set not containing the origin, we must simply establish the antipolar of the particular constraint set  $\mathcal{C} = \{x \mid \rho(Ax - b) \leq \epsilon\}$ . We begin by showing how linear and polar transformations on  $\mathcal{C}$  commute under certain assumptions.

**Proposition 3.3.1.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be finite-dimensional Euclidean spaces,  $\mathcal{C} \subseteq \mathcal{X}$  a closed, convex set which does not contain the origin,  $A : \mathcal{X} \rightarrow \mathcal{Y}$  a linear operator, and  $A^* : \mathcal{Y} \rightarrow \mathcal{X}$  its adjoint. Then*

$$(AC)' = (A^*)^{-1}\mathcal{C}'. \quad (3.3.1)$$

*Additionally, assume  $\mathcal{C}$  is polyhedral or  $\text{ri}(\mathcal{C}) \cap \text{range}(A) \neq \emptyset$ , and  $A^{-1}\mathcal{C}$  is not empty. Then the  $(A^{-1}\mathcal{C})'$  is nonempty and the following set equality holds*

$$(A^{-1}\mathcal{C})' = A^*\mathcal{C}'. \quad (3.3.2)$$

*Proof.* The first result is proved in [FMP14, Proposition 3.3] and the second in [FMP14, Proposition 3.4, 3.5].  $\square$

The previous propositions allows us to construct the antipolar of the constraint set  $\mathcal{C} = \{x \mid \rho(Ax - b) \leq \epsilon\}$  in Proposition 3.3.2.

**Proposition 3.3.2.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be finite-dimensional Euclidean spaces,  $A : \mathcal{X} \rightarrow \mathcal{Y}$  a linear operator, and  $A^* : \mathcal{Y} \rightarrow \mathcal{X}$  its adjoint. Let  $\mathcal{C} = \{x \mid \rho(Ax - b) \leq \epsilon\}$  with  $0 < \epsilon < \rho(b)$ . Also assume  $\text{ri}(\mathcal{C}) \cap \text{range}(A)$  and  $A^{-1}\mathcal{C}$  are not empty. Then*

$$\mathcal{C}' = \{A^*y \mid \langle b, y \rangle - \epsilon\rho^\circ(y) \geq 1\}. \quad (3.3.3)$$

*Proof.* Since the arguments of  $\rho$  lie in  $\mathcal{Y}$ , we first consider the antipolar of  $\mathcal{D} = A\mathcal{C} \subseteq \mathcal{Y}$  to establish the constraint  $\langle b, y \rangle - \epsilon\rho^\circ(y) \geq 1$ . Note that  $y \in \mathcal{D}'$  is equivalent to  $\langle Ax, y \rangle \geq 1$  for all  $Ax \in A\mathcal{C}$ . This is again equivalent to

$$\langle b - Ax, y \rangle \leq \langle b, y \rangle - 1 \quad \forall x \text{ such that } \rho(Ax - b) \leq \epsilon. \quad (3.3.4)$$

Next apply definition (3.2.14) for the antipolar  $\rho^\circ$  and take the supremum over  $u = \frac{b-Ax}{\epsilon}$ , giving

$$\begin{aligned} \epsilon\rho^\circ(y) &= \epsilon \sup_u \{ \langle u, y \rangle \mid \rho(u) \leq 1, u = \frac{b-Ax}{\epsilon} \} \\ &= \sup_x \{ \langle b - Ax, y \rangle \mid \rho\left(\frac{b-Ax}{\epsilon}\right) \leq 1 \} \\ &= \sup_x \{ \langle b - Ax, y \rangle \mid \rho(b - Ax) \leq \epsilon \}, \end{aligned} \quad (3.3.5)$$

where the last equality uses the postive homogeneity of the gauge  $\rho$ . Using equation (3.3.5), we see that equation (3.3.4) is equivalent to the desired constraint  $\epsilon\rho^\circ(y) \leq \langle b, y \rangle - 1$ . Thus  $\mathcal{D}' = \{y \mid \langle b, y \rangle - \epsilon\rho^\circ(y) \geq 1\}$ .

Finally, note that  $\mathcal{C}$  and  $\mathcal{C}'$  both lie in  $\mathcal{X}$ , while  $\mathcal{D} = A\mathcal{C}$  and  $\mathcal{D}'$  lie in  $\mathcal{Y}$ . Then by Proposition 3.3.1, the antipolar  $\mathcal{C}'$  has the form

$$\begin{aligned} \mathcal{C}' &= (A^{-1}\mathcal{D})' = A^*\mathcal{D}' \\ &= \{A^*y \mid \langle b, y \rangle - \epsilon\rho^\circ(y) \geq 1\}. \end{aligned} \quad (3.3.6)$$

$\square$

The next four propositions allow us to derive the inequality-constrained primal-gauge dual pair (3.2.21) from the linearly-constrained pair (3.2.20) by transforming the P-ineq model into a linearly-constrained gauge model of the form P-lin (see the dependency map (Figure 3.1) for reference). This process uses an indicator function to embed the constraint set  $\mathcal{C} = \{x \mid \rho(Ax - b) \leq \epsilon\}$  into the primal objective function, resulting in a P-lin model. Finally, we determine the gauge dual of the resulting model using polar relations and properties established in Section 3.2.

The first two propositions show that we may discuss the polar of a sum of gauges in terms of sets induced by Minkowski functions. This allows us to determine the polar of a sum of gauges in Proposition 3.3.5.

**Proposition 3.3.3.** *Let  $\mathcal{C} \subseteq \mathcal{X}$  be a closed, convex set containing the origin and  $\kappa = \gamma_{\mathcal{C}}$  the Minkowski function induced by  $\mathcal{C}$ . Then  $\kappa$  is a gauge,  $\mathcal{C} = \{x \mid \kappa(x) \leq 1\}$ , and  $\mathcal{C}$  is the unique closed, convex set containing the origin such that  $\kappa = \gamma_{\mathcal{C}}$ .*

*Proof.* To verify that  $\kappa$  is a gauge, first note that positive homogeneity and  $\kappa(0) = 0$  are direct results of  $\kappa$  being a Minkowski function. To show convexity, let  $x, y \in \mathcal{X}$  and  $0 \leq \alpha \leq 1$ . Then  $x \in \kappa(x)\mathcal{C}$  means  $\alpha x \in \alpha\kappa(x)\mathcal{C}$ , and likewise  $(1 - \alpha)y \in (1 - \alpha)\kappa(y)\mathcal{C}$ . Thus  $\alpha x + (1 - \alpha)y \in (\alpha\kappa(x) + (1 - \alpha)\kappa(y))\mathcal{C}$ . Then by the infimum of the Minkowski function,  $\kappa(\alpha x + (1 - \alpha)y) \leq \alpha\kappa(x) + (1 - \alpha)\kappa(y)$  and  $\kappa$  is convex.

Additionally,  $\kappa(x) \leq 1$  is equivalent to  $x \in \mathcal{C}$ , and thus  $\mathcal{C} = \{x \mid \kappa(x) \leq 1\}$ .

Finally, assume there is some closed, convex set  $\mathcal{D} \subseteq \mathcal{X}$  such that  $\kappa = \gamma_{\mathcal{D}}$ . Then  $\kappa(x) = \gamma_{\mathcal{C}}(x) = \gamma_{\mathcal{D}}(x) \leq 1$  is equivalent to  $x$  being in both  $\mathcal{C}$  and  $\mathcal{D}$ , since both sets are closed and convex. Likewise,  $\kappa(x) > 1$  indicates  $x$  is in neither set. Thus  $\mathcal{C} = \mathcal{D}$ .  $\square$

**Proposition 3.3.4.** *Let  $\kappa_1$  and  $\kappa_2$  be gauges. Let  $\kappa(x_1, x_2) = \kappa_1(x_1) + \kappa_2(x_2)$ ,  $\mathcal{C}_1 = \{z_1 \mid \kappa^\circ(z_1) \leq 1\}$ ,  $\mathcal{C}_2 = \{z_2 \mid \kappa^\circ(z_2) \leq 1\}$ , and  $\mathcal{C} = \{(z_1, z_2) \mid \kappa^\circ(z_1, z_2) \leq 1\}$ . Then  $\kappa$  and  $\kappa^\circ$  are gauges,  $\kappa^\circ = \gamma_{\mathcal{C}}$ , and  $\mathcal{C} = \mathcal{C}_1 \times \mathcal{C}_2$ .*

*Proof.* Since  $\kappa$  is the sum of gauges, it is also convex, nonnegative, positively homogeneous, and zero at the origin, and thus a gauge. Then  $\kappa^\circ$  is also a gauge and by Proposition 3.3.3,  $\mathcal{C}$  is the unique set such that  $\kappa = \gamma_{\mathcal{C}}$ . If  $(z_1, z_2) \in \mathcal{C}$ , then  $\langle x_1, z_1 \rangle + \langle x_2, z_2 \rangle \leq \kappa(x_1, x_2)\kappa^\circ(z_1, z_2) \leq \kappa_1(x_1) + \kappa_2(x_2)$  for all  $x_1 \in \text{dom } \kappa_1$  and  $x_2 \in \text{dom } \kappa_2$ . In particular, if  $x_2 = 0$  then  $\langle x_1, z_1 \rangle \leq \kappa_1(x_1)$  for all  $x_1 \in \text{dom } \kappa_1$ , indicating  $z_1 \in \mathcal{C}_1$ . Similarly  $z_2 \in \mathcal{C}_2$  and thus  $\mathcal{C}_1 \times \mathcal{C}_2 \subseteq \mathcal{C}$ . For the reverse inclusion,  $(z_1, z_2) \in \mathcal{C}_1 \times \mathcal{C}_2$  means  $\langle x_1, z_1 \rangle \leq \kappa_1(x_1)$  for all  $x_1 \in \text{dom } \kappa_1$  and  $\langle x_2, z_2 \rangle \leq \kappa_2(x_2)$  for all  $x_2 \in \text{dom } \kappa_2$ . Adding these inequalities, we have  $\langle x_1, z_1 \rangle + \langle x_2, z_2 \rangle \leq \kappa_1(x_1) + \kappa_2(x_2)$  for all  $x_1 \in \text{dom } \kappa_1$  and  $x_2 \in \text{dom } \kappa_2$ , and thus  $\mathcal{C} = \mathcal{C}_1 \times \mathcal{C}_2$ .  $\square$

Given the two propositions above, we now show that the polar of a sum of gauges is the max of the set of polars.

**Proposition 3.3.5.** *Let  $\kappa_1$  and  $\kappa_2$  be gauges. Then the sum  $\kappa(x_1, x_2) := \kappa_1(x_1) + \kappa_2(x_2)$  is a gauge and has the polar*

$$\kappa^\circ(z_1, z_2) = \max\{\kappa_1^\circ(z_1), \kappa_2^\circ(z_2)\}. \quad (3.3.7)$$

*Proof.* Proposition 3.3.4 shows  $\kappa$  and its polar  $\kappa^\circ$  are gauges. Setting  $\mathcal{C}_1 = \{z_1 \mid \kappa_1^\circ(z_1) \leq 1\}$ ,  $\mathcal{C}_2 = \{z_2 \mid \kappa_2^\circ(z_2) \leq 1\}$ , and  $\mathcal{C} = \{(z_1, z_2) \mid \kappa^\circ(z_1, z_2) \leq 1\}$ , we may express the gauge polars as Minkowski functions  $\kappa_1^\circ = \gamma_{\mathcal{C}_1}$ ,  $\kappa_2^\circ = \gamma_{\mathcal{C}_2}$ , and  $\kappa^\circ = \gamma_{\mathcal{C}}$ . Since  $\mathcal{C}_1$ ,  $\mathcal{C}_2$ , and  $\mathcal{C}$  are closed, convex sets containing the origin, Proposition 3.3.3 tells us these are the unique such sets defining  $\kappa_1^\circ$ ,  $\kappa_2^\circ$ , and  $\kappa^\circ$ . Additionally, Proposition 3.3.4 implies  $\mathcal{C} = \mathcal{C}_1 \times \mathcal{C}_2$  and  $\kappa^\circ(z_1, z_2) = \gamma_{\mathcal{C}}(z_1, z_2) = \gamma_{\mathcal{C}_1 \times \mathcal{C}_2}(z_1, z_2)$ .

Then we have

$$\begin{aligned}
\kappa^\circ(z_1, z_2) &= \gamma_{\mathcal{C}_1 \times \mathcal{C}_2}(z_1, z_2) \\
&= \inf \{ \lambda \geq 0 \mid z_1 \in \lambda \mathcal{C}_1, z_2 \in \lambda \mathcal{C}_2 \} \\
&= \max \{ \inf \{ \lambda \geq 0 \mid z_1 \in \lambda \mathcal{C}_1 \}, \inf \{ \lambda \geq 0 \mid z_2 \in \lambda \mathcal{C}_2 \} \} \\
&= \max \{ \gamma_{\mathcal{C}_1}(z_1), \gamma_{\mathcal{C}_2}(z_2) \} \\
&= \max \{ \kappa_1^\circ(z_1), \kappa_2^\circ(z_2) \}.
\end{aligned} \tag{3.3.8}$$

□

The following proposition establishes two equalities which allow us to compute the polar of an objective function which includes an indicator function. This strategy allows us to embed the inequality  $\rho(Ax - b) \leq \epsilon$  from (3.2.21) into the objective function of (3.2.20) in the proof of Theorem 3.3.1, (ii).

**Proposition 3.3.6.** *Let  $\rho$  be a gauge. Then for all  $y \in \mathcal{X}$  and  $\tau \geq 0$  the following equalities hold.*

$$(i) \quad (\delta_{\text{epi } \rho})^\circ(y, \tau) = \delta_{(\text{epi } \rho)^\circ}(y, \tau),$$

$$(ii) \quad \delta_{(\text{epi } \rho)^\circ}(y, \tau) = \delta_{\text{epi}(\rho^\circ)}(y, -\tau).$$

*Proof.* To show (i) holds, consider the expansion of the expressions

$$\begin{aligned}
(\delta_{\text{epi } \rho})^\circ(y, \tau) &= \inf \{ \mu > 0 \mid \langle (x, \sigma), (y, \tau) \rangle \leq \mu \delta_{\text{epi } \rho}(x, \sigma) \quad \forall (x, \sigma) \}, \\
\delta_{(\text{epi } \rho)^\circ}(y, \tau) &= \begin{cases} 0 & \text{if } \langle (x, \sigma), (y, \tau) \rangle \leq 1 \quad \forall (x, \sigma) \in \text{epi } \rho \\ +\infty & \text{else.} \end{cases}
\end{aligned}$$

Note that the value of  $\delta_{(\text{epi } \rho)^\circ}(y, \tau)$  is either 0 or  $+\infty$ . If  $\delta_{(\text{epi } \rho)^\circ}(y, \tau) = 0$ , then  $\langle (x, \sigma), (y, \tau) \rangle \leq 1$  for all  $(x, \sigma) \in \text{epi } \rho$ . Since  $\rho$  is a gauge,  $(x, \sigma) \in \text{epi } \rho$  is equivalent to  $\rho(\alpha x) \leq \alpha \sigma$  for all  $\alpha > 0$ . Then  $(\alpha x, \alpha \sigma) \in \text{epi } \rho$  for all  $\alpha > 0$ . Dividing the inequality by  $\alpha$  and taking the limit, we have

$$\langle (x, \sigma), (y, \tau) \rangle \leq \lim_{\alpha \rightarrow +\infty} \frac{1}{\alpha} = 0 \quad \text{for all } (x, \alpha) \in \text{epi } \rho.$$

Thus  $(\delta_{\text{epi } \rho})^\circ(y, \tau) = 0$ . On the other hand, if  $\delta_{(\text{epi } \rho)^\circ}(y, \tau) = +\infty$  then there is some  $(x, \sigma) \in \text{epi } \rho$  with  $\langle (x, \sigma), (y, \tau) \rangle > 1$ . Then  $\delta_{\text{epi } \rho}(x, \sigma) = 0$  and there is no  $\mu > 0$  such that  $\langle (x, \sigma), (y, \tau) \rangle \leq \mu \delta_{\text{epi } \rho}(x, \sigma)$ . Thus  $(\delta_{\text{epi } \rho})^\circ(y, \tau) = +\infty$ .

To show (ii) holds, we have the following set of equivalences

$$\begin{aligned}
(y, \tau) \in (\text{epi } \rho)^\circ &\iff \langle (x, \sigma), (y, \tau) \rangle \leq 0 \quad \forall x \text{ and } \rho(x) \leq \sigma \\
&\iff \langle x, y \rangle \leq -\tau \rho(x) \quad \forall x \\
&\iff \rho^\circ(y) = \inf \{ \mu > 0 \mid \langle x, y \rangle \leq \mu \rho(x) \quad \forall x \} \leq -\tau \\
&\iff (y, -\tau) \in \text{epi}(\rho^\circ).
\end{aligned}$$

Here, the first equivalence was established in the proof of (i) and the second comes from setting  $\sigma$  as the minimal value  $\sigma = \rho(x)$ . Thus  $\delta_{(\text{epi } \rho)^\circ}(y, \tau) = \delta_{\text{epi}(\rho^\circ)}(y, -\tau)$ . □

We are now prepared to prove the essential results which underly our phase retrieval method and signal recovery strategy. We begin with the following gauge duality result.

**Theorem 3.3.1.** *Let P-ineq and GD-ineq be the inequality-constrained pair of models from (3.2.21). Also let  $\mathcal{C} = \{x \mid \rho(Ax - b) \leq \epsilon\}$  with  $0 < \epsilon < \rho(b)$  and assume  $\text{ri}(\mathcal{C}) \cap \text{range}(A)$  and  $A^{-1}\mathcal{C}$  are not empty. Then the model GD-ineq is the gauge dual of the primal P-ineq based on the gauge duality of*

- (i) the nonlinear gauge dual pair (3.2.2), and
- (ii) the linear gauge dual pair (3.2.20).

*Proof.* Item (i) is a direct result of Proposition 3.3.2, which gives the antipolar set  $\mathcal{C}' = \{A^*y \mid \langle b, y \rangle - \epsilon\rho^\circ(y) \geq 1\}$ . As a result, the argument of the GD-nonline objective  $\kappa^\circ$  is  $A^*y$  and GD-ineq is equivalent to GD-nonline.

To prove item (ii), we must first express P-ineq as a linear gauge model of the form P-lin. Define the function  $\phi(x, r, \sigma) = \kappa(x) + \delta_{\text{epi}\rho}(r, \sigma)$ . Since the epigraph of a gauge is closed under positive scaling, the indicator function  $\delta_{\text{epi}\rho}$  is a gauge. By Proposition 3.3.4,  $\phi$  is a gauge, and thus P-ineq is equivalent to the linear gauge model

$$\begin{aligned} \min_{x, r, \sigma} \quad & \phi(x, r, \sigma) \\ \text{s.t.} \quad & r = b - Ax \\ & \sigma = \epsilon. \end{aligned} \tag{3.3.9}$$

To combine these linear constraints, define the following extended matrix and vectors

$$\bar{A} = \begin{bmatrix} A & I & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \bar{x} = \begin{bmatrix} x \\ r \\ \sigma \end{bmatrix}, \quad \bar{b} = \begin{bmatrix} b \\ \epsilon \end{bmatrix}, \tag{3.3.10}$$

and define the spaces  $\bar{\mathcal{X}} = \mathcal{X} \times \mathcal{Y} \times \{\mathbb{R} \cup +\infty\}$  and  $\bar{\mathcal{Y}} = \mathcal{Y} \times \{\mathbb{R} \cup +\infty\}$ .

Then (3.3.9) has the linear form equivalent to P-lin

$$\min_{\bar{x} \in \bar{\mathcal{X}}} \phi(\bar{x}) \quad \text{s.t.} \quad \bar{A}\bar{x} = \bar{b}. \tag{3.3.11}$$

This model has the following gauge dual per (3.2.20)

$$\min_{\bar{y} \in \bar{\mathcal{Y}}} \phi^\circ(\bar{A}^* \bar{y}) \quad \text{s.t.} \quad \langle \bar{b}, \bar{y} \rangle = 1, \tag{3.3.12}$$

where  $\bar{A}^* \bar{y} = (A^*y, y, \tau)$  and  $\langle \bar{b}, \bar{y} \rangle = \langle b, y \rangle + \sigma\tau$ .

Taking the polar of  $\phi$ , we have

$$\begin{aligned} \phi^\circ(A^*y, y, \tau) &= \max \{ \kappa^\circ(A^*y), (\delta_{\text{epi}\rho})^\circ(y, \tau) \} \\ &= \max \{ \kappa^\circ(A^*y), \delta_{(\text{epi}\rho)^\circ}(y, \tau) \} \\ &= \kappa^\circ(A^*y) + \delta_{(\text{epi}\rho)^\circ}(y, \tau) \\ &= \kappa^\circ(A^*y) + \delta_{\text{epi}(\rho^\circ)}(y, -\tau). \end{aligned} \tag{3.3.13}$$

The first equality is a result of Proposition 3.3.5. The second and last equalities are results of Proposition 3.3.6, (i) and (ii), respectively. And the third equality is a consequence of the indicator function mapping to  $\{0, +\infty\}$ .

The indicator function  $\delta_{\text{epi}(\rho^\circ)}(y, -\tau)$  corresponds to the constraint  $\rho^\circ(y) \leq -\tau$ . This constraint and the equality constraint  $\langle b, y \rangle + \sigma\tau = 1$  may be combined as

$$\langle b, y \rangle - \sigma\rho^\circ(y) \geq \langle b, y \rangle + \sigma\tau = 1.$$

Thus we eliminate  $\tau$  and recover GD-ineq.  $\square$

With the gauge duality of (3.2.21) established, we proceed by establishing weak duality, strong duality, and optimality conditions for this pair of models. These properties will play a central role in the first-order method for optimizing the PLGD model, allowing us to develop termination conditions as well as a recovery method for the primal signal  $x$  from a dual iterate  $y$ . The weak duality property of (3.2.21) is a straightforward consequence of the definition of the polar of a function and inequality constraints in this model.

**Proposition 3.3.7.** (weak duality) *Assume the primal model P-ineq in (3.2.21) is feasible and  $0 \leq \epsilon < \rho(b)$ . Then for all primal and gauge dual feasible pairs  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  we have*

$$\kappa(x)\kappa^\circ(\mathcal{A}^*y) \geq 1. \quad (3.3.14)$$

*Proof.* Since  $x$  and  $y$  are feasible for (3.2.21), we have

$$\begin{aligned} 1 &\leq \langle b, y \rangle - \epsilon\rho^\circ(y) \\ &\leq \langle b, y \rangle - \rho(b - Ax)\rho^\circ(y) \\ &\leq \langle b, y \rangle - \langle b - Ax, y \rangle \\ &= \langle x, A^*y \rangle \\ &\leq \kappa(x)\kappa^\circ(A^*y). \end{aligned} \quad (3.3.15)$$

Here, the first and second inequalities are due to GD-ineq and P-ineq feasibility, respectively. The third and fourth inequalities are a consequence of polar functions.  $\square$

In order to show strong duality holds for the pair (3.2.21), we appeal to the strong duality properties of the Lagrange primal and dual pair (see [Roc70, Section 28]).

$$\begin{array}{ll} \min_{x \in \mathcal{X}} & \kappa(x) \\ \text{(P-ineq) s.t.} & \rho(Ax - b) \leq \epsilon \end{array} \quad \begin{array}{ll} \max_{y \in \mathcal{Y}} & \langle b, y \rangle - \epsilon\rho^\circ(y) \\ \text{(LD-ineq) s.t.} & \kappa^\circ(A^*y) \leq 1. \end{array} \quad (3.3.16)$$

This proof strategy also requires that  $\rho^\circ$  be continuous in order to identify an appropriate Lagrange dual feasible sequence which may be rescaled to a gauge dual feasible sequence.

**Proposition 3.3.8.** (strong duality) *Assume the primal model P-ineq in (3.2.21) is feasible,  $\rho^\circ$  is continuous, and  $0 \leq \epsilon < \rho(b)$ . Then P-ineq admits an optimal variable  $x_\star$  and*

$$\kappa(x_\star)\nu_{GD} = 1 \quad (3.3.17)$$

where  $\nu_{GD}$  is the optimal GD-ineq value

$$\nu_{GD} := \inf_{\langle b, y \rangle - \epsilon \rho^\circ(y) \geq 1} \kappa^\circ(A^*y). \quad (3.3.18)$$

Additionally, if P-ineq is strictly feasible then P-ineq and GD-ineq admit an optimal pair  $(x_\star, y_\star) \in \mathcal{X} \times \mathcal{Y}$  and

$$\kappa(x_\star)\kappa^\circ(A^*y_\star) = 1, \quad (3.3.19)$$

where  $y_\star$  is a rescaling of the optimal LD-ineq variable  $\hat{y}$  such that

$$y_\star = \frac{\hat{y}}{\langle b, \hat{y} \rangle - \epsilon \rho^\circ(\hat{y})}. \quad (3.3.20)$$

*Proof.* Since P-ineq is feasible, LD-ineq is bounded above. Additionally, since LD-ineq is strictly feasible for  $y = 0$  (i.e.,  $\kappa^\circ(0) = 0 < 1$  and 0 is in the relative interior of the Lagrange dual feasible set), [Roc70, Theorem 28.2] indicates that P-ineq admits an optimal variable  $\hat{x}$  and LD-ineq has finite optimal objective value

$$\nu_{LD} := \sup_{\kappa^\circ(A^*y) \leq 1} \langle b, y \rangle - \epsilon \rho^\circ(y) < +\infty.$$

Furthermore, [Roc70, Theorem 28.4] tells us this pair has zero Lagrange duality gap.

Since  $\hat{x}$  is optimal for P-ineq,  $x_\star = \hat{x}$  is our desired primal solution. By strong Lagrange duality  $\nu_{LD} = \kappa(x_\star)$ , and by weak gauge duality (Proposition 3.3.7) this value is strictly greater than zero. Let  $\{y_i\}$  be a LD-ineq feasible sequence such that  $\langle b, y_i \rangle - \epsilon \rho^\circ(y_i) \rightarrow \nu_{LD} > 0$ . Since  $\rho^\circ$  is continuous, there exists a subsequence  $\{y_{i_j}\}$  such that  $\langle b, y_{i_j} \rangle - \epsilon \rho^\circ(y_{i_j})$  is bounded above zero for all  $j$ . Rescaling this sequence such that

$$\bar{y}_j = \frac{y_{i_j}}{\langle b, y_{i_j} \rangle - \epsilon \rho^\circ(y_{i_j})},$$

we have  $\langle b, \bar{y}_j \rangle - \epsilon \rho^\circ(\bar{y}_j) = 1$  for all  $j$ , and thus  $\{\bar{y}_j\}$  is a GD-ineq feasible sequence. Then we have

$$\begin{aligned} \nu_{GD} &\leq \lim_{j \rightarrow \infty} \kappa^\circ(A^*\bar{y}_j) \\ &= \lim_{j \rightarrow \infty} \frac{1}{\langle b, y_{i_j} \rangle - \epsilon \rho^\circ(y_{i_j})} \kappa^\circ(A^*y_{i_j}) \\ &= \frac{1}{\kappa(x_\star)} \cdot \lim_{j \rightarrow \infty} \kappa^\circ(A^*y_{i_j}) \leq \frac{1}{\kappa(x_\star)}, \end{aligned}$$

where the last inequality is due to  $\{y_{i_j}\}$  being LD-ineq feasible (i.e.,  $\kappa^\circ(A^*y_{i_j}) \leq 1$  for all  $j$ ). And by weak gauge duality (Proposition 3.3.7) we have  $\nu_{GD} \geq \frac{1}{\kappa(x_\star)}$ , giving (3.3.17).

If P-ineq is strictly feasible, then by [Roc70, Theorem 28.2] LD-ineq will admit a finite optimal solution  $\hat{y}$ . This variable may be rescaled to obtain the optimal GD-ineq variable (3.3.20) and the same subsequence argument gives (3.3.19).  $\square$

We close this section by establishing the optimality conditions of (3.2.21) which are primarily a consequence of the strong and weak duality results above.

**Proposition 3.3.9.** (optimality conditions) *If the primal model P-ineq in (3.2.21) is feasible,  $\rho^\circ$  is continuous, and  $0 \leq \epsilon < \rho(b)$ , then the pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is primal-dual optimal if and only if the following conditions hold:*

$$\rho(Ax - b) = \epsilon \quad \text{primal feasibility} \quad (3.3.21a)$$

$$\langle b, y \rangle - \epsilon \rho^\circ(y) = 1 \quad \text{dual feasibility} \quad (3.3.21b)$$

$$\langle b - Ax, y \rangle = \rho(b - Ax) \rho^\circ(y) \quad \text{complementarity} \quad (3.3.21c)$$

$$\langle x, A^*y \rangle = \kappa(x) \kappa^\circ(A^*y) = 1 \quad \text{strong duality} \quad (3.3.21d)$$

$$\frac{1}{\rho^\circ(y)} y \in \partial \rho(\cdot)(b - Ax) \quad \text{primal subdifferential} \quad (3.3.21e)$$

$$\frac{1}{\epsilon} (b - Ax) \in \partial \rho^\circ(\cdot)(y) \quad \text{dual subdifferential} \quad (3.3.21f)$$

*Proof.* If these conditions hold, then  $(x, y)$  are primal-dual feasible and have zero duality gap (i.e.,  $\kappa(x) \kappa^\circ(A^*y) = 1$ ). Thus by strong duality (Proposition 3.3.8)  $(x, y)$  are an optimal pair.

Likewise, if  $(x, y)$  are an optimal pair, then strong duality implies  $\kappa(x) \kappa^\circ(A^*y) = 1$ . Thus the set of inequalities in the weak duality proof (3.3.15) are tight and the first four conditions hold.

The last two conditions are a consequence of the first four conditions and polar relations. To show the primal subdifferential condition, note that for all  $z$  we have  $\langle z, y \rangle \leq \rho(z) \rho^\circ(y)$ . The primal feasibility and complementarity conditions give  $\langle b - Ax, y \rangle = \epsilon \rho^\circ(y)$ . Then for all  $z$  we have

$$\langle z - (b - Ax), y \rangle \leq \rho(z) \rho^\circ(y) - \epsilon \rho^\circ(y).$$

Replacing  $\epsilon = \rho(b - Ax)$  and dividing by  $\rho^\circ(y)$ , we have for all  $z$

$$\rho(z) \geq \rho(b - Ax) + \langle z - (b - Ax), y / \rho^\circ(y) \rangle,$$

and thus  $\frac{1}{\rho^\circ(y)} y \in \partial \rho(\cdot)(b - Ax)$ .

For the dual subdifferential condition, we have  $\langle b - Ax, z \rangle \leq \rho(b - Ax) \rho^\circ(z) = \epsilon \rho^\circ(z)$  for all  $z$ . Subtracting  $\langle b - Ax, y \rangle = \epsilon \rho^\circ(y)$  gives  $\langle b - Ax, z - y \rangle \leq \epsilon \rho^\circ(z) - \epsilon \rho^\circ(y)$  for all  $z$ , or

$$\rho^\circ(z) \geq \rho^\circ(y) + \langle (b - Ax) / \epsilon, z - y \rangle,$$



and thus  $\frac{1}{\epsilon}(b - Ax) \in \partial\rho^\circ(\cdot)(y)$ .

□

### 3.4 Optimality conditions and primal recovery for the PLGD model

This section establishes optimality conditions and a primal signal recovery strategy for the PLGD model (3.2.1), two essential tools for developing the method presented in Chapter 4. We see first that the duality of the PLP-PLGD pair (3.2.1) is a direct result of the duality previously established for the inequality-constrained primal-dual model (3.2.21). Thus Proposition 3.3.9 applies to the PLGD model, supporting the development of termination conditions and a primal signal recovery method for any given PLGD optimization method. We also present a more general optimality condition which applies to the first-order method established in the next chapter.

We begin by establishing the duality of the PLP-PLGD pair (3.2.1) based on the duality established in Theorem 3.3.1. In the PLP model, the linear operator  $\mathcal{A}$  is defined as a map from the space of Hermitian matrices  $\mathcal{H}^n$  to  $\mathbb{R}^m$  (see (2.3.6) for details). We may pass the PLP constraint  $X \succeq 0$  into the objective by setting  $\kappa(X) := \|X\|_1 + \delta_{(\cdot) \succeq 0}(X)$ . The polar of  $\kappa$  has the form

$$\kappa^\circ(Z) = \inf \{ \mu > 0 \mid \langle X, Z \rangle \leq \mu \|X\|_1 \ \forall X \succeq 0 \},$$

which can be simplified by considering two cases. If  $\lambda_1(Z)$  is positive then  $\kappa^\circ(Z)$  operates as the dual norm of the Schatten 1-norm restricted to the positive eigenspace of  $Z$ . Otherwise, if  $Z \preceq 0$  then  $\kappa^\circ(Z)$  is zero, giving

$$\kappa^\circ(Z) = \begin{cases} \lambda_1(Z) & \text{if } Z \succeq 0 \\ 0 & \text{else.} \end{cases} \quad (3.4.1)$$

Next we set  $\rho(y) := \|y\| = \|y\|_2$ , which has the polar (or dual norm)  $\rho^\circ(y) = \|y\|_* = \|y\|$ . Thus PLP has the form of P-ineq (3.2.21)

$$\begin{aligned} \min_{X \in \mathcal{H}^n} \quad & \kappa(X) = \|X\|_1 + \delta_{(\cdot) \succeq 0}(X) \\ \text{s.t.} \quad & \rho(\mathcal{A}(X) - b) = \|\mathcal{A}(X) - b\| \leq \epsilon. \end{aligned}$$

Then by Theorem 3.3.1, the GD-ineq model has the objective  $\kappa^\circ(\mathcal{A}^*y)$  and constraint  $\langle b, y \rangle - \epsilon \|y\| \geq 1$ . We may further simplify  $\kappa^\circ$  by noting that  $\kappa(X)$  is strictly positive and finite for all feasible  $X \in \mathcal{H}^n$ . Then by weak duality (Proposition 3.3.7)  $\kappa^\circ(\mathcal{A}^*y)$  is likewise strictly positive and finite for all feasible  $y$ . Thus  $\kappa^\circ(\mathcal{A}^*y) = \lambda_1(\mathcal{A}^*y)$  and we recover the PLGD model from the GD-ineq model.

Next we see that the PLP-PLGD (3.2.1) optimality conditions and primal recovery property are a direct result of the propositions in Section 3.3 applied to primal-gauge dual

space  $\mathcal{X} \times \mathcal{Y} = \mathcal{H}^n \times \mathbb{R}^m$ . These two corollaries rely on the von Neumann trace inequality, which states that for all  $n \times n$  complex matrices  $A$  and  $B$ ,

$$\langle A, B \rangle \leq \sum_{i=1}^n \sigma_i(A) \sigma_i(B), \quad (3.4.2)$$

and equality holds if and only if  $A$  and  $B$  admit a simultaneous ordered SVD [Gri91].

**Corollary 3.4.1.** (PLP-PLGD optimality conditions) *If PLP in (3.2.1) is feasible and  $0 \leq \epsilon < \|b\|$ , then  $(X, y) \in \mathcal{H}^n \times \mathbb{R}^m$  is primal-dual optimal if and only if the following conditions hold*

- (a)  $X \succeq 0$  and  $\|\mathcal{A}(X) - b\| = \epsilon$ ,
- (b)  $\langle b, y \rangle - \epsilon \|y\| = 1$ ,
- (c)  $\langle b - \mathcal{A}(X), y \rangle = \|\mathcal{A}(X) - b\| \cdot \|y\|$ ,
- (d)  $\langle X, \mathcal{A}^* y \rangle = \|X\|_1 \cdot \lambda_1(\mathcal{A}^* y) = 1$ ,
- (e)  $\lambda_i(X) \cdot (\lambda_1(\mathcal{A}^* y) - \lambda_i(\mathcal{A}^* y)) = 0$  for all  $i = 1, \dots, n$ ;
- (f)  $X$  and  $\mathcal{A}^* y$  admit a simultaneous ordered eigendecomposition.

*Proof.* Since  $\rho^\circ(\cdot) = \|\cdot\|$  is continuous, we may invoke Proposition 3.3.9. The first four conditions are identical to those discussed in Proposition 3.3.9 for the more general model (3.2.21). Thus these conditions holding is equivalent to  $(X, y)$  being primal-dual optimal.

Then we must simply show the first four conditions imply the last two. Applying the von Neumann trace inequality to the matrices  $X$  and  $\mathcal{A}^* y$ , we have

$$1 = \langle X, \mathcal{A}^* y \rangle \leq \sum_{i=1}^n \sigma_i(X) \sigma_i(\mathcal{A}^* y) \leq \|X\|_1 \cdot \lambda_1(\mathcal{A}^* y) = 1.$$

Thus for all  $i$ , if  $\lambda_i(X) > 0$  then we must have  $\lambda_i(\mathcal{A}^* y) = \lambda_1(\mathcal{A}^* y)$  and the matrices  $X$  and  $\mathcal{A}^* y$  admit a simultaneous ordered eigendecomposition.  $\square$

In addition to Corollary 3.4.1, the following first-order optimality condition for convex optimization also applies to the PLGD model (3.2.1). Recall that convex optimization is the minimization of any convex function  $f$  over a convex set  $\mathcal{C}$ , and a first-order method is any iteration  $y_+ = \Pi_{\mathcal{C}}(y - \alpha g)$ , where  $g$  is some descent direction in  $\partial f(y)$  and  $\alpha > 0$  is a steplength chosen with some linesearch method. In the next chapter we will develop a first-order method for optimizing (3.2.1), but for now consider any first-order PLGD method for  $f(y) = \lambda_1(\mathcal{A}^* y)$  and  $\mathcal{C} = \{y \in \mathbb{R}^m \mid \langle b, y \rangle - \epsilon \|y\| \geq 1\}$ .

The following result can be viewed as a fixed-point optimality condition, where  $y$  is optimal if  $y = \Pi_{\mathcal{C}}(y - \alpha g)$  for some  $g \in \partial f(y)$  and all  $\alpha > 0$ . To prove this result, we frame

(3.2.1) as an unconstrained convex optimization problem by defining  $F(y) := f(y) + \delta_{\mathcal{C}}(y)$ , where  $\delta_{\mathcal{C}}(y)$  is the indicator function (1.2.6) of  $\mathcal{C}$ . Then optimizing (3.2.1) is equivalent to minimizing the unconstrained function  $F$  over  $\mathbb{R}^m$ .

**Proposition 3.4.2.** *Let  $f$  be a convex, subdifferentiable function over a finite-dimensional Euclidean space  $\mathcal{Y}$  and  $\mathcal{C} \subseteq \mathcal{Y}$  a convex set. Also assume  $\text{ri}(\text{dom}(f)) \cap \text{ri}(\mathcal{C})$  is not empty. Then  $y_0$  is optimal for the minimization problem*

$$\min_y F(y) := f(y) + \delta_{\mathcal{C}}(y) \quad (3.4.3)$$

if  $y_0 = \Pi_{\mathcal{C}}(y_0 - \alpha g)$  for some nonzero  $g \in \partial f(y_0)$  and all  $\alpha > 0$ .

*Proof.* We begin by demonstrating the first-order optimality condition for unconstrained convex subdifferentiable functions. A variable  $y_*$  is optimal for (3.4.3) if and only if

$$F(y_*) = F(y_*) + 0^T(y - y_*) \leq F(y) \quad \forall y \in \mathcal{Y}. \quad (3.4.4)$$

Then zero is in the subdifferential of  $F$  at  $y_*$ , i.e.,

$$y_* \text{ is optimal if and only if } 0 \in \partial F(y_*). \quad (3.4.5)$$

Thus it suffices to show  $0 \in \partial F(y_0)$ . Note that  $-g$  is not a descent direction for  $f$ , since  $f(y_0) = f(\Pi_{\mathcal{C}}(y_0 - \alpha g))$  for all  $\alpha > 0$ . Then  $-g$  is in the normal cone (1.2.7) of  $\mathcal{C}$  at  $y_0$ . Additionally, the normal cone  $N_{\mathcal{C}}(y_0)$  is equivalent to the subdifferential of the indicator function  $\partial \delta_{\mathcal{C}}(y_0)$ , since

$$\begin{aligned} \partial \delta_{\mathcal{C}}(y_0) &= \{g \mid \delta_{\mathcal{C}}(y) \geq \delta_{\mathcal{C}}(y_0) + \langle g, y - y_0 \rangle \quad \forall y\} \\ &= \{g \mid 0 \geq \langle g, y - y_0 \rangle \quad \forall y \in \mathcal{C}\} \\ &= N_{\mathcal{C}}(y_0). \end{aligned} \quad (3.4.6)$$

Finally, since  $f$  and  $\mathcal{C}$  are convex and  $\text{ri}(\text{dom}(f)) \cap \text{ri}(\mathcal{C}) \neq \emptyset$ , we have the set equality  $\partial F(y_0) = \partial f(y_0) + \partial \delta_{\mathcal{C}}(y_0)$  [Roc70, Theorem 23.8]. Then  $g \in \partial f(y_0)$  and  $-g \in \partial \delta_{\mathcal{C}}(y_0)$  imply that  $0 \in \partial F(y_0)$ , and thus  $y_0$  is optimal for (3.4.3). □

Thus  $y = \Pi_{\mathcal{C}}(y - \alpha g)$  is another optimality condition for the PLGD model (3.2.1).

Next we apply Corollary 3.4.1 to establish a primal recovery strategy for (3.2.1).

**Corollary 3.4.3.** (PLP-PLGD primal recovery)

*Assume the optimality conditions of Corollary 3.4.1 hold. Let  $y \in \mathbb{R}^m$  be an arbitrary optimal solution for (3.2.1),  $U \in \mathbb{C}^{n \times r}$  the eigenvectors corresponding to the algebraically largest eigenvalue  $\lambda_1$  of  $\mathcal{A}^*y$ , and  $r$  the multiplicity of  $\lambda_1$ . Then  $X \in \mathbb{C}^{n \times n}$  is a solution to (3.2.1) if and only if there exists an  $r \times r$  matrix  $S \succeq 0$  such that*

(a)  $X = USU^*$ ,

(b)  $b - \mathcal{A}(X) \in \epsilon \partial \|\cdot\|(y)$ .

*Proof.* If  $X$  is a solution to (3.2.1), then (e) of Corollary 3.4.1 implies that rank of  $X$  is no larger than the multiplicity  $r$  of  $\lambda_1(\mathcal{A}^*y)$ , and (f) implies that  $X$  has the form  $X = USU^*$  for an  $r \times r$  matrix  $S \succeq 0$ . Also, the dual subdifferential condition of Proposition 3.3.9 implies that (b) holds.

Conversely, if (a) and (b) hold then it can be show that the optimality conditions (a)-(d) of Corollary 3.4.1 also hold. The optimality of  $y$  gives  $\langle b, y \rangle - \epsilon \|y\| = 1$ . The next two conditions may be derived from definitions of the dual norm  $\rho^\circ(\cdot) = \|\cdot\|$ , which we temporarily denote as  $\|\cdot\|_*$  for clarity. The subdifferential of the dual norm (3.2.7) is the set of maximizers which attain the dual norm value, that is  $\partial \|\cdot\|_*(y) = \{x \mid \|y\|_* = \langle x, y \rangle\}$ . Then we have

$$\langle b - \mathcal{A}(X), y \rangle = \epsilon \|y\|_*. \quad (3.4.7)$$

Like the polar (3.2.5), the dual norm  $\|\cdot\|_*$  may also be defined as the function which most tightly satisfies the Cauchy-Schwarz inequality  $|\langle x, y \rangle| \leq \|x\| \cdot \|y\|_*$ . This definition gives

$$\langle b - \mathcal{A}(X), y \rangle = \|b - \mathcal{A}(X)\| \cdot \|y\|_*. \quad (3.4.8)$$

The pair (3.4.7) and (3.4.8) give optimality conditions (a) and (c).

The following equality establishes optimality condition (d):

$$1 = \langle X, \mathcal{A}^*y \rangle = \sum_{i=1}^n \sigma_i(X) \sigma_i(\mathcal{A}^*y) = \|X\|_1 \cdot \lambda_1(\mathcal{A}^*y).$$

In this expression, the first equality is a result of optimality conditions (a)-(c), which cause the first four lines of (3.3.15) to hold with equality. The second equality is a consequence of the von Neumann trace inequality (3.4.2) holding tightly. The columns of  $U \in \mathbb{C}^{n \times r}$  span the eigenspace of the algebraically largest eigenvalue  $\lambda_1$  of  $\mathcal{A}^*y$ . Since the diagonalization of  $X = USU^*$  only requires a unitary transformation of  $U$ , this transformation will not affect the eigenspace of  $\mathcal{A}^*y$ . Thus  $X$  and  $\mathcal{A}^*y$  admit a simultaneous ordered SVD. Finally, the third equality is a result of  $X$  having rank at most  $r$  and  $\lambda_1$  having multiplicity  $r$ .

Thus optimality conditions (a)-(d) hold, and  $X$  is a solution to (3.2.1). □

This primal recovery property allows us to develop a recovery strategy for obtaining a signal  $x$  from a dual variable  $y$ . If  $y \neq 0$  then the 2-norm used in the PLP-PLGD pair (3.2.1) gives  $\partial \|\cdot\|(y) = \nabla \|\cdot\|_2(y) = y/\|y\|$ . For a given PLGD model, if the lifted true signal  $X = \bar{x}\bar{x}^*$  is in the set of optimal matrices, then Corollary 3.4.3 indicates that the corresponding dual optimal variable  $y_\star$  will have the relation

$$\eta = (\bar{b} + \eta) - \bar{b} = b - \mathcal{A}(X_\star) = \epsilon \nabla \|\cdot\|_2(y_\star) = \epsilon \frac{y_\star}{\|y_\star\|}. \quad (3.4.9)$$

Thus the noise term  $\eta$  will be in the set of rescaled optimal dual variables.

In the case of noisy phase retrieval (2.1.1), there is no guarantee that the lifted true signal  $X = \bar{x}\bar{x}^*$  will be an optimal matrix for the PLGD model (3.2.1). In this case, the gauge dual variable  $y$  can be viewed instead as a parameter which learns the noise term  $\eta$  with some degree of inaccuracy. As we will see in Section 5.3, the tendency of  $y$  to learn the noise term persists even when the rank of the optimal matrix is greater than one and a given first-order method is unable to converge to  $y_*$ .

Given the optimality and primal recovery properties established above, we now proceed to develop a first-order method for optimizing the PLGD model (3.2.1) and recovering a primal signal  $x$  from the PLGD solution  $y$ .

## Chapter 4

# A first-order method for the PLGD model

### 4.1 Introduction

In this chapter we present a first-order algorithm for optimizing the PLGD model (3.2.1). Section 4.2 develops the computational steps necessary for a first-order PLGD algorithm. We then present the resulting algorithm and discuss specific implementation details. Section 4.3 demonstrates the effectiveness of this algorithm for noiseless phase retrieval problems. We close with a summary of the challenges to this algorithm which are addressed in the remainder of this dissertation.

The algorithm presented in this chapter was first established in [FM16]. As we will see in Section 5.3, the challenges posed by noisy phase retrieval are intrinsic to the PLGD model (3.2.1) itself, and independent of the choice of first-order method. Thus, due to the effectiveness of this particular method for noiseless phase retrieval and the existence of a comprehensive, specialized software package, we use this method to examine the behavior of first-order methods for optimizing (3.2.1). All implementation details in this chapter are identical to the original software package unless otherwise noted.

### 4.2 Algorithm and implementation details

We begin by examining the features of the PLGD model (3.2.1) which are essential to developing a first-order descent method. Recall that first-order methods like gradient descent have a basic iterate update  $y_+ = \Pi_{\mathcal{C}}(y - \alpha g)$ , where  $-g$  is a descent direction based on first-order information,  $\alpha$  is a steplength determined with some policy, and  $\Pi_{\mathcal{C}}$  is projection onto the feasible region  $\mathcal{C}$ . Since we are primarily concerned with large-scale signal recovery, the descent direction must be constructed using only first-order information. Each objective function evaluation  $\lambda_1(\mathcal{A}^*y)$  requires the solution of a costly eigenvalue problem, so its com-

putation should be kept to a minimum. Additionally, the PLGD model (3.2.1) is a convex problem, having both convex objective function  $\lambda_1(\mathcal{A}^*y)$  and convex constraint domain  $\mathcal{C} = \{y \in \mathbb{R}^m \mid \langle b, y \rangle - \epsilon \|y\| \geq 1\}$ . Thus, any descent method should take advantage of this convexity and determine the steplength  $\alpha$  using a linesearch with minimal backtracking (or evaluations of  $\lambda_1(\mathcal{A}^*y)$ ) and guaranteed convergence. Some applicable methods include projected gradient, quasi-Newton, and spectral bundle methods. Our first-order method of choice for solving (3.2.1) is a projected gradient descent method with adaptive steplength based on the local differentiability of  $\lambda_1(\mathcal{A}^*y)$ .

Construction of this projected gradient method requires the following sequence of computational steps. First, a descent direction is chosen from the gradient or subdifferential (1.2.12) of the objective function. Next, an initial steplength is computed and used to initialize a linesearch (backtracking) method for determining the update  $y_+$ . The linesearch and update both require a method for projecting onto the feasible region  $\mathcal{C}$ . Finally, a recovery method is used to recover the primal signal update  $x_+$  from the dual update  $y_+$ .

The steps described above are sufficient for a projected gradient method. However, the authors of [FM16] found that two additional refinement steps tend to accelerate convergence greatly for this descent method. Thus our algorithm follows the primal signal recovery with a primal refinement step and a dual refinement step.

To determine the descent vector, we consider the subdifferential of the function  $\lambda_1(\mathcal{A}^*y)$

$$\partial\lambda_1(\mathcal{A}^*y) = \{\mathcal{A}(V_1TV_1^*) \mid T \succeq 0, \text{tr}(T) = 1\}, \quad (4.2.1)$$

where  $V_1 \in \mathbb{C}^{n \times r_1}$  are the eigenvectors corresponding to the algebraically largest eigenvalue  $\lambda_1$  of  $\mathcal{A}^*y$  and  $r_1$  is the multiplicity of  $\lambda_1$  [PB<sup>+</sup>14, Section 6.7]. Note that evaluation of the objective function  $\lambda_1(\mathcal{A}^*y)$  likewise returns the (sub)gradient of this function as a product of the eigenvector(s) corresponding to  $\lambda_1$ . If the eigenvalue  $\lambda_1$  has separation from the second eigenvalue  $\lambda_2$ , then the function  $\lambda_1(\mathcal{A}^*y)$  is differentiable at  $y$  and we have the descent vector

$$g = \nabla\lambda_1(\mathcal{A}^*y) = \mathcal{A}(v_1v_1^*). \quad (4.2.2)$$

However, if the multiplicity of  $\lambda_1$  is greater than one, then  $\lambda_1(\mathcal{A}^*y)$  is nondifferentiable and we may select any  $g \in \partial\lambda_1(\mathcal{A}^*y)$ . In practice, we consider the function  $\lambda_1(\mathcal{A}^*y)$  differentiable if

$$\frac{\lambda_1 - \lambda_2}{\lambda_1} \geq \text{tol}_{\text{diff}}, \quad (4.2.3)$$

for some tolerance  $\text{tol}_{\text{diff}}$ .

Next, the descent vector  $g$  is used to identify an initial steplength and perform a line-search. If  $\lambda_1(\mathcal{A}^*y)$  is differentiable, then we take an initial step with Barzilai-Borwein steplength [BB88]

$$\alpha = \frac{\langle dy, dy \rangle}{\langle dy, dg \rangle}, \quad (4.2.4)$$

where  $dy = y_k - y_{k-1}$  and  $dg = \nabla \lambda_1(\mathcal{A}^* y_k) - \nabla \lambda_1(\mathcal{A}^* y_{k-1})$ . We then perform a linesearch with Wolfe conditions (see [NW06, Section 3.1]) on the problem

$$\min_{\alpha} \lambda_1(\mathcal{A}^* y(\alpha)), \quad y(\alpha) = \Pi_{\mathcal{C}}(y - \alpha g). \quad (4.2.5)$$

This method converges R-linearly for strongly convex functions and is found to outperform standard gradient descent significantly on differentiable functions [ZH04]. However, the linesearch has the added cost of additional objective evaluations  $\lambda_1(\mathcal{A}^* y)$  if the initial value for  $\alpha$  does not satisfy the Wolfe conditions.

If instead (4.2.3) fails and  $\lambda_1(\mathcal{A}^* y)$  is nondifferentiable, then we revert to a decreasing steplength sequence  $\{\alpha_k\}$ . For convex models like the PLGD model, it is known that any sequence of steplengths satisfying the conditions

$$\lim_{k \rightarrow \infty} \alpha_k = 0 \quad \text{and} \quad \sum_{k=0}^{\infty} \alpha_k = \infty \quad (4.2.6)$$

will generate a sequence  $y_k$  converging to an optimal solution (see, e.g., [Ber16, Proposition 1.2.3 and Section 3.3.1]). A typical choice for this steplength sequence is  $\alpha_k = \mathcal{O}(1/k)$ .

The linesearch subproblem (4.2.5) and the resulting update  $\hat{y} = y - \alpha g$  require projection onto the feasible region  $\mathcal{C} = \{y \mid \langle b, y \rangle - \epsilon \|y\| \geq 1\}$ . If  $\epsilon = 0$ , this projection has the closed form expression

$$\Pi_{\mathcal{C}}(\hat{y}) = \begin{cases} \hat{y} + \frac{1 - \langle b, \hat{y} \rangle}{\|b\|^2} b & \text{if } \langle b, \hat{y} \rangle < 1 \\ \hat{y} & \text{else.} \end{cases} \quad (4.2.7)$$

If  $\epsilon > 0$ , then the projection is the solution to the problem

$$\min_{y \in \mathbb{R}^m} \frac{1}{2} \|y - \hat{y}\|^2 \quad \text{subject to} \quad \langle b, y \rangle - \epsilon \|y\| \geq 1. \quad (4.2.8)$$

The KKT conditions of this problem can be simplified into a one-dimensional degree-4 polynomial whose largest real root is used to express  $\Pi_{\mathcal{C}}(\hat{y})$  again as a linear combination of  $\hat{y}$  and  $b$  (details omitted, see [BV04, Chapter 5]).

Given the updated dual variable  $y$ , we must recover a primal variable (signal)  $x$  to test for primal feasibility ( $\|\mathcal{A}(xx^*) - b\| \leq \epsilon$ ). Corollary 3.4.3 indicates that the general primal recovery problem is

$$\min_{S \succeq 0} \|\mathcal{A}(V_1 S V_1^*) - b_{\epsilon}\|^2, \quad b_{\epsilon} := b - \epsilon \frac{y}{\|y\|}, \quad (4.2.9)$$

where  $V_1 \in \mathbb{C}^{r_1 \times n}$  are the eigenvectors corresponding to  $\lambda_1$ ,  $r_1$  is the multiplicity of  $\lambda_1$ ,  $S \in \mathbb{C}^{r_1 \times r_1}$  is positive semidefinite, and  $b_{\epsilon}$  is the noisy observation  $b$  with a removal of the current approximation  $\epsilon y / \|y\|$  to the noise term  $\eta$  as shown in (3.4.9). If  $\lambda_1(\mathcal{A}^* y)$  is differentiable, then  $\lambda_1$  is unique and (4.2.9) simplifies to the closed-form expression

$$\hat{x} = [\langle \mathcal{A}(v_1 v_1^*), b_{\epsilon} \rangle]_+ / \|\mathcal{A}(v_1 v_1^*)\|^2. \quad (4.2.10)$$



The previous steps constitute a complete projected gradient descent method for optimizing the PLGD model. We now discuss a pair of refinement steps which exploit the PLP-PLGD optimality conditions (Corollaries 3.4.1, 3.4.3) to accelerate the convergence of this algorithm. A primal refinement step aims to recover a better approximate signal  $x$  than the primal recovery solution (4.2.10). If the phase retrieval model has no noise, then a dual refinement step further accelerates the convergence rate of the descent method.

The primal refinement step makes further use of the fact that  $\epsilon y/\|y\|$  approximates the noise term  $\eta$ , as indicated by Corollary 3.4.3. Setting  $b_\epsilon = b - \epsilon y/\|y\|$ , we then solve the unconstrained nonconvex problem

$$\min_{x \in \mathbb{C}^n} h(x) := \frac{1}{4} \|\mathcal{A}(xx^*) - b_\epsilon\|^2, \quad (4.2.11)$$

which is initialized with the primal recovery solution  $\hat{x}$  from (4.2.10). This problem can be interpreted as a partially denoised version of the wflow least-squares problem (2.3.14), where  $b$  has been replaced by  $b_\epsilon$ .

For noiseless phase retrieval problems, an additional dual refinement method promotes convergence of the dual iterate  $y$ . In essence, this step uses optimality conditions and the refined signal  $x$  from (4.2.11) to identify the corresponding dual variable  $y$ . If the optimal matrix  $X_\star = \bar{x}\bar{x}^*$  is rank-one, then Corollary 3.4.3 (a) indicates that the algebraically largest eigenpair  $(\lambda_1^*, v_1^*)$  of  $\mathcal{A}^*y_\star$  will be unique. Thus  $v_1^*$  will be a rescaling of the optimal signal  $\bar{x}$ , and Corollary 3.4.1 (d) (strong duality) gives this rescaling

$$\|X_\star\|_1 = \|\bar{x}\|_2^2 = 1/\lambda_1^* \implies \bar{x} = v_1^*/\sqrt{\lambda_1^*}.$$

As a result, the dual refinement step attempts to find the update  $y \in \mathcal{C}$  which satisfies  $[\mathcal{A}^*y]x = \lambda_1 x$  by solving the constrained linear-least-squares problem

$$\min_{y \in \mathbb{R}^m} \frac{1}{2} \|[\mathcal{A}^*y]x - \lambda_1 x\|^2 \quad \text{subject to} \quad \langle b, y \rangle - \epsilon \|y\| \geq 1, \quad (4.2.12)$$

where  $\lambda_1 = \lambda_1(\mathcal{A}^*y)$  is the most recent PLGD objective evaluation and  $x$  is the solution to (4.2.11). If the refined iterate  $\hat{y}$  improves the dual objective, i.e.,  $\lambda_1(\mathcal{A}^*\hat{y}) < \lambda_1(\mathcal{A}^*y)$ , then  $\hat{y}$  replaces  $y$ . This spacer iterate, as described in [Ber16, Proposition 1.2.5], is guaranteed not to interfere with the convergence behavior of the underlying method.

The above sequence of steps and methods leads to Algorithm 3, a projected gradient descent method for optimizing the PLGD model (3.2.1):

---

**Algorithm 3** Projected gradient descent method with refinement

---

**Input:** Sensing operator  $\mathcal{A}$  and adjoint  $\mathcal{A}^*$ , measurement vector  $b$ , initial dual variable  $y$ , estimate of total noise level  $\epsilon$ , convergence tolerances.

**Output:** Approximate solution signal  $x$ .

```
1: while not converged do
2:   Compute algebraically largest eigenpairs:  $(\lambda_1, v_1)$  and  $(\lambda_2, v_2)$  from  $\mathcal{A}^*y$ .
3:   Compute (sub)gradient:  $g = \mathcal{A}(v_1 v_1^*)$  based on (4.2.1).
4:   Determine differentiability of  $\lambda_1(\mathcal{A}^*y)$  based on (4.2.3).
5:   if  $\lambda_1(\mathcal{A}^*y)$  is differentiable then
6:     Linesearch: Perform linesearch (4.2.5) with initial step  $\alpha$  (4.2.4) to obtain  $y_+$ .
7:   else
8:     Projected subgradient step: Set  $y_+ = \Pi_{\mathcal{C}}(y - \alpha g)$  with  $\alpha$  from (4.2.6).
9:   end if
10:  Primal recovery: Compute  $\hat{x}$  based on (4.2.10).
11:  Primal refinement: Find  $x_+$  as the solution to (4.2.11) initialized with  $\hat{x}$  and  $y_+$ .
12:  if  $\epsilon = 0$  then
13:    Dual refinement: Find  $\hat{y}$  based on (4.2.12).
14:    if  $\lambda_1(\mathcal{A}^*\hat{y}) < \lambda_1$ , set  $y_+ = \hat{y}$ 
15:  end if
16:  Update:  $x = x_+$ ,  $y = y_+$ .
17: end while
18: Return:  $x$ .
```

---

We close this section by describing the implementation details of Algorithm 3. This algorithm was originally established in [FM16] and implemented in the accompanying software package `low-rank-opt`<sup>1</sup> (with the identical clone<sup>2</sup>). All tests, experiments, and new methods in this dissertation are available for reproduction in a branched package<sup>3</sup>. This branch includes a few minor changes to the original algorithm which are noncritical to the behavior of this algorithm (see `README.md` in the main directory for details).

In the case of noisy phase retrieval, the input  $\epsilon$  is a measure of the total noise

$$\epsilon = \|\eta\| = \|b - \bar{b}\|$$

of the model (2.1.1). Convergence of Algorithm 3 (step 1) is based on strong duality (Corollary 3.4.1, d) along with primal and dual feasibility. Since dual feasibility is maintained throughout the entire algorithm, only primal feasibility and strong duality must be measured. Thus the algorithm terminates when both of the following conditions are met:

$$\|\mathcal{A}(xx^*) - b\| \leq \epsilon + \text{tol}_{\text{feas}}(1 + \|b\|), \quad (4.2.13)$$

---

<sup>1</sup><https://www.cs.ubc.ca/~mpf/pubs/low-rank-spectral-optimization-via-gauge-duality/>

<sup>2</sup><https://github.com/Will-Wright/low-rank-opt-orig>

<sup>3</sup><https://github.com/Will-Wright/low-rank-opt-rapid-eig>

$$\|xx^*\|_1 \cdot \lambda_1(\mathcal{A}^*y) \leq 1 + \text{tol}_{\text{gap}}. \quad (4.2.14)$$

The convergence tolerances are set to  $\text{tol}_{\text{feas}} = \text{tol}_{\text{gap}} = 2 \times 10^{-4}$  for noisy phase retrieval, and  $\text{tol}_{\text{feas}} = \text{tol}_{\text{gap}} = 1 \times 10^{-5}$  for noiseless.

If the dual variable  $y$  is not provided, then it is initialized as  $y = \Pi_{\mathcal{C}}(b)$ . On the 0-th iteration of this algorithm, steps 5-9 are skipped. This strategy has the result of avoiding an additional eigenvalue computation (step 6) during initialization. Additionally, the algorithm has the default setting of skipping dual refinement in the presence of noise (as dual refinement inhibits convergence for noisy phase retrieval models, see Section 5.3).

The (sub)gradient computation (Algorithm 3, line 3) as described in (4.2.1) is always set to  $g = \mathcal{A}(v_1 v_1^*)$  for the eigenvector  $v_1$  regardless of the multiplicity of  $\lambda_1$ . Since noisy phase retrieval problems often deal with nondifferentiable objectives (see Section 5.3), the algorithm includes a test for differentiability (4.2.3) with the default tolerance  $\text{tol}_{\text{diff}} = 10^{-5}$ . If this condition fails during some iteration  $k$ , the algorithm performs a projected subgradient step (step 8) with a specific steplength  $\alpha_k$ . Testing indicates that projecting the Barzilai-Borwein steplength (4.2.4) onto an interval with monotonically decreasing bounds maintains progress of the algorithm across test cases. Thus the steplength is set to

$$\alpha_k = \min \left\{ u_k, \max \left\{ l_k, \frac{\langle dy, dy \rangle}{\langle dy, dg \rangle} \right\} \right\}, \quad (4.2.15)$$

where  $u_k = 200/k$  and  $l_k = 1/k$ .

Computationally, most of the steps in Algorithm 3 are very simple. The three key exceptions are the eigenvalue computation (steps 2, 6, and 14), the primal refinement (step 11), and the dual refinement (step 13). In the original implementation, the eigenvalue computation is performed using the MATLAB function `eigs` (see Section 6.3 for details about this method), where only the first two algebraically largest eigenvalues  $\lambda_1, \lambda_2$  are requested. The primal refinement (step 11) is computed using `minFunc`, a quasi-Newton solver for unstrained optimization, with the descent direction determined using the limited-memory (l-)BFGS method for Hessian approximation [Sch05]. Dual refinement (step 13) is computed with `minConf`, another quasi-Newton solver which is optimized for problems like (4.2.12) with expensive objective functions and simple constraints [Sch08], [SBFM09].

In each of these three subroutines, the primary computational cost comes from  $\mathcal{A}$ -products (2.3.7), where each  $\mathcal{A}(xx^*)$  product requires  $L$  DFTs and each  $[\mathcal{A}^*y]x$  product requires  $2L$  DFTs. Thus we measure computational costs in terms of number of DFTs, following the convention of [CLS15] and [FM16].

Theoretically, the steplength strategy in Algorithm 3 is guaranteed to converge for all feasible problems (i.e., find an optimal pair  $(X_\star, y_\star)$  which satisfies the (PLGD) optimality conditions of Corollary 3.4.1) [ZH04], [Ber16, Proposition 1.2.3 and Section 3.3.1]. Yet the rate of convergence differs greatly for noiseless and noisy phase retrieval problems. This algorithm is particularly efficient for noiseless problems, as we discuss briefly in Section 4.3. However, noisy problems pose a unique set of challenges which are intrinsic to the PLGD model (3.2.1) which are addressed in the remainder of this dissertation.

### 4.3 Efficient recovery for noiseless phase retrieval

This section examines the behavior of Algorithm 3 for noiseless phase retrieval problems. The efficiency of this algorithm for noiseless problems is well-established in [FM16, Sections 5.1.1, 5.1.3]. The purpose of this section is to establish the role of the primal (4.2.11) and dual refinement (4.2.12) steps for noiseless phase retrieval so we may examine the utility of these steps for noisy phase retrieval in Chapter 5.

We begin by discussing the general behavior of each refinement step and the relationship between the two. For noiseless problems ( $\epsilon = 0$ ), the primal refinement step (4.2.11) is initialized with  $b_\epsilon = b$  and is therefore independent of the dual variable  $y$  other than the fact that (4.2.11) is initialized with  $\hat{x}$  from (4.2.10). In this case, (4.2.11) is equivalent to the wflow least-squares problem (2.3.14) discussed in Section 2.3.3. As proved in [SQW16] and restated at the end of Section 2.3.3, if this problem has a sufficient oversampling  $L$ , then the objective function of (4.2.11) will have no spurious local minima, only one global minima (up to a phase constant), and negative directional curvature at all saddle points. As a result, for noiseless problems the primal refinement step effectively solves the phase retrieval problem on the first iteration of Algorithm 3 (i.e.,  $x$  typically has a relative error  $\|\bar{x}\bar{x}^* - xx^*\|_F / \|\bar{x}\|_2^2$  on the order of  $10^{-7}$  or smaller). Thus the primal refinement step serves as an effective standalone method for promoting the quick convergence of the primal signal  $x$ .

In contrast to primal refinement, the dual refinement problem (4.2.12) is directly dependent on the accuracy of the primal signal  $x$ . If the dual refinement is initialized with an inaccurate signal, then the returned dual variable  $\hat{y}$  may be inferior to the initial update  $y_+$ . Yet  $\hat{y}$  may also satisfy the condition for replacing the previous variable (step 14), preventing Algorithm 3 from converging.

Table 4.1 compares the behavior of these refinement steps for a random noiseless phase retrieval problem.

	Primal & dual ref.	Primal ref. only	Dual ref. only	No ref.
Iterations	1	62	1,000	1,000
DFTs	21,800	349,420	273,817,590	14,686,690
Relative error	$7.469 \times 10^{-9}$	$8.097 \times 10^{-9}$	$1.166 \times 10^2$	$1.176 \times 10^2$

Table 4.1: Algorithm 3 with and without primal (4.2.11) and dual refinement (4.2.12). This experiment involves a random gaussian signal of size  $n = 128$  with  $L = 10$  observations and no noise. Both termination conditions (4.2.13) and (4.2.14) are required. Signal relative error is measured as  $\|\bar{x}\bar{x}^* - x_0x_0^*\|_F / \|\bar{x}\|_2^2$ .

Table 4.1 demonstrates that primal refinement is necessary for Algorithm 3 to converge, and dual refinement is essential to the efficiency of this algorithm. With only primal refinement, this algorithm exhibits a convergence rate typical of a steepest descent method.

In this case, primal feasibility (4.2.13) is achieved in the first few iterates. Yet the duality gap condition (4.2.14) is not satisfied until the dual variable  $y$  is sufficiently close to  $y_*$ . With only the dual refinement, this algorithm fails to converge as the dual problem (4.2.12) is initialized with an inaccurate signal  $x$ . And when Algorithm 3 is run without either refinement step, it again fails to converge.

However, when both the primal and dual refinement subroutines are used, Algorithm 3 rarely takes more than a couple iterations. Since steps 5-9 are skipped during the 0-th iterate, Algorithm 3 serves as an improvement to the wflow algorithm. Both algorithms set the initial signal as the eigenvector corresponding to the algebraically largest eigenvalue of  $\mathcal{A}^*b$  (where Algorithm 3 first rescales  $b$  by projecting it onto  $\mathcal{C}$ ). Yet Algorithm 3 uses better search directions generated by the l-BFGS method rather than the Wirtinger derivative. Next, the dual refinement problem (4.2.12) is initialized with a sufficient pair of terms  $(\lambda_1, x)$  to recover a nearly optimal dual iterate  $y$ . Thus the primal feasibility (4.2.13) and duality gap (4.2.14) conditions are typically met within a few iterations.

Algorithm 3 with primal and dual refinement is very effective for noiseless phase retrieval, acting as a robust, efficient competitor to the wflow algorithm with the added benefit of greater signal accuracy [FM16, Section 5.1.1, 5.1.3]. Note that the optimality of  $y$  is unnecessary if we simply want to minimize the noiseless problem  $\|\mathcal{A}(xx^*) - b\|$ . To this end, Algorithm 3 as implemented includes a setting specifically for noiseless problems, where the strong duality condition (4.2.14) is dropped and only primal feasibility (4.2.13) is required. This primal feasibility version of Algorithm 3 typically converges in 0-1 iterations with about the same computational cost as the wflow algorithm.

This dissertation addresses the challenges noisy phase retrieval poses for Algorithm 3 with the following contributions.

Chapter 5 addresses the convergence challenges Algorithm 3 experiences for noisy problems. We begin by demonstrating that this algorithm stagnates prior to convergence for noisy problems and determining the cause of this stagnation. Although Algorithm 3 does not generally converge for noisy phase retrieval problems, we demonstrate that this algorithm is more robust and accurate than the wflow method 2 for noisy problems. Section 5.4 establishes new termination conditions which indicate stagnation, so we may treat Algorithm 3 as a black-box and focus our attention on handling the EMEP.

Chapters 6 and 7 examine the computational challenges Algorithm 3 experiences for noisy problems. We identify the EMEP as the dominant computational cost and develop new methods for handling this problem. Chapter 6 provides a detailed explanation of three common methods for large-scale eigenvalue problems. The unique structure and properties of each method inform our strategy for handling the EMEP. Chapter 7 then examines the evolving spectral structure of the matrices  $\mathcal{A}^*y_k$  across sequences of iterates  $k$  to develop an optimal method for handling the EMEP. Chapter 7 then presents a revised version of Algorithm 3 with an optimized method for handling the EMEP and new termination conditions for identifying stagnation.

## Chapter 5

# New termination conditions for the first-order PLGD algorithm

### 5.1 Introduction

In the previous chapter we developed Algorithm 3 to optimize the PLGD model (3.2.1) and saw that this method is efficient and accurate for noiseless phase retrieval. We now examine the tendency of Algorithm 3 to fail to converge for noisy phase retrieval and establish new termination conditions for this algorithm.

We begin in Section 5.2 by describing how we construct experimental noisy PLGD models and discussing potential residuals for measuring the progress of Algorithm 3. Section 5.3 then examines the tendency of Algorithm 3 not to converge for noisy phase retrieval problems. We see that signals observed with nontrivial noise tend to have optimal PLGD matrices  $X_*$  with rank greater than one, preventing first-order methods like Algorithm 3 from converging. To handle this issue, Section 5.4 establishes new termination conditions based on heuristic evidence indicating Algorithm 3 has stopped making signal recovery progress. These new conditions allow us to treat Algorithm 3 as a black-box solver in the remainder of this dissertation so we may develop optimal methods for handling the EMEP, the computational bottleneck of this algorithm.

### 5.2 Experimental models and potential residuals

This section describes two methods for creating experimental noisy phase retrieval problems and presents a set of potential residuals to measure the progress of Algorithm 3. The *phase retrieval problem with Gaussian noise* imitates the typical phase retrieval scenario and is used throughout this dissertation as the default method for creating noisy phase retrieval problems. The *phase retrieval problem with synthetic noise* is constructed around the PLGD primal recovery conditions (Corollary 3.4.3) and is used exclusively in Section 5.3

to examine the convergence behavior of Algorithm 3. The set of potential residuals is a combination of those discussed in Chapter 4 and new residuals based on the variables available in Algorithm 3.

We begin by describing experimental noisy phase retrieval problems which have *Gaussian noise*. Recall that the noisy phase retrieval problem (2.1.1) involves an observation  $b = \bar{b} + \eta$ , where the true observation  $\bar{b}$  is contaminated by some nontrivial noise  $\eta$ . To mimic the typical experimental noisy phase retrieval scenario, we begin with an unknown signal  $\bar{x}$ . A sensing operator  $\mathcal{A}$  (2.3.7) is then chosen with diagonal mask matrices  $C_j$  whose diagonal elements have complex standard Gaussian distribution (1.2.17) (see Section 2.2 for an explanation of masks). The sensing operator is then used to create the true observation  $\bar{b} = \mathcal{A}(\bar{x}\bar{x}^*)$ , and a noise term  $\eta$  is chosen with real standard Gaussian distribution (1.2.16). Finally, the noisy observation is set as  $b = \bar{b} + \eta$ . Altogether, given a signal  $\bar{x}$ , a sensing operator  $\mathcal{A}$  (2.3.7), and noise ratio  $\epsilon_{\text{rel}}$ , we create the phase retrieval problem with Gaussian noise experimentally with the steps

$$\begin{aligned} 1) \quad & \bar{b} = \mathcal{A}(\bar{x}\bar{x}^*), \\ 2) \quad & \eta \sim \mathcal{N}(0, 1), \\ 3) \quad & b = \bar{b} + \eta, \end{aligned} \tag{5.2.1}$$

where  $\eta$  is rescaled in the third step to satisfy the noise ratio  $\epsilon_{\text{rel}} = \|\eta\|/\|b\|$ . Thus we define the *phase retrieval problem with Gaussian noise* as

$$\begin{aligned} \text{find} \quad & x \\ \text{s.t.} \quad & \|\mathcal{A}(xx^*) - b\|_2 \leq \epsilon, \end{aligned} \tag{5.2.2}$$

where  $b = \bar{b} + \eta$  is constructed experimentally using the steps (5.2.1). Likewise, the *PLGD model with Gaussian noise* refers to the the PLGD model (3.2.1)

$$\begin{aligned} \min_y \quad & \lambda_1(\mathcal{A}^*y) \\ \text{(PLGD) s.t.} \quad & \langle b, y \rangle - \epsilon\|y\| \geq 1, \end{aligned} \tag{5.2.3}$$

where again  $b = \bar{b} + \eta$  is constructed experimentally using the steps (5.2.1).

To demonstrate how the differentiability of the dual objective  $\lambda_1(\mathcal{A}^*y)$  impacts the behavior of Algorithm 3, Section 5.3 also considers experimental noisy phase retrieval problems which we say have *synthetic noise*. The purpose of this synthetic noise is to create a PLGD model where the dual objective  $\lambda_1(\mathcal{A}^*y)$  is differentiable at  $y_\star$ . To achieve this goal, the phase retrieval problem with synthetic noise is constructed to satisfy the PLGD primal recovery conditions of Corollary 3.4.3. First we choose a random variable with standard Gaussian distribution (1.2.16) to be the optimal dual variable  $y_\star$  and use this vector to construct the rest of the noisy phase retrieval problem. The true signal  $\bar{x}$  is set as the eigenvector corresponding to the algebraically largest eigenvalue of  $\mathcal{A}^*y_\star$  per Corollary 3.4.3 (a). This signal is then used to create a true observation  $\bar{b}$  which is noised

with a rescaling of  $y_\star$  per Corollary 3.4.3 (b) and (3.4.9). Altogether, given a noise ratio  $\epsilon_{\text{rel}}$  we have the following steps for constructing the phase retrieval problem with synthetic noise

$$\begin{aligned}
& 1) \quad y_\star \sim \mathcal{N}(0, 1), \\
& 2) \quad \bar{x} \text{ is the eigenvector corresponding to} \\
& \quad \text{the algebraically largest eigenvalue of } \mathcal{A}^* y_\star, \\
& 3) \quad \bar{b} = \mathcal{A}(\bar{x} \bar{x}^*), \\
& 4) \quad b = \bar{b} + \epsilon \frac{y_\star}{\|y_\star\|},
\end{aligned} \tag{5.2.4}$$

where  $\epsilon = \epsilon_{\text{rel}} \|b\|$ . Steps one and two in (5.2.4) guarantee the PLP-PLGD optimality conditions (Corollary 3.4.1) will hold for a known pair  $(X_\star = \bar{x} \bar{x}^*, y_\star)$  with rank-one optimal matrix  $X_\star$ . Additionally, step four in (5.2.4) satisfies the primal recovery equation (3.4.9), guaranteeing the primal refinement strategy (4.2.11) used in Algorithm 3 will recovery the true signal  $\bar{x}$  when initialized with  $y_\star$ . We define the *phase retrieval problem with synthetic noise* as

$$\begin{aligned}
& \text{find } x \\
& \text{s.t. } \quad \|\mathcal{A}(xx^*) - b\|_2 \leq \epsilon,
\end{aligned} \tag{5.2.5}$$

where  $b = \bar{b} + \epsilon \frac{y_\star}{\|y_\star\|}$  is constructed experimentally using the steps (5.2.4). Likewise, the *PLGD model with synthetic noise* refers to the the PLGD model (3.2.1)

$$\begin{aligned}
& \min_y \quad \lambda_1(\mathcal{A}^* y) \\
& \text{(PLGD) s.t. } \quad \langle b, y \rangle - \epsilon \|y\| \geq 1,
\end{aligned} \tag{5.2.6}$$

where again  $b = \bar{b} + \epsilon \frac{y_\star}{\|y_\star\|}$  is constructed experimentally using the steps (5.2.4).

Next we establish an exhaustive list of residuals which are available for measuring the progress of Algorithm 3. Here the PLGD dual matrix is defined  $A := \mathcal{A}^*(y)$  and previous iterates are denoted with a hat (e.g.,  $\hat{y} = y_{k-1}$ ). Note that this set of residuals does not contain a residual for dual feasibility because Algorithm 3 maintains dual feasibility (each gradient descent step projects the dual iterate  $y$  onto the feasible set).



signal relative error	$\ xx^* - \bar{x}\bar{x}^*\ _F / \ \bar{x}\bar{x}^*\ _F$	(5.2.7a)
dual relative error	$\ y - y_\star\  / \ y_\star\ $	(5.2.7b)
primal true relative error	$\ \mathcal{A}(xx^*) - \bar{b}\  / \ \bar{b}\ $	(5.2.7c)
primal relative error	$\rho := \ \mathcal{A}(xx^*) - b\  / \ b\ $	(5.2.7d)
primal difference	$ \rho - \hat{\rho}  /  \rho $	(5.2.7e)
duality gap	$\gamma := \ xx^*\ _1 \cdot \lambda_1 - 1$	(5.2.7f)
duality gap difference	$ \gamma - \hat{\gamma}  /  \gamma $	(5.2.7g)
dual objective difference	$ \lambda_1 - \hat{\lambda}_1  /  \lambda_1 $	(5.2.7h)
dual variable difference	$\ y - \hat{y}\  / \ y\ $	(5.2.7i)
dual matrix difference	$\ A - \hat{A}\ _F / \ A\ _F$	(5.2.7j)

The residuals in (5.2.7) can be separated into three groups. The first group contains ideal error measurements, including the signal relative error (5.2.7a), dual relative error (5.2.7b), and primal true relative error (5.2.7c). The next group of residuals are those used as termination conditions for Algorithm 3 in Chapter 4: the primal relative error (5.2.7d) and the duality gap (5.2.7f), used in conditions (4.2.13) and (4.2.14), respectively. The final group of residuals are the five difference residuals (5.2.7e, g-j). We present these difference residuals as a new set of measurements for determining when Algorithm 3 is no longer making signal recovery progress and termination should be declared. As we will see in Section 5.4, the primal difference (5.2.7e) and dual variable difference (5.2.7i) are the two difference residuals best suited for determining stagnation of Algorithm 3, and thus used to establish new termination conditions in that section.

### 5.3 Stagnation of the first-order PLGD algorithm

In this section we demonstrate the tendency of Algorithm 3 not to converge when solving PLGD models with Gaussian noise (5.2.3). We see that PLGD models with Gaussian noise (5.2.3) typically have an optimal matrix  $X_\star$  with rank greater than one. In this circumstance, Corollary 3.4.1 (e) indicates that the algebraically largest eigenvalue of the optimal dual matrix  $\mathcal{A}^*y_\star$  will also be greater than one, causing the dual objective  $\lambda_1(\mathcal{A}^*y)$  to be nondifferentiable in a neighborhood of  $y_\star$ . Additionally, since  $X_\star \neq \bar{x}\bar{x}^*$ , the primal refinement strategy (4.2.11) used in Algorithm 3 is not guaranteed to recover the true signal  $\bar{x}$ .

In order to describe the behavior of Algorithm 3 for noisy phase retrieval problems, we use the following terminology to make the distinction between Algorithm 3 converging to an optimal solution and converging to a nonoptimal variable. Recall that Algorithm 3 *converges* at a given iterate if the conditions (4.2.13) and (4.2.14) are satisfied. We

say that Algorithm 3 *fails to converge* if conditions (4.2.13) or (4.2.14) are not satisfied for any iterate within a maximum number of iterations (typically 1,000). In contrast with converging, we say Algorithm 3 *stagnates* at a given iterate if the progress from the previous iterate is trivial. Specifically, Algorithm 3 stagnates at a given iterate if a chosen subset of difference residuals (5.2.7e, g-j) are below a required set of tolerances.

The authors of [FM16] use PLGD models with synthetic noise (5.2.6) to demonstrate that Algorithm 3 is able to converge and accurately recover the true signal  $\bar{x}$  when the optimal matrix in the PLGD model is rank-one (i.e.,  $X_\star = \bar{x}\bar{x}^*$ ). Table 5.1 depicts the results of Algorithm 3 applied to the PLGD model with synthetic noise (5.2.6), corroborating the results in [FM16, Section 5.1.2].

	$L = 5$		$L = 10$		$L = 15$	
	Iterations	xErr	Iterations	xErr	Iterations	xErr
$\epsilon_{\text{rel}} = 0.05$	154	$8.581 \times 10^{-3}$	112	$8.442 \times 10^{-3}$	355	$1.618 \times 10^{-3}$
$\epsilon_{\text{rel}} = 0.15$	207	$2.772 \times 10^{-3}$	255	$3.500 \times 10^{-4}$	82	$1.196 \times 10^{-2}$
$\epsilon_{\text{rel}} = 0.30$	186	$1.636 \times 10^{-3}$	104	$2.047 \times 10^{-3}$	111	$4.606 \times 10^{-3}$

Table 5.1: Number of iterations and signal relative error (5.2.7a) (xErr) for Algorithm 3 applied to the PLGD model with synthetic noise (5.2.6). Signal size is  $n = 128$ , with various noise ratios  $\epsilon_{\text{rel}}$  and oversampling values  $L$ . The convergence conditions (4.2.13) and (4.2.14) are set to tolerances  $\text{tol}_{\text{feas}} = \text{tol}_{\text{gap}} = 2 \times 10^{-4}$ . Note that each signal relative error in Table 5.1 was 1-3 orders of magnitude smaller than the relative error of the signal returned by wflow (Algorithm 2) for the same problem.

Table 5.1 demonstrates that Algorithm 3 tends to converge within a few hundred iterations when solving PLGD models with synthetic noise (5.2.6). However, the convergence behavior depicted in Table 5.1 does not occur when Algorithm 3 solves PLGD models with Gaussian noise (5.2.3). In Table 5.2, the experiments from Table 5.1 are replaced with phase retrieval problems with Gaussian noise (5.2.2). Note that the iteration count is replaced by the final duality gap value, since Algorithm 3 failed to converge after 1,000 iterations for all cases.

	$L = 5$		$L = 10$		$L = 15$	
	duGap	xErr	duGap	xErr	duGap	xErr
$\epsilon_{\text{rel}} = 0.05$	323.03	$9.072 \times 10^{-2}$	101.56	$4.053 \times 10^{-2}$	93.08	$3.112 \times 10^{-2}$
$\epsilon_{\text{rel}} = 0.15$	448.68	$4.200 \times 10^{-1}$	364.07	$1.249 \times 10^{-1}$	374.10	$9.443 \times 10^{-2}$
$\epsilon_{\text{rel}} = 0.30$	492.56	$1.096e \times 10^0$	665.73	$2.973 \times 10^{-1}$	747.68	$1.958 \times 10^{-1}$

Table 5.2: Final duality gap (5.2.7f) (duGap) and signal relative error (5.2.7a) (xErr) for Algorithm 3 applied to PLGD models with Gaussian noise (5.2.3). Algorithm 3 failed to converge after 1,000 iterations for all problems. Signal size is  $n = 128$ , with various noise ratios  $\epsilon_{\text{rel}}$  and oversampling values  $L$ . The convergence conditions (4.2.13) and (4.2.14) are set to tolerances  $\text{tol}_{\text{feas}} = \text{tol}_{\text{gap}} = 2 \times 10^{-4}$ .

As we see in Table 5.2, the duality gap value (5.2.7f) remains several orders of magnitude above the duality gap convergence tolerance (4.2.14). As a result, Algorithm 3 fails to converge for PLGD models with Gaussian noise (5.2.3).

The next experiment demonstrates that Algorithm 3 applied to PLGD models with Gaussian noise (5.2.3) tends to achieve primal feasibility (4.2.13) during the early iterations and then stagnates within a few hundred iterations. Since the model (5.2.3) is constructed without knowledge of an optimal PLP-PLGD pair  $(X_*, y_*)$ , we use the interior-point solver SDPT3 [TTT99] to obtain the pair  $(X_*, y_*)$  within square-root machine-precision. In order to use this interior-point solver, the models in Figure 5.1 are very small.

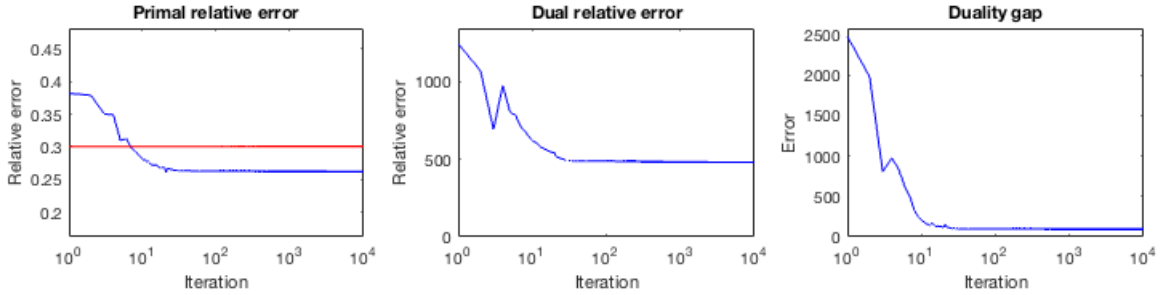


Figure 5.1: Primal relative error (5.2.7d), dual relative error (5.2.7b), and duality gap (5.2.7f) for 10,000 iterates of Algorithm 3 applied to a natural noise model with  $n = 16$ ,  $L = 6$  observations, and noise ratio 0.30. The horizontal axis is log-plotted to highlight early progress along with later stagnation. The pair  $(X_*, y_*)$  are computed with SDPT3.

Figure 5.1 demonstrates that Algorithm 3 stagnates when attempting to solve a PLGD model with Gaussian noise (5.2.3). In this problem, primal feasibility (4.2.13) is achieved at iterate 6, and further progress is made over the following early iterations, as indicated by the primal relative error (5.2.7d) plot. However, the strong duality condition (4.2.14) is never satisfied, and the final duality gap at iterate 10,000 is 97.63. Likewise, the dual

variable  $y_k$  does not approach the optimal dual variable  $y_*$ , as indicated by the dual relative error (5.2.7b) plot. In this particular PLGD model,  $X_*$  is rank three, causing the dual objective  $\lambda_1(\mathcal{A}^*y)$  to be nondifferentiable near  $y_*$  per Corollary 3.4.1 (e). Thus Algorithm 3 identifies the dual objective  $\lambda_1(\mathcal{A}^*y)$  as nondifferentiable at iterate 84 reverts to the monotone stepsize sequence (4.2.15).

The next experiment establishes that PLGD models with Gaussian noise (5.2.3) almost always have an optimal matrix  $X_*$  with rank greater than one, and this rank problem causes the stagnation of Algorithm 3. Table 5.3 compares PLGD models with synthetic (5.2.6) and Gaussian noise (5.2.3), depicting the rank of the optimal matrix  $X_*$  and the behavior of Algorithm 3.

		Synthetic noise			Gaussian noise		
		$L = 4$	$L = 6$	$L = 8$	$L = 4$	$L = 6$	$L = 8$
$\epsilon_{\text{rel}} = 0.05$	$\text{rank}(X_*)$	1	1	1	3.41	3.28	3.27
	Algo. 3 itns.	125.09	144.20	182.26	1,000	1,000	1,000
	duGap	1.17 <sub>-4</sub>	1.18 <sub>-4</sub>	1.18 <sub>-4</sub>	9.60	3.88	3.86
$\epsilon_{\text{rel}} = 0.15$	$\text{rank}(X_*)$	1	1	1	2.99	3.00	3.04
	Algo. 3 itns.	62.48	85.32	95.58	1,000	1,000	1,000
	duGap	1.20 <sub>-4</sub>	1.35 <sub>-4</sub>	1.33 <sub>-4</sub>	15.9	14.0	15.8
$\epsilon_{\text{rel}} = 0.30$	$\text{rank}(X_*)$	1	1	1	2.64	2.70	2.89
	Algo. 3 itns.	40.34	50.88	64.28	1,000	1,000	1,000
	duGap	1.17 <sub>-4</sub>	1.23 <sub>-4</sub>	1.27 <sub>-4</sub>	21.2	26.7	31.5

Table 5.3: Algorithm 3 results and optimal matrix rank for PLGD models with synthetic (5.2.6) and Gaussian noise (5.2.3). This table depicts the mean rank of  $X_*$ , number of Algorithm 3 iterations, and final duality gap (5.2.7f) (duGap) for 100 random phase retrieval models with signal size  $n = 16$ , noise ratio  $\epsilon_{\text{rel}}$ , and oversampling  $L$ . In all synthetic noise models (5.2.6), the solution  $X_*$  was rank-one and Algorithm 3 converged. In all Gaussian noise models (5.2.3), the algorithm reached the maximum of 1,000 iterations without attaining the termination condition (4.2.14). Note that pairs  $(X_*, y_*)$  are computed with SDPT3 and numbers  $n_{-k}$  are shorthand for  $n \times 10^{-k}$ .

Table 5.3 demonstrates that PLGD models with Gaussian noise (5.2.3) typically have optimal matrices with rank greater than one, causing Algorithm 3 to stagnate. Algorithm 3 cannot attain an optimal primal matrix  $X_*$  with rank greater than one because the primal recovery (4.2.10) and refinement (4.2.11) steps used in this algorithm only return a rank-one matrix  $X = xx^*$ . Additionally, Corollary 3.4.1, (e) indicates that  $\text{rank}(X_*)$  is a lower bound on the multiplicity of the algebraically largest eigenvalue of  $\mathcal{A}^*y_*$ . Thus the objective function  $\lambda_1(\mathcal{A}^*y)$  will be nondifferentiable in some neighborhood around  $y_*$  and Algorithm 3 will stagnate prior to approaching  $y_*$ . As a result, this algorithm cannot attain a pair  $(X, y)$  that are sufficiently close to the optimal pair  $(X_*, y_*)$  and fails to satisfy the

duality gap condition (4.2.14).

In contrast to the Gaussian models (5.2.3), we see that PLGD models with synthetic noise (5.2.6) consistently have rank-one optimal matrices  $X_\star$ , allowing Algorithm 3 to converge. Each synthetic noise model (5.2.6) in Table 5.3 had an optimal dual matrix  $\mathcal{A}^*y_\star$  with a unique algebraically largest eigenvalue, making the dual objective function  $\lambda_1(\mathcal{A}^*y)$  differentiable at  $y_\star$  and allowing Algorithm 3 to approach the optimal dual variable  $y_\star$ . Once the pair  $(X, y)$  are sufficiently close to the optimal pair, the duality gap condition (4.2.14) is met and convergence is declared. At this point, the primal recovery equation (3.4.9) successfully denoises the noisy observation  $b = \bar{b} + \epsilon y_\star / \|y_\star\|$  because the synthetic noise steps (5.2.4) guarantee an exact relation  $\eta = \epsilon y_\star / \|y_\star\|$  between the noise term  $\eta$  and optimal dual variable  $y_\star$ . Thus Algorithm 3 is able to return a matrix  $X$  which accurately approximates the optimal matrix  $X_\star = \bar{x}\bar{x}^*$ .

Table 5.3 also helps to explain why dual refinement (Algorithm 3, step 13) is not beneficial for PLGD models with Gaussian noise (5.2.3). As we saw in Table 4.1, an inaccurate signal  $x$  can cause the dual refinement problem (4.2.12) to return an unreliable update  $\hat{y}$ . Figure 5.2 depicts the progress made by Algorithm 3 with and without dual refinement for three PLGD models with Gaussian noise (5.2.3).

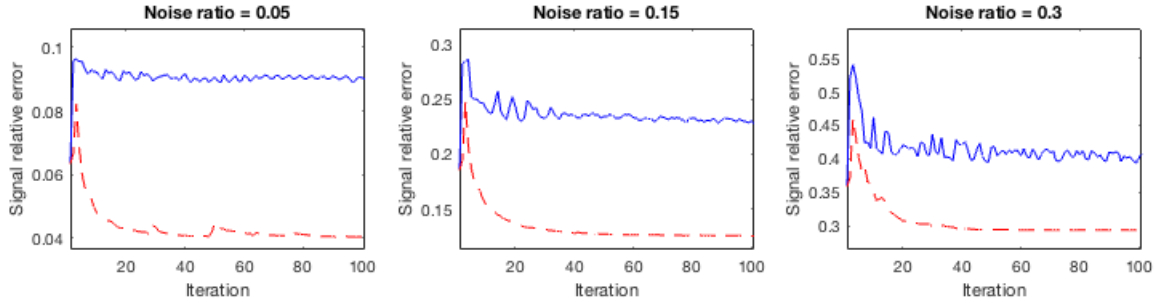


Figure 5.2: A comparison of Algorithm 3 with and without dual refinement (4.2.12) for PLGD models with Gaussian noise (5.2.3) with various noise levels, with random Gaussian signal of size  $n = 128$  and  $L = 10$  observations. The solid blue line indicates dual refinement, and the dashed red line indicates no dual refinement.

Figure 5.2 demonstrates that the dual refinement step (4.2.12) inhibits the progress otherwise made by Algorithm 3. This step is initialized with a signal  $x$  which corresponds to the rank-one matrix iterate  $X = xx^*$ . Since the optimal matrix  $X_\star$  tends to have rank greater than one, the dual refinement problem (4.2.12) will not be properly initialized and may return a poor update  $\hat{y}$ . Thus the dual refinement step in Algorithm 3 is not beneficial for signal recovery in phase retrieval problems with Gaussian noise (5.2.2).

## 5.4 New termination conditions

In Section 5.3, we saw that Algorithm 3 tends to stagnate when solving phase retrieval problems with Gaussian noise (5.2.2), failing to satisfy the duality gap termination condition (4.2.14) originally established in [FM16]. This section establishes new termination conditions for Algorithm 3 for PLGD models with Gaussian noise (5.2.3) and demonstrates the effectiveness of these conditions for identifying stagnation of Algorithm 3. (For a comparison of the unused residuals from (5.2.7), see Appendix II, Chapter B.)

We begin by presenting the new termination conditions for Algorithm 3 for PLGD models with Gaussian noise (5.2.3). To identify stagnation of Algorithm 3 and declare termination, the primal difference (5.2.7e) is set to

$$\frac{|\rho - \hat{\rho}|}{\rho} \leq \text{tol}_{\text{primal}} = 10^{-5}, \quad \rho = \|\mathcal{A}(xx^*) - b\| \quad (5.4.1)$$

and the dual variable difference (5.2.7i) is set to

$$\frac{\|y - \hat{y}\|}{\|y\|} \leq \text{tol}_{\text{dual}} = 10^{-4}, \quad (5.4.2)$$

where a hat indicates the previous iterate (e.g.,  $\hat{y} = y_{k-1}$ ). Termination is declared when (5.4.1) and (5.4.2) are satisfied or after a maximum of 300 iterations are performed. After termination, the signal returned corresponds to the signal among the previous 20 with the smallest duality gap (5.2.7f) value.

To demonstrate the effectiveness of these termination conditions, we consider a set of six PLGD models with Gaussian noise (5.2.3) solved with Algorithm 3. The signal relative error (5.2.7a) serves as a control measurement to identify when Algorithm 3 has stagnated and should terminate. Figure 5.3 depicts the signal relative error (5.2.7a) for each of these models.

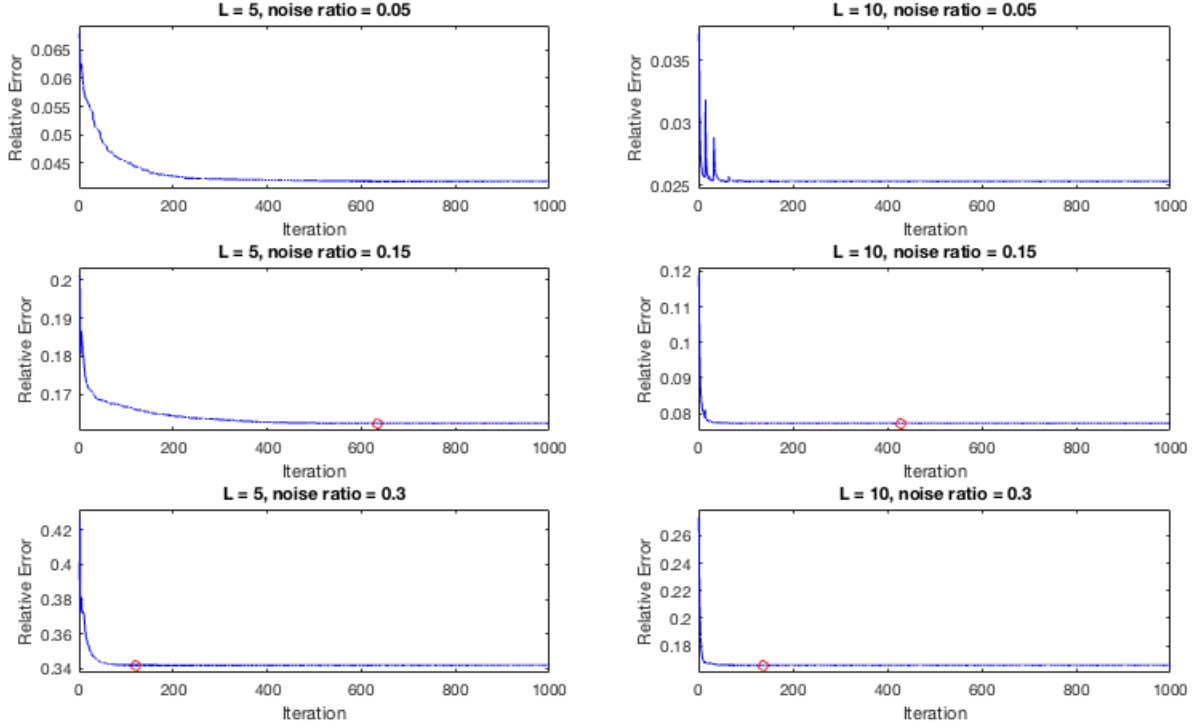


Figure 5.3: Signal relative errors (5.2.7a) for six PLGD models with Gaussian noise (5.2.3) solved with Algorithm 3. The image in Figure 2.2 was resized to  $64 \times 64$  pixels, made grayscale, and six models were generated with oversampling rates of  $L = 5, 10$  and noise ratios  $\epsilon_{\text{rel}} = 0.05, 0.15, 0.30$ . These models were then solved using Algorithm 3 set to terminate after 1,000 iterations. The red circle (where present) indicates when the dual objective was determined nondifferentiable.

Figure 5.3 depicts the early progress and quick stagnation of Algorithm 3 for these PLGD models with Gaussian noise (5.2.3). For each model, virtually no measurable progress was made after iteration 500. Yet the point at which each model stagnates appears to differ, and in one case ( $L = 10$  and  $\epsilon_{\text{rel}} = 0.05$ ) the signal quality appears to vary greatly for neighboring iterates. Thus we examine the behavior in Figure 5.3 to identify the desired interval of iterates in which Algorithm 3 should terminate for each model.

The results in Figure 5.3 suggest that models with a larger oversampling rate and those with a larger noise ratio make progress faster and stagnate earlier. One indicator that the point of stagnation differs for these models is the point at which Algorithm 3 identifies the objective function  $\lambda_1(\mathcal{A}^*y)$  as nondifferentiable (i.e., the condition (4.2.3) fails). The models in Figure 5.3 with the least noise never identified a nondifferentiable objective, whereas the noisiest models identified nondifferentiable objectives very early (at iterate

122 for  $L = 5$  and 137 for  $L = 10$ ). Another indicator that the models of Figure 3 stagnate at different rates is the point at which the primal objective is found to be feasible, as described in Table 5.4.

	$L = 5$	$L = 10$
$\epsilon_{\text{rel}} = 0.05$	33	3
$\epsilon_{\text{rel}} = 0.15$	N/A	3
$\epsilon_{\text{rel}} = 0.30$	22	3

Table 5.4: Iterate at which Algorithm 3 became primal feasible for models from Figure 5.3.

Table 5.4 shows that all three models with oversampling  $L = 10$  found feasible signals after just 3 iterations, yet the models with  $L = 5$  were a bit slower, and in particular the model with  $L = 5$  and  $\epsilon_{\text{rel}} = 0.15$  never identified a feasible signal.

Based on these observations, Table 5.5 proposes iterate intervals in which Algorithm 3 appears to have stagnated for each model in Figure 5.3, providing a guideline for assessing the new termination conditions (5.4.1) and (5.4.2).

	$L = 5$	$L = 10$
$\epsilon_{\text{rel}} = 0.05$	200-400	50-200
$\epsilon_{\text{rel}} = 0.15$	200-400	50-100
$\epsilon_{\text{rel}} = 0.30$	100-200	50-100

Table 5.5: Intervals of iterates at which Algorithm 3 appears to stagnate for models from Figure 5.3.

In order to demonstrate that the primal difference (5.2.7e) and dual variable difference (5.2.7i) are indicators of the stagnation of Algorithm 3, Figure 5.4 depicts the behavior of these residuals for the models in Figure 5.3. Note that the vertical axis indicates specific tolerances and the horizontal axis indicates the first iterate at which Algorithm 3 would satisfy this tolerance.



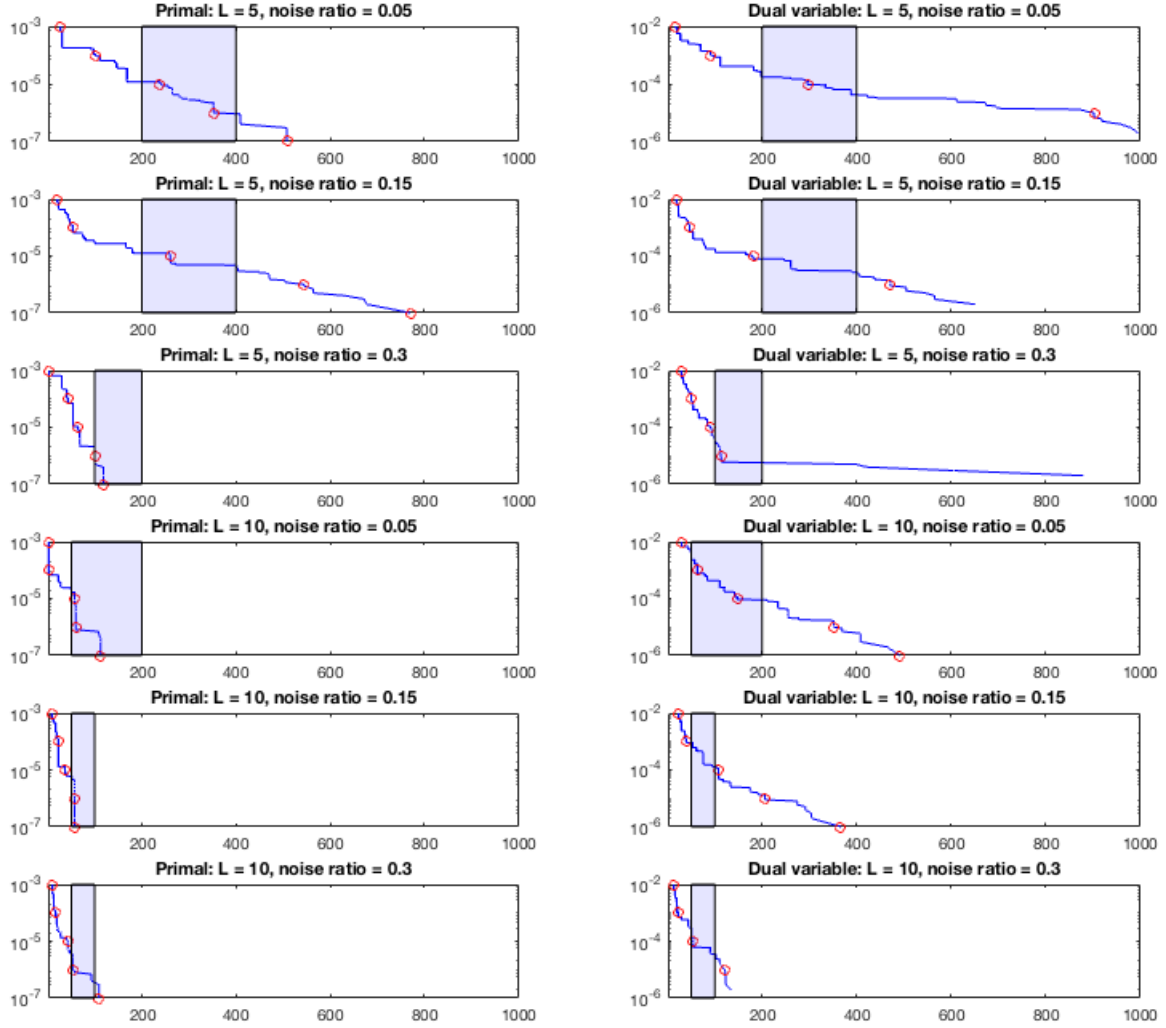


Figure 5.4: Plots of tolerance values against the iterate at which Algorithm 3 first satisfies this tolerance for the models discussed in Figure 5.3. Tolerances depicted are difference values for the primal objective (5.2.7e) and dual variable (5.2.7i). Red circles are placed at tolerances  $10^{-n}$ . The blue rectangles indicate the proposed intervals of termination from Table 5.5.

For each of the models in Figure 5.4, the termination conditions (5.4.1) and (5.4.2) are satisfied within or close to each respective interval of stagnation from Table 5.5. For the models with oversampling  $L = 5$ , Algorithm 3 satisfies the tolerances  $\text{tol}_{\text{primal}} = 10^{-5}$  from

(5.4.1) and  $\text{tol}_{\text{dual}} = 10^{-4}$  from (5.4.2) within or slightly prior to each respective desired interval. Likewise, Algorithm 3 achieve these desired tolerances for the models with  $L = 10$  within or slightly after each desired interval. As an additional precaution, the maximum iteration count of 300 iterates will prevent running Algorithm 3 after it has stagnated.

Note that there is also theoretical justification for selecting the dual variable difference (5.2.7i) as a termination condition. Proposition 3.4.2 showed that the variable  $y$  is optimal for the PLGD model (3.2.1) if  $y = \Pi_{\mathcal{C}}(y - \alpha g)$  for some  $g \in \partial f(y)$  and all  $\alpha > 0$ . This property corresponds to the termination condition

$$\|\Pi_{\mathcal{C}}(y - \alpha g) - y\| \leq \text{tol}.$$

If we make this condition relative by setting  $\text{tol} = 10^{-4} \|\Pi_{\mathcal{C}}(y - \alpha g)\|$ , then we recover the new dual variable difference condition (5.4.2).

One additional concern for termination conditions is the nonmonotonic nature of Algorithm 3. In Figure 5.3, the model with  $L = 10$  and  $\epsilon_{\text{rel}} = 0.05$  demonstrates that Algorithm 3 may produce neighboring signal iterates with varying accuracy. Since this algorithm is nonmonotonic and relies on a subroutine to recover the current approximate signal, the accuracy of the sequence of recovered signals can vary dramatically. Thus we need a reliable indicator to select a sufficiently accurate signal among the previous iterates. Figure 5.5 depicts the signal relative error (5.2.7a) and duality gap (5.2.7f) for the  $L = 10$ ,  $\epsilon_{\text{rel}} = 0.05$  model. Note that the new termination conditions  $\text{tol}_{\text{primal}} = 10^{-5}$  and  $\text{tol}_{\text{dual}} = 10^{-4}$  would select the 148th iterate as the candidate solution.

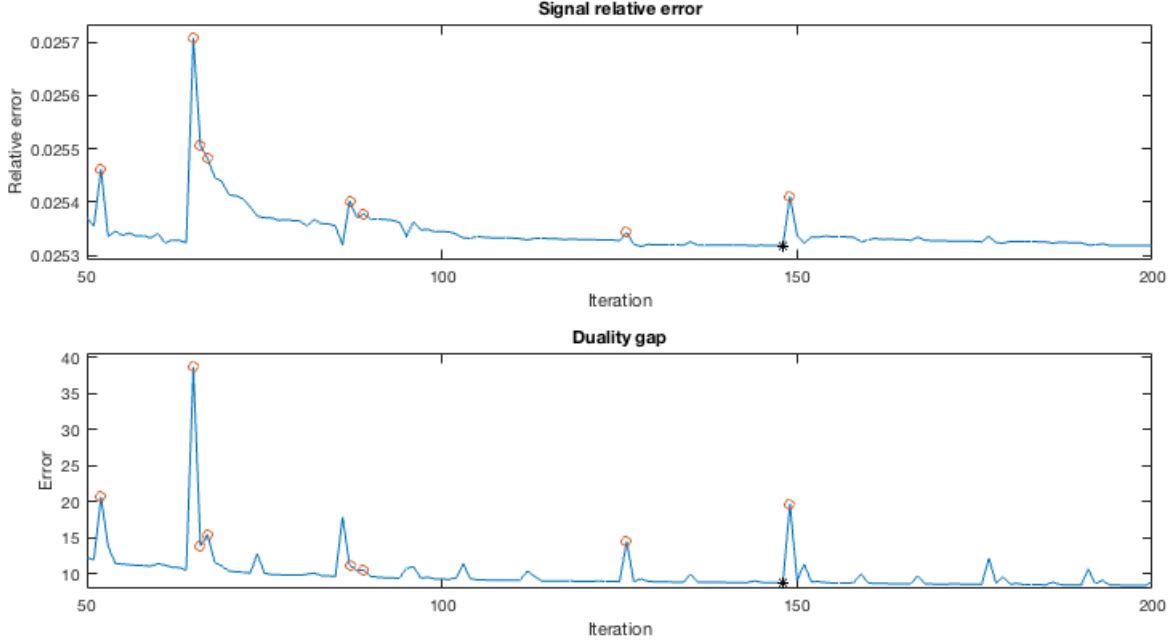


Figure 5.5: Signal relative error (5.2.7a) and duality gap (5.2.7f) for the model from Figure 5.3 with oversampling  $L = 10$  and noise ratio  $\epsilon_{\text{rel}} = 0.05$ . Red circles indicate signals with relative error at least 0.05% above the mean of their ten neighbors, and the black asterisk indicates both the final iterate based on new termination conditions and the optimal iterate based on minimum duality gap.

Figure 5.5 demonstrates that the duality gap (5.2.7f) violation closely matches the signal relative error (5.2.7a). Among the previous iterates, those with the least accurate signal (indicated by the red circle) are accurately identified as having a duality gap value greater than the minimum (black asterisk). Thus once Algorithm 3 terminates, we select the signal among the last 20 with the lowest duality gap value.

Given the new termination conditions (5.4.1), (5.4.2), the maximum iteration count of 300, and the signal selection method discussed above, the Figure 5.6 depicts the iterate at which Algorithm 3 would terminate for each PLGD model with Gaussian noise from Figure 5.3.

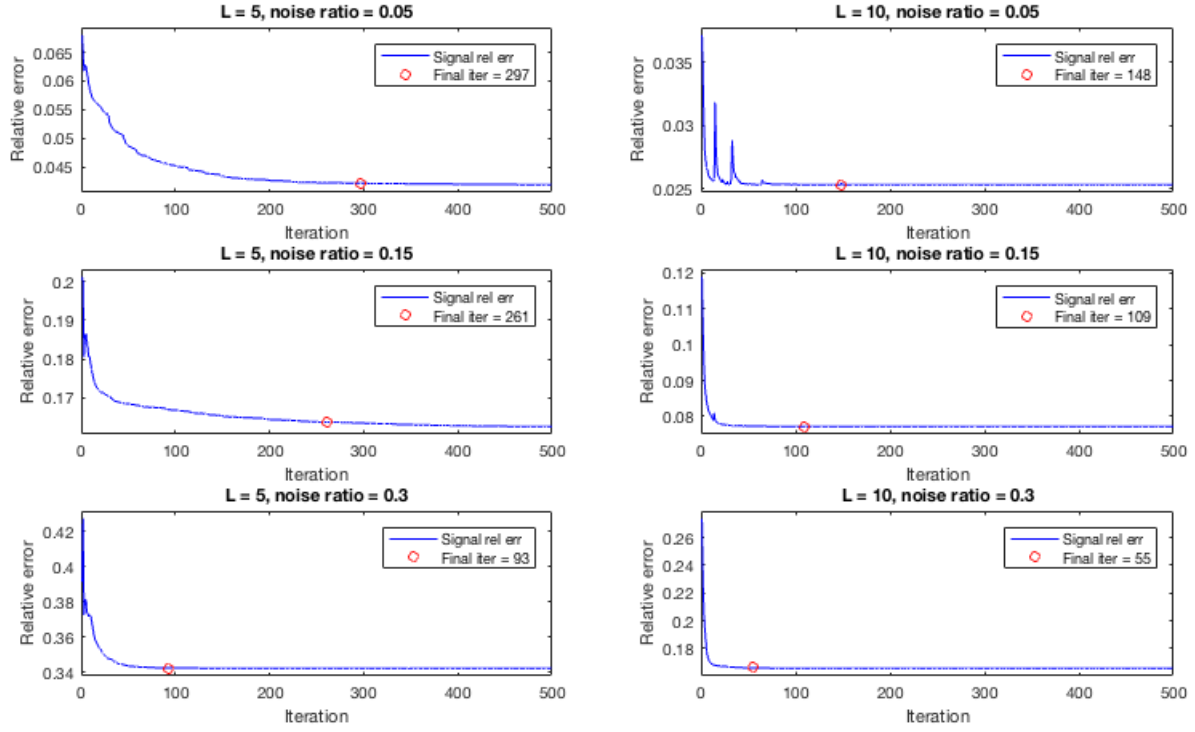


Figure 5.6: Iterate at which Algorithm 3 terminates for the models from Figure 5.3 based on new termination conditions.

For each model in Figure 5.6, Algorithm 3 successfully terminates within or nearby the interval of stagnation established in Table 5.5.

## Chapter 6

# The evolving matrix eigenvalue problem

### 6.1 Introduction

In this chapter we examine the *evolving matrix eigenvalue problem* (EMEP) in Algorithm 3. Section 6.2 defines the EMEP and examines its evolving spectrum across matrix iterates  $A_k$ . In particular, we observe that the algebraically largest eigenvalues tend to cluster as the algorithm proceeds, leading to more difficult eigenvalue problems for later matrix iterates. To handle these eigenvalue problems, we review the *implicitly restarted Arnoldi method* (IRAM) [Sor92] in Section 6.3. The convergence behavior of the IRAM is based on the spectrum of the given eigenvalue problem as well as the IRAM parameters chosen by the user. To understand and exploit the convergence behavior of the IRAM, we will review this method systematically from its component methods.

### 6.2 The EMEP: computational costs and spectral properties

In this section we examine the sequence of eigenvalue problems in Algorithm 3, which we define as the *evolving matrix eigenvalue problem* (EMEP). We will see that the EMEP is the most computationally expensive subroutine in Algorithm 3. We also observe that the EMEP matrix iterates have a spectrum which evolves in a predictable way from early to later iterates.

We begin by formally defining the EMEP. Generally speaking, we are concerned with a sequence of eigenvalue problems in which each eigenvalue problem is dependent on the results of the previous problems. For each iterate  $k$  in this sequence of problems, we have some Hermitian matrix iterate  $A_k \in \mathcal{H}^n$  and seek its  $j$  algebraically largest eigenvalues  $\Lambda_j^{(k)}$  and the corresponding eigenvectors  $V_j^{(k)}$ . Some examples of this problem include subspace tracking in signal processing (see, e.g., [CG90], [Ste92], [Yan95], [DM08]), matrix

completion (e.g., [NS12]), and the Kohn-Sham equation in density functional theory (e.g., [SCS10]).

To see that each eigenvalue problem in Algorithm 3 is dependent on the results of the previous eigenvalue problems, first note that each eigenvalue problem in Algorithm 3 (steps 2, 6, and 14) requires the two algebraically largest eigenvalues of the matrix iterate  $A_k = \mathcal{A}^*y_k$ . (Also note that the iterate  $k$  does not correspond to the iterate number in Algorithm 3 because step 6 involves a linesearch which may involve more than one eigenvalue problem.) We will show that each matrix iterate  $A_k = \mathcal{A}^*y_k$  in Algorithm 3 is computed using a variable  $y_k$  which is dependent on the previous set of basis vectors  $\mathcal{V}_{k-1} = \{v_1^{(i)}\}_{i=0}^{k-1}$ . It can be shown inductively that  $y_k$  is a function of  $\mathcal{V}_{k-1}$ . The initial iterate  $k = 0$  corresponds to the first eigenvalue problem in Algorithm 3, where  $y_0 = \Pi_{\mathcal{C}}(b)$  is initialized using the observation vector  $b$  and is independent of any eigenvector. Thus we initialize  $v_1^{(0)}$  as the empty set. If  $k > 0$  then the update  $y_k$  is computed as  $y_k = \Pi_{\mathcal{C}}(y_{k-1} - \alpha_{k-1}g_{k-1})$ , where  $g_{k-1} = \mathcal{A}(v_{k-1}v_{k-1}^*)$  is a function of  $v_{k-1}$  and  $\alpha_{k-1}$  is determined using a linesearch on the minimization problem

$$\min_{\alpha} \lambda_1(\mathcal{A}^*(\Pi_{\mathcal{C}}(y_{k-1} - \alpha g_{k-1}))). \quad (6.2.1)$$

As a result,  $y_k$  is dependent on  $v_{k-1}$  and  $y_{k-1}$ ; and thus each eigenvalue problem in Algorithm 3 is dependent on the results of the previous eigenvalue problems.

Thus we define the sequence of eigenvalue problems generated by Algorithm 3 as the *evolving matrix eigenvalue problem* (EMEP)

$$\begin{aligned} &\text{for } k = 1, 2, \dots, K \\ &\text{find } \left( \lambda_1^{(k)}, v_1^{(k)} \right) \text{ and } \left( \lambda_2^{(k)}, v_2^{(k)} \right) \text{ of } A_k, \end{aligned} \quad (6.2.2)$$

where  $\lambda_1^{(k)}$  and  $\lambda_2^{(k)}$  are the two algebraically largest eigenvalues of the *matrix iterate*  $A_k = \mathcal{A}^*y_k$ , and  $y_k$  is the previous dual variable generated by Algorithm 3 (from step 2, 6, or 14).

Next, we examine the overall computational costs of the EMEP (6.2.2). As discussed in Section 4.2, the main computational costs in Algorithm 3 are the EMEP (step 2, 6, and 14), the primal refinement (step 11), and the dual refinement (step 13). Since we are focused on PLGD models with Gaussian noise (5.2.3), the dual refinement step of Algorithm 3 is skipped (see the end of Section 5.3 for an explanation). In both the EMEP (6.2.2) and the primal refinement step, the primary computational cost comes from  $\mathcal{A}$ -products (2.3.7), where each  $\mathcal{A}(xx^*)$  product requires  $L$  DFTs and each  $[\mathcal{A}^*y]x$  product requires  $2L$  DFTs. Thus we measure computational costs in terms of the number of DFTs, following the convention of [CLS15] and [FM16]. Also note that the computation of the eigenpair  $(\lambda_1, v_1)$  in the EMEP (6.2.2) must be very accurate in order to determine an accurate descent step  $g = \mathcal{A}(v_1v_1^*)$  in Algorithm 3. Table 6.1 depicts the number of DFTs in Algorithm 3 for a variety of noisy problems.

$n$	$L$	$\epsilon_{\text{rel}}$	EMEP			Primal refinement		All other steps	
			eigs calls	Minutes	DFTs	Minutes	DFTs	Minutes	DFTs
4,096	5	0.05	228	13.13 (0.94)	51,935 (0.97)	0.73 (0.05)	1,516 (0.03)	0.04	17
4,096	5	0.15	120	6.63 (0.94)	31,085 (0.97)	0.45 (0.06)	1,076 (0.03)	0.01	10
4,096	5	0.30	52	3.56 (0.89)	16,410 (0.95)	0.45 (0.11)	854 (0.05)	0.01	4
4,096	10	0.05	190	12.06 (0.96)	72,587 (0.98)	0.45 (0.04)	1,819 (0.02)	0.03	29
4,096	10	0.15	106	8.60 (0.96)	51,450 (0.98)	0.30 (0.03)	1,194 (0.02)	0.02	17
4,096	10	0.30	111	17.95 (0.98)	107,936 (0.99)	0.36 (0.02)	1,420 (0.01)	0.01	18
16,384	5	0.05	199	46.09 (0.95)	69,745 (0.98)	2.13 (0.04)	1,468 (0.02)	0.06	16
16,384	5	0.15	91	27.71 (0.95)	41,880 (0.98)	1.34 (0.05)	853 (0.02)	0.03	8
16,384	5	0.30	61	30.95 (0.94)	45,834 (0.98)	2.04 (0.06)	1,026 (0.02)	0.02	5
16,384	10	0.05	160	56.73 (0.97)	92,391 (0.98)	1.64 (0.03)	1,560 (0.02)	0.07	25
16,384	10	0.15	103	36.30 (0.97)	60,189 (0.98)	1.21 (0.03)	1,167 (0.02)	0.05	17
16,384	10	0.30	47	18.48 (0.96)	30,498 (0.98)	0.65 (0.03)	617 (0.02)	0.02	8

Table 6.1: Algorithm 3 runtime and number of DFTs (with percentage of the total in parentheses) for the EMEP (6.2.2), primal refinement (solving (4.2.11) in step 11) and all other operations. Here  $n$  is signal size (i.e., number of pixels squared in the image from Figure 2.2),  $L$  is number of observations, and  $\epsilon_{\text{rel}}$  is the noise ratio.

The results in Table 6.1 demonstrate the essential computational challenges of Algorithm 3. First, the EMEP is the dominant computational cost in the algorithm, and its proportion to other costs (in both runtime and number of DFTs) increases as the size of the model increases. Additionally, the primal refinement step requires a small but nontrivial amount of computation. All other operations accounted for 0.00% of the overall runtime.

We will now profile the cost of the EMEP (6.2.2). Figure 6.1 depicts the number of matrix-vector products  $[\mathcal{A}^* y_k]x$  in the EMEP (6.2.2) for each of the six smaller models from Table 6.1.

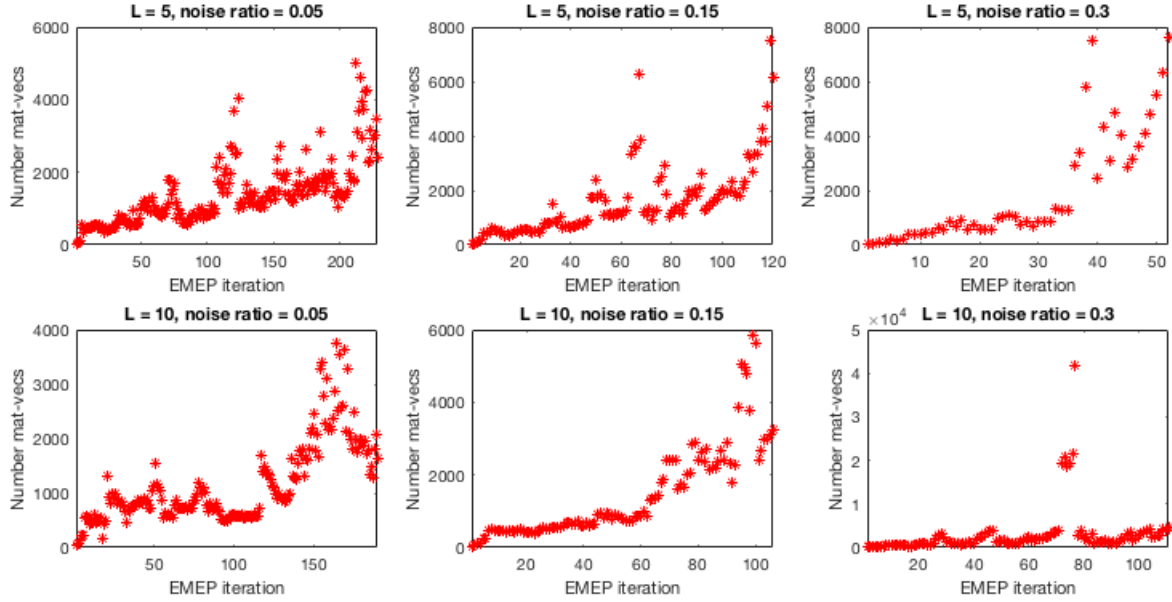


Figure 6.1: Number of matrix-vector products for each iteration in the EMEP (6.2.2) for the six smaller models from Table 6.1 .

Figure 6.1 demonstrates that the number of matrix-vector products required for each EMEP (6.2.2) iterate varies greatly from earlier to later iterates. In each model, the later iterates account for the majority of the computational cost of the EMEP (6.2.2).

To understand why the number of matrix-vector products increases as the EMEP (6.2.2) progresses, we close this section by examining the evolving spectrum distribution of these eigenvalue problems. Figure 6.2 depicts the spectrum of earlier and later iterates for a particular model from Table 6.1.



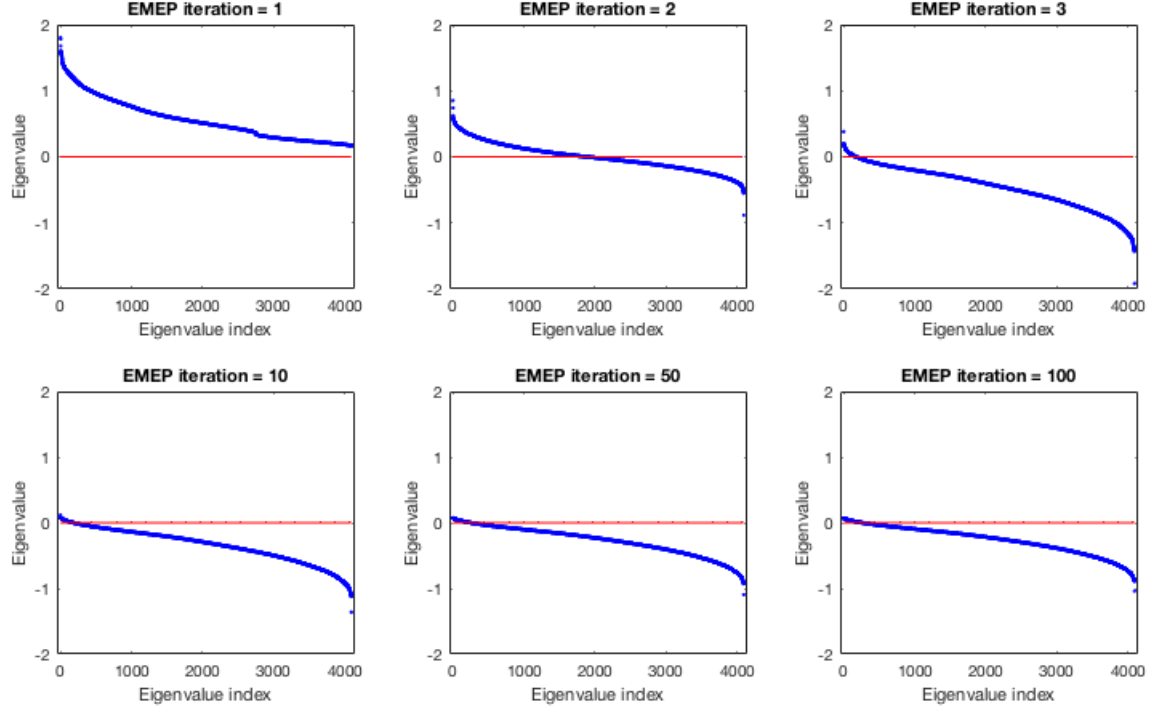


Figure 6.2: Spectrum of specific EMEP (6.2.2) matrix iterates  $A_k$  for the model from Table 6.1 with signal size  $n = 4,096$ , oversampling  $L = 5$ , and noise ratio  $\epsilon_{\text{rel}} = 0.15$ .

As we see in Figure 6.2, the spectrum of the matrix iterates  $A_k$  in the EMEP (6.2.2) shifts from completely positive for  $A_1$  to mostly negative for later iterates. This shift in spectrum is a consequence of optimizing the PLGD model (3.2.1). The first matrix iterate  $A_1 = \mathcal{A}^*b$  will always be positive-semidefinite because the components of the observation  $b = [b_1; b_2; \dots; b_L]$  are all nonnegative and thus for all  $x$  we have

$$x^*[\mathcal{A}^*b]x = \sum_{j=1}^L [FC_j x]^* \text{Diag}(b_j) FC_j x \geq 0.$$

Since Algorithm 3 minimizes the objective function  $\lambda_1(\mathcal{A}^*y_k)$ , the algebraically largest eigenvalue  $\lambda_1^{(k)}$  of  $\mathcal{A}^*y_k$  can be expected to decrease for later iterates  $k$ .

As we will see in Section 6.3, the convergence rate of eigenvalue methods often depends on the distance between the desired eigenvalue  $\lambda_j$  and the next algebraically largest eigenvalue  $\lambda_{j+1}$ . Figure 6.3 depicts the 20 algebraically largest eigenvalues of the EMEP (6.2.2) iterates from Figure 6.2.

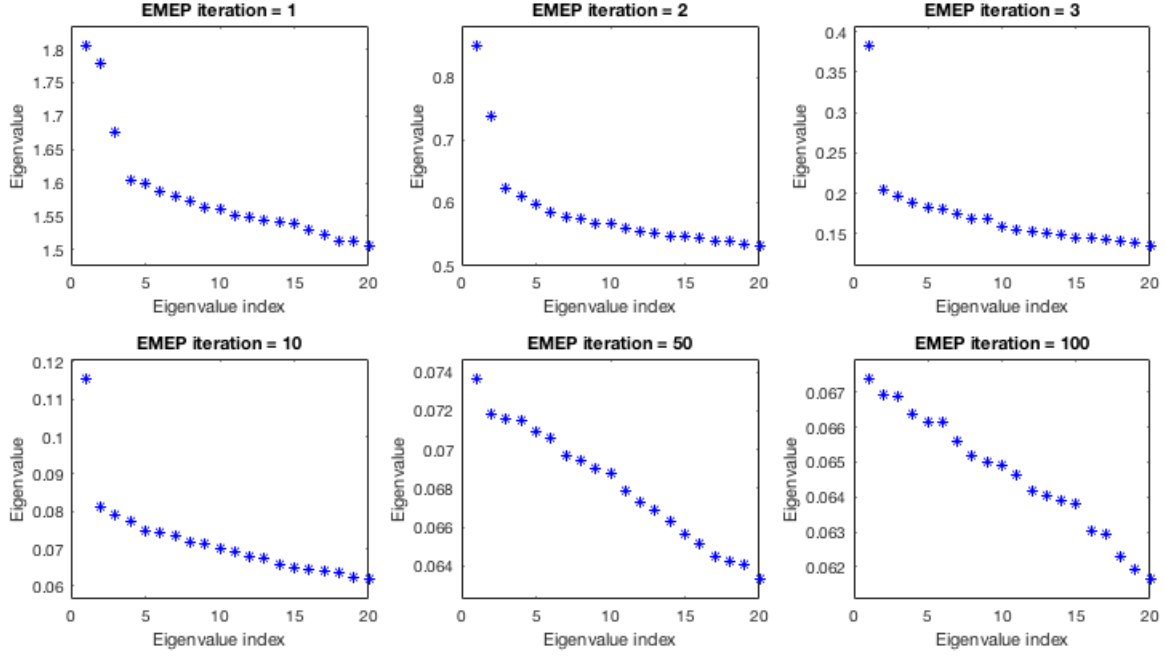


Figure 6.3: Twenty algebraically largest eigenvalues of specific EMEP (6.2.2) matrix iterates  $A_k$  for the model from Table 6.1 with signal size  $n = 4,096$ , oversampling  $L = 5$ , and noise ratio  $\epsilon_{\text{rel}} = 0.15$ .

Figure 6.3 demonstrates that the algebraically largest eigenvalues of the matrix iterates  $A_k$  cluster together as the EMEP (6.2.2) progresses. In general, this clustering in the EMEP (6.2.2) spectrum is expected. Section 5.3 demonstrated that PLGD models with Gaussian noise (5.2.3) typically have optimal primal matrices  $X_\star$  with rank greater than one (see Table 5.3). And Corollary 3.4.1 (e) indicates that  $\text{rank}(X_\star)$  is a lower bound on the multiplicity  $r$  of the algebraically largest eigenvalue of the optimal dual matrix  $\mathcal{A}^*y_\star$ . Thus for later EMEP (6.2.2) iterates  $k$ , we can expect some  $r$  algebraically largest eigenvalues  $\lambda_1^{(k)}, \lambda_2^{(k)}, \dots, \lambda_r^{(k)}$  to have a decreasing relative difference

$$\frac{\lambda_i^{(k)} - \lambda_{i+1}^{(k)}}{\lambda_i^{(k)}}$$

for  $i = 1, 2, \dots, r - 1$ .

The clustering of the algebraically largest eigenvalues as depicted in Figure 6.3 is known to slow the convergence rate of eigenvalue methods like the IRAM. Yet the choice of the parameters in the IRAM can also have a significant effect on the convergence rate of the IRAM. A thorough understanding of the IRAM and its subroutines will help us to develop

a new, adaptive strategy for choosing IRAM parameters for the EMEP (6.2.2) in Chapter 7. Thus we proceed to review the IRAM.

### 6.3 The implicitly restarted Arnoldi method

In this section we review the *implicitly restarted Arnoldi method* (IRAM), a common large-scale eigenvalue method. First proposed by Sorensen [Sor92], [Sor97], the IRAM is a combination of two essential algorithms. The *m-step Arnoldi iteration* is used to build a matrix  $Q_m$  of  $m$  basis vectors which approximates the desired eigenspace. The *p-step shifted QR iteration* restarts the matrix  $Q_m$  with a specific strategy to damp unwanted eigenvalues, resulting in a smaller matrix  $Q_j$  of  $j < m$  basis vectors. Since the *m-step Arnoldi iteration* is an extension of the *power method*, we first discuss the Power method before developing the IRAM. Altogether, the algorithms in this section are presented in the order depicted in Figure 6.4.

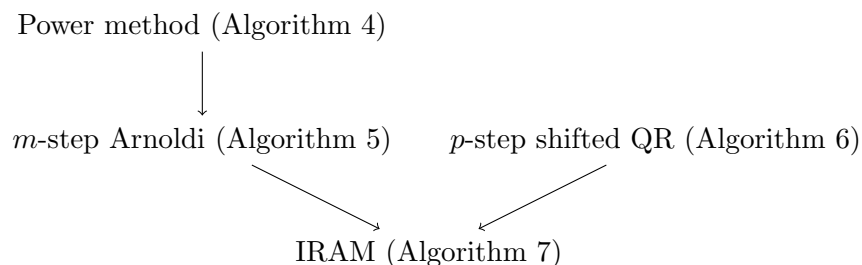


Figure 6.4: Dependency chart for the IRAM.

This section follows the treatment found in [GVL12, Chapters 8, 10], with occasional minor changes in notation.

The first method we consider is the *power method*, a method for determining the largest magnitude eigenvalue  $\lambda_1$  and corresponding eigenvector  $v_1$  of a Hermitian matrix  $A$ . The power method is based on the property that if  $\lambda_1$  is strictly larger in magnitude than the next largest magnitude eigenvalue and the initial vector  $q^{(0)}$  has a nonzero component in the direction of  $v_1$  (i.e.,  $v_1^* q^{(0)} \neq 0$ ), then the sequence

$$q^{(0)}, \frac{Aq^{(0)}}{\|Aq^{(0)}\|}, \frac{A^2q^{(0)}}{\|A^2q^{(0)}\|}, \frac{A^3q^{(0)}}{\|A^3q^{(0)}\|}, \dots$$

will have  $v_1$  as its limit.

Formally, Algorithm 4 presents the power method as seen in [GVL12, Section 8.2.1].

---

**Algorithm 4** Power method

---

**Input:** Hermitian matrix  $A$ , initial approximate eigenvector  $q^{(0)}$ , relative tolerance  $\text{tol}_{\text{rel}} > 0$ .

**Output:** Approximate largest magnitude eigenvalue  $\lambda$  and the corresponding eigenvector  $v$ .

- 1: *Initialize:*  $q^{(0)} = q^{(0)} / \|q^{(0)}\|$ ,  $\rho^{(0)} = [q^{(0)}]^* A q^{(0)}$ ,  $r^{(0)} = A u^{(0)} - \rho^{(0)} q^{(0)}$ ,  $i = 1$ .
  - 2: **while** not converged:  $\|r^{(i)}\| / (\|A q^{(i)}\| + |\rho^{(i)}|) > \text{tol}_{\text{rel}}$  **do**
  - 3:      $z^{(i)} = A q^{(i-1)}$
  - 4:      $q^{(i)} = z^{(i)} / \|z^{(i)}\|$
  - 5:      $\rho^{(i)} = [q^{(i)}]^* z^{(i)}$
  - 6:      $r^{(i)} = A q^{(i)} - \rho^{(i)} q^{(i)}$ ,  $i = i + 1$
  - 7: **end while**
  - 8: *Return:*  $(\lambda, v) = (\rho^{(i-1)}, q^{(i-1)})$ .
- 

The simplicity of power method allows for insightful convergence results like the Theorem 6.3.1, in which we assume the matrix  $A$  is real for clarity.

**Theorem 6.3.1.** *Suppose  $A \in \mathbb{R}^{n \times n}$  is symmetric with an eigenvalue decomposition*

$$V^* A V = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_n),$$

*where  $V = [v_1 \mid v_2 \mid \dots \mid v_n]$  is orthogonal and  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$ . Let the vectors  $q^{(i)}$  be generated by Algorithm 4 and define  $\theta_i \in [0, \pi/2]$  as*

$$\cos(\theta_i) = \left| v_1^T q^{(i)} \right|.$$

*If  $\cos(\theta_0) \neq 0$ , then for  $i = 0, 1, \dots$  we have*

$$|\sin(\theta_i)| \leq \tan(\theta_0) \left| \frac{\lambda_2}{\lambda_1} \right|^i, \quad (6.3.1)$$

$$\left| \lambda^{(i)} - \lambda_1 \right| \leq \max_{2 \leq j \leq n} |\lambda_1 - \lambda_j| \tan(\theta_0)^2 \left| \frac{\lambda_2}{\lambda_1} \right|^{2i}. \quad (6.3.2)$$

*Proof.* See [GVL12, Theorem 8.2.1]. □

Theorem 6.3.1 establishes that the convergence rate of the power method (Algorithm 4) is dependent on the distance between  $|\lambda_1|$  and  $|\lambda_2|$ . If this distance  $\epsilon = |\lambda_1| - |\lambda_2|$  is very small relative to  $|\lambda_1|$ , then we have

$$\left| \frac{\lambda_2}{\lambda_1} \right| = \frac{|\lambda_1| - \epsilon}{|\lambda_1|} = 1 - \frac{\epsilon}{|\lambda_1|} \approx 1,$$

and  $|\sin(\theta_i)|$  in (6.3.1) may decreases very slowly.

The next method we consider is the *m-step Arnoldi iteration* which extends the power method (Algorithm 4) to achieve a superior convergence rate. An iteration of the power method generates a new approximate eigenvector ( $q^{(i)}$  from steps 3 and 4) by normalizing the matrix-vector product of the previous vector. In essence, the power method searches for the largest magnitude eigenvalue  $\lambda_1$  and corresponding eigenvector  $v_1$  of a matrix  $A$  in the one-dimensional subspace  $V = \text{span}\{Aq_1\}$ . The *m-step Arnoldi iteration* extends the power method by searching for the Ritz pair (1.2.10)  $(\theta_1, u_1)$  for  $A$  with respect to the *m-dimensional Krylov subspace*

$$\mathcal{K}_m(A, q_1) = \text{span}\{q_1, Aq_1, A^2q_1, \dots, A^{m-1}q_1\}. \quad (6.3.3)$$

Algorithm 5 (as described in [GVL12, Algorithm 10.5.1]) builds a unitary basis  $Q_m$  of  $\mathcal{K}_m(A, q_1)$  which may be used to find the Ritz pair  $(\theta_1, u_1)$ .

---

**Algorithm 5** *m-step Arnoldi iteration*

---

**Input:** Matrix  $A \in \mathbb{C}^{n \times n}$ , number of Arnoldi steps  $m$ , initial approximate eigenvector  $q_1$ .

**Output:** Hessenberg matrix  $H_m$ , basis  $Q_m$ , residual  $r_m$ .

- 1: *Initialize:*  $q_1 = q_1/||q_1||$ ,  $z = Aq_1$ ,  $\alpha_1 = q_1^*z$ ,  $r_1 = z - \alpha_1q_1$ ,  $Q_1 = [q_1]$ ,  $H_1 = [\alpha_1]$ .
  - 2: **for**  $i = 1, \dots, m-1$  **do**
  - 3:    $\beta_i = ||r_i||$ ,  $q_{i+1} = r_i/\beta_i$ .
  - 4:    $Q_{i+1} = [Q_i \mid q_{i+1}]$ ,  $\hat{H}_i = \begin{bmatrix} H_i \\ \beta_i e_i^T \end{bmatrix}$ .
  - 5:    $z = Aq_{i+1}$ .
  - 6:    $h = Q_{i+1}^*z$ ,  $r_{i+1} = z - Q_{i+1}h$ .
  - 7:    $H_{i+1} = [\hat{H}_i \mid h]$ .
  - 8: **end for**
  - 9: *Return:*  $H_m, Q_m, r_m$ .
- 

In order to obtain a Ritz pair  $(\theta, u)$  for  $A$  with respect to  $\mathcal{K}_m(A, q_1)$ , the *m-step Arnoldi iteration* generates an *m-step Arnoldi decomposition*

$$AQ_m = Q_m H_m + r_m e_m^*, \quad (6.3.4)$$

where  $H_m$  is an upper Hessenberg matrix. If  $(\theta, w)$  is an eigenpair for  $H_m$  and  $u = Q_m w$  then (6.3.4) implies

$$(AQ_m - Q_m H_m)w = (A - \theta I)u = (e_m^* w)r_m. \quad (6.3.5)$$

Additionally, steps 5 and 6 of Algorithm 5 indicate that  $r_m$  is orthogonal to  $\mathcal{K}_m(A, q_1)$ , and thus  $(\theta, u)$  is a Ritz pair for  $A$  with respect to  $\mathcal{K}_m(A, q_1)$ .

The use of  $\mathcal{K}_m(A, q_1)$  in Algorithm 5 allows for superior convergence to Algorithm 4. Note that the largest magnitude Ritz pair  $(\theta_1, u_1)$  for  $A$  with respect to  $\mathcal{K}_m(A, q_1)$  generated by Algorithm 5 is guaranteed to be at least comparable to the  $m$ -th iterate of Algorithm 4 since  $A^{m-1}q_1 \in \mathcal{K}_m(A, q_1)$ . To compare the convergence rates of Algorithms 4 and 5 more precisely, assume the matrix  $A$  is real and symmetric. Then the matrix  $H_m$  returned by Algorithm 5 is tridiagonal and this algorithm is equivalent to the  $m$ -step *Lanczos iteration* [GVL12, Algorithm 10.1.1]. In this case, we have the Theorem 6.3.2. Note that this theorem involves *Chebyshev polynomials* [Saa11, Section 4.4], a sequence of polynomials defined recursively as

$$c_k(x) = 2xc_{k-1}(x) - c_{k-2}(x) \quad (6.3.6)$$

for  $k \geq 2$ , with  $c_0 = 1$  and  $c_2 = x$ .

**Theorem 6.3.2.** *Let  $A \in \mathbb{R}^{n \times n}$  be symmetric with an eigenvalue decomposition*

$$V^*AV = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_n),$$

*where  $V = [v_1 \mid v_2 \mid \dots \mid v_n]$  is orthogonal and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ . Suppose the  $m$ -step Arnoldi iteration (Algorithm 5) is performed and  $H_k$  is the tridiagonal matrix returned by this algorithm. If  $\theta_1$  is the algebraically largest eigenvalue of  $H_m$ , then*

$$\lambda_1 \geq \theta_1 \geq \lambda_1 - (\lambda_1 - \lambda_n) \left( \frac{\tan(\phi_1)}{c_{m-1}(1 + 2\rho_1)} \right)^2, \quad (6.3.7)$$

*where  $\cos(\phi_1) = |q_1^T v_1|$ ,*

$$\rho_1 = \frac{\lambda_1 - \lambda_2}{\lambda_2 - \lambda_n}, \quad (6.3.8)$$

*and  $c_{m-1}(x)$  is the Chebyshev polynomial of degree  $m - 1$ .*

*Proof.* See [GVL12, Theorem 10.1.2]. □

The convergence rate established in Theorem 6.3.2 may also be applied to Algorithm 4, giving the Corollary 6.3.3.

**Corollary 6.3.3.** *Let  $A \in \mathbb{R}^{n \times n}$  be symmetric and positive semidefinite with an eigenvalue decomposition*

$$V^*AV = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_n),$$

*where  $V = [v_1 \mid v_2 \mid \dots \mid v_n]$  is orthogonal and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ . Suppose  $m$  steps of the power method (Algorithm 4) are performed and  $\gamma_1 = \rho^{(m)}$  is the returned Ritz value. Then*

$$\lambda_1 \geq \gamma_1 \geq \lambda_1 - (\lambda_1 - \lambda_n) \tan^2(\phi_1) \left( \frac{\lambda_2}{\lambda_1} \right)^{2(m-1)}, \quad (6.3.9)$$

*where  $\cos(\phi_1) = |q_1^T v_1|$ .*

*Proof.* See [GVL12, Theorem 10.1.2] and replace the Chebyshev polynomial in this proof with  $p(x) = x^{k-1}$ .  $\square$

The lower bounds in Theorem 6.3.2 and Corollary 6.3.3 may be used to compare the expected convergence rates for Algorithms 4 and 5. The following comparison is based on [GVL12, Section 10.1.6]. Assume  $A \in \mathbb{R}^{n \times n}$  is symmetric and also positive semidefinite for clarity. Assume Algorithms 4 and 5 have been run for  $m$  steps with the same initial vector  $q_1$ . Let  $\gamma_1 = \rho^{(m)}$  be the Ritz value for  $A$  generated by step 5 of Algorithm 4. And let  $\theta_1$  be the Ritz value for  $A$  with respect to  $\mathcal{K}_m(A, q_1)$  generated by the algebraically largest eigenvalue of  $H_m$  from Algorithm 5. Then we may compare the lower bounds (6.3.9) for  $\gamma_1$  and (6.3.7) for  $\theta_1$  by comparing the values

$$P_{m-1} = \left( \frac{\lambda_2}{\lambda_1} \right)^{2(m-1)}, \quad (6.3.10)$$

$$L_{m-1} = \frac{1}{\left[ c_{m-1} \left( 2 \frac{\lambda_1}{\lambda_2} - 1 \right) \right]^2} \geq \frac{1}{[c_{m-1} (1 + 2\rho_1)]^2}. \quad (6.3.11)$$

Table 6.2 compares  $P_{m-1}$  and  $L_{m-1}$  for a few values of  $m$  and  $\lambda_1/\lambda_2$ .

$\lambda_1/\lambda_2$	$m = 10$		$m = 20$	
	$P_{m-1}$	$L_{m-1}$	$P_{m-1}$	$L_{m-1}$
1.10	$1.8 \times 10^{-1}$	$5.5 \times 10^{-5}$	$2.7 \times 10^{-2}$	$2.1 \times 10^{-10}$
1.01	$8.4 \times 10^{-1}$	$1.0 \times 10^{-1}$	$6.9 \times 10^{-1}$	$2.0 \times 10^{-3}$

Table 6.2: Lower bound terms (6.3.10) and (6.3.11) for Ritz values generated by Algorithms 4 and 5

Table 6.2 demonstrates that the use of the Krylov subspace  $\mathcal{K}_m(A, q_1)$  in Algorithm 5 allows for superior convergence to Algorithm 4. Yet this superior convergence rate is slowed somewhat when the desired eigenvalue  $\lambda_1$  is close to  $\lambda_2$ . For eigenvalue problems where the value  $\lambda_1/\lambda_2$  is very small (like later iterates of the EMEP (6.2.2), as demonstrated in Figure 6.3), we may seek to increase the number of steps  $m$  in Algorithm 5. Yet increasing  $m$  can be computationally prohibitive if the eigenvalue problem is very large, requiring significant memory to store  $Q_m$  and significant computation to compute the eigenvalue decomposition of  $H_m$ .

To take advantage of the convergence rate of Algorithm 5 for larger eigenvalue problems, the Arnoldi decomposition (6.3.4) may be restarted with the *p-step shifted QR iteration* developed by [Sor92] and discussed in [GVL12, Sections 10.5.2-3]. To develop this algorithm, assume we are seeking the  $j$  algebraically largest eigenvalues of a Hermitian matrix  $A \in \mathbb{C}^{n \times n}$  and we require that the  $m$ -step Arnoldi decomposition (6.3.4)  $AQ_m = Q_m H_m + r_m e_m^*$  has a fixed size  $m > j$ . First we run Algorithm 5 with the initial

vector  $q_1$  to obtain  $AQ_m = Q_m H_m + r_m e_m^*$ . Next, recall that the matrix  $H_m$  may be used to identify the desired Ritz pairs  $\{(\theta_i, u_i)\}_{i=1}^j$  for  $A$  with respect to  $\mathcal{K}_m(A, q_1)$ , as described in (6.3.5). Yet  $H_m$  also contains Ritz values  $\theta_{j+1}, \dots, \theta_m$  which correspond to unwanted eigenvalues of  $A$ . To damp these unwanted Ritz values, we may select an appropriate degree  $p = m - j$  filter polynomial  $p(\lambda)$ . The  $p$ -step shifted QR iteration uses the filter polynomial

$$p(\lambda) = c \cdot (\lambda - \mu_1)(\lambda - \mu_2) \cdots (\lambda - \mu_p), \quad (6.3.12)$$

where  $c$  is a constant and the shift values  $\mu_1 = \theta_{j+1}, \dots, \mu_p = \theta_m$  are the  $p$  unwanted Ritz values of  $A$  with respect to  $\mathcal{K}_m(A, q_1)$ . Algorithm 6 (as described in [GVL12, Section 10.5.3]) uses the filter polynomial (6.3.12) implicitly by applying  $p$  shifted QR steps to  $H_m$ .

---

**Algorithm 6**  $p$ -step shifted QR iteration (implicit polynomial filtering)

---

**Input:** Hessenberg matrix  $H_m \in \mathbb{C}^{m \times m}$  and shift values  $\mu_1, \dots, \mu_p$ .

**Output:** Processed Hessenberg matrix  $H_m^+ \in \mathbb{C}^{m \times m}$  and change of basis  $V \in \mathbb{C}^{m \times m}$ , with  $H_m^+ = V^* H_m V$ .

- 1: Set  $H^{(0)} = H_m$ .
  - 2: **for**  $i = 1, \dots, p$  **do**
  - 3:     QR factorization:  $H^{(i-1)} - \mu_i I = V_i R_i$ .
  - 4:     Update:  $H^{(i)} = R_i V_i + \mu_i I$ .
  - 5: **end for**
  - 6: Set  $H_m^+ = H^{(p)}$ ,  $V = V_1 V_2 \cdots V_p$ .
  - 7: Return:  $H_m^+, V$ .
- 

Proposition 6.3.1 establishes that Algorithm 6 implicitly applies the filter polynomial  $p(\lambda)$  from (6.3.12) to the initial vector  $q_1$  used to create the  $m$ -step Arnoldi decomposition (6.3.4).

**Proposition 6.3.1.** *Let  $A \in \mathbb{C}^{n \times n}$  be Hermitian and  $AQ_m = Q_m H_m + r_m e_m^*$  be the  $m$ -step Arnoldi decomposition (6.3.4) returned by Algorithm 5 with initial vector  $q_1$ . And let  $\mu_1, \dots, \mu_p$  be the  $p$  smallest algebraic eigenvalues of  $H_m$ . Run Algorithm 6 with  $H_m$  and  $\mu_1, \dots, \mu_p$  as inputs, and return  $H_m^+$ ,  $V = V_1 \cdots V_p$ , and  $R = R_p \cdots R_1$ .*

*Then the restarted matrix  $Q_+ = Q_m V$  will have the first column*

$$q_+ = Q_m V(:, 1) = p(A)q_1,$$

*where  $p(\lambda)$  is the filter polynomial (6.3.12) with constant  $c = 1/R(1, 1)$ .*

*Proof.* See [GVL12, Section 10.5.3]. □

After performing Algorithm 6, we have the transformed  $m$ -step Arnoldi decomposition (6.3.4)

$$AQ_+ = Q_+ H_+ + r_m e_m^* V, \quad (6.3.13)$$



where  $V = V_1 \cdots V_p$  from Algorithm 6 and  $Q_+ = Q_m V$ . As a consequence of the QR steps used in Algorithm 6, we can also show that the first  $j = m - p$  columns of (6.3.13) form a new  $j$ -step Arnoldi decomposition. Note that  $V_1, \dots, V_p$  are all upper Hessenberg due to the QR factorization in step 3 of Algorithm 6. Then  $V$  has a lower band  $p$  and  $V(m, 1 : m - p - 1) = V(m, 1 : j - 1) = 0$ , giving

$$r_m e_m^* V(:, 1 : j) = V(m, j) r_m e_j^*. \quad (6.3.14)$$

Also,  $H_+$  is upper Hessenberg and thus  $H_+(j + 1 : m, 1 : j) = H_+(j + 1, j) e_1 e_j^*$ , giving

$$\begin{aligned} Q_+ H_+(:, 1 : j) &= Q_+(:, 1 : j) H_+(1 : j, 1 : j) + Q_+(:, j + 1 : m) H_+(j + 1 : m, 1 : j) \\ &= Q_+(:, 1 : j) H_+(1 : j, 1 : j) + H_+(j + 1, j) Q_+(:, j + 1) e_j^*. \end{aligned} \quad (6.3.15)$$

Therefore, if we set  $Q_j = Q_+(:, 1 : j) = Q_m V(:, 1 : j)$ ,  $H_j = H_+(1 : j, 1 : j)$ , and  $r_j = V(m, j) r_m + H_+(j + 1, j) Q_+(:, j + 1)$ , then equations (6.3.13-6.3.15) give the new  $j$ -step Arnoldi decomposition

$$\begin{aligned} A Q_j &= A Q_+(:, 1 : j) \\ &= Q_+(:, 1 : j) H_+(1 : j, 1 : j) + [V(m, j) r_m + H_+(j + 1, j) Q_+(:, j + 1)] e_j^* \\ &= Q_j H_j + r_j e_j^*, \end{aligned} \quad (6.3.16)$$

and we may resume Algorithm 5 at step  $j + 1$ .

Combining Algorithms 5 and 6 as described above, we have Algorithm 7 as presented in [GVL12, Section 10.5.3].

---

**Algorithm 7** Implicitly restarted Arnoldi method (IRAM)

---

- Input:** Matrix  $A \in \mathbb{C}^{n \times n}$ , initial approximate eigenvector  $q_1$ , number of requested algebraically largest eigenvalues  $j$ , maximum Arnoldi decomposition (6.3.4) size  $m$ .
- Output:** Approximate algebraically largest eigenpairs  $(\Lambda_j, V_j)$ .
- 1: *Initialize with Algorithm 5:* Perform the  $m$ -step Arnoldi iteration with initial vector  $q_1$  to obtain  $AQ_m = Q_m H_m + r_m e_m^*$ .
  - 2: **while** not converged **do**
  - 3:   Compute the eigenvalues  $\theta_1, \dots, \theta_m$  of  $H_m$  and identify the  $p = m - j$  (unwanted) shift values  $\mu_1 = \theta_{j+1}, \dots, \mu_p = \theta_m$ .
  - 4:   *Algorithm 6:* Perform the  $p$ -step shifted QR iteration to obtain the Hessenberg matrix  $H_+$  and change of basis  $V$ .
  - 5:   *Restart the Arnoldi factorization:* Set  $Q_j = Q_m V(:, 1 : j)$ ,  $H_j = H_+(1 : j, 1 : j)$ , and  $r_j = V(m, j)r_m + H_+(j + 1, j)Q_+(:, j + 1)$  per (6.3.16).
  - 6:   *Algorithm 5:* Beginning with  $AQ_j = Q_j H_j + r_j e_j^*$ , perform steps  $j + 1, \dots, m$  of the Arnoldi iteration to obtain  $AQ_m = Q_m H_m + r_m e_m^*$ .
  - 7: **end while**
  - 8: Compute the  $j$  algebraically largest eigenvalues  $\Lambda_j = \{\lambda_1, \dots, \lambda_j\}$  and corresponding eigenvectors  $u_1, \dots, u_j$  of  $H_m$ . Set  $V_j = [Q_m u_1 \mid \dots \mid Q_m u_j]$ .
  - 9: *Return:*  $(\Lambda_j, V_j)$ .
- 

The IRAM is the eigenvalue method we use in Chapter 7 to handle the EMEP (6.2.2). The choice of parameters  $m$  (the Arnoldi decomposition size) and  $j$  (number of requested eigenvalues) can greatly impact the efficiency of IRAM (see Section 7.2). For many large-scale eigenvalue problems, the IRAM is a very effective and convenient method. Due to the implicit polynomial filtering in step 4 of IRAM, this method is particularly effective when the  $j$  algebraically largest eigenvalues have modest separation from  $\lambda_{j+1}$ . And since the IRAM only has two parameter choices, there is little optimization required by the user.

However, when  $\lambda_j \approx \lambda_{j+1}$  and the Arnoldi decomposition size  $m$  is not sufficiently large, the IRAM may require many iterations to achieve the desired tolerance. As we will see in Section 7.2, the appropriate choice of  $m$  and  $j$  in this circumstance may make the IRAM far more competitive. Yet choosing  $m$  and  $j$  without prior knowledge of the eigenvalue distribution is inherently heuristic. Additionally, if the inputted matrix  $A$  is very large, then it may be prohibitive to store the  $Q_m \in \mathbb{C}^{n \times m}$  in active memory. In particular, if the image or signal  $x$  being recovered in the PLGD problem has  $n$  pixels, then  $Q_m$  will require  $m$ -times as much storage space. Thus we proceed in the next section by considering an alternative Krylov subspace method which does not require parameter tuning, nor a large subspace to be held in memory.

The numerical software package **ARPACK** (the **AR**noldi **PACK**age) is an implementation of IRAM in FORTRAN 77 [LSY98]. Many numerical computing environments include large-scale eigenvalue methods which having bindings to ARPACK, including **eigs**

in MATLAB, `eigs` and `eigsh` in the Python package SciPy, `eigs` in R, and `eigs` in the Julia package Arpack.jl.

## Chapter 7

# The adaptive parameter method for the EMEP

### 7.1 Introduction

This chapter presents the *adaptive parameter method* for the EMEP (6.2.2), a new strategy for solving the EMEP with the IRAM (Algorithm 7). Section 7.2 develops the adaptive parameter method by examining the changing behavior of the IRAM across EMEP iterates and states the formal algorithm for the adaptive parameter method. Section 7.3 examines a few EMEPs more closely to demonstrate that the adaptive parameter method effectively increases the number of requested eigenvalues  $r$  as the algebraically largest eigenvalues begin to cluster, thus allowing the IRAM to converge more quickly. Section 7.4 then presents a variety of numerical results comparing the adaptive parameter method with the original parameters used for the EMEP. As compared with the original parameters, we see that the adaptive parameter method decreases the number of matrix-vector products by 50-90% for EMEPs with minimal oversampling in the phase retrieval problem (2.1.1).

All EMEP (6.2.2) experiments in this section are performed with Algorithm 3 using the new termination conditions (5.4.1) and (5.4.2) established in Chapter 5. All experiments in this chapter are available for reproduction.<sup>1</sup>

### 7.2 The adaptive parameter method for the EMEP

In this section we develop a new strategy for solving the EMEP (6.2.2). This strategy uses the IRAM (Algorithm 7) to handle each EMEP matrix iterate  $A_k$  while adaptively changing one of the IRAM parameters based on the results from the EMEP iterates. As discussed in Section 6.3, the IRAM has only two key parameters: the number of requested eigenvalues

---

<sup>1</sup><https://github.com/Will-Wright/low-rank-opt-rapid-eig>

$r$  and the Arnoldi decomposition (6.3.4) size  $m$ . As we will see, the proper choice of these parameters can greatly reduce the number of matrix-vector products required for the EMEP.

We begin by examining the change in computational costs (as measured by matrix-vector products) for various EMEP (6.2.2) iterates and IRAM parameters  $r$  in Figure 7.1. In the original implementation of Algorithm 3, all EMEP iterates were handled using the IRAM with  $r = 2$  requested eigenvalues and Arnoldi decomposition size  $m = \min\{\max\{2r, 20\}, n\}$ , where  $n$  is the size of the desired signal  $x$ . This choice of  $m$  is equivalent to the default parameter setting in the IRAM solver `eigs` for MATLAB and evaluates to  $m = 20$  for  $n \geq 20$ . Yet the plots in Figure 7.1 demonstrate that choosing a fixed parameter  $r = 2$  can result in far more matrix-vector products than is optimal.

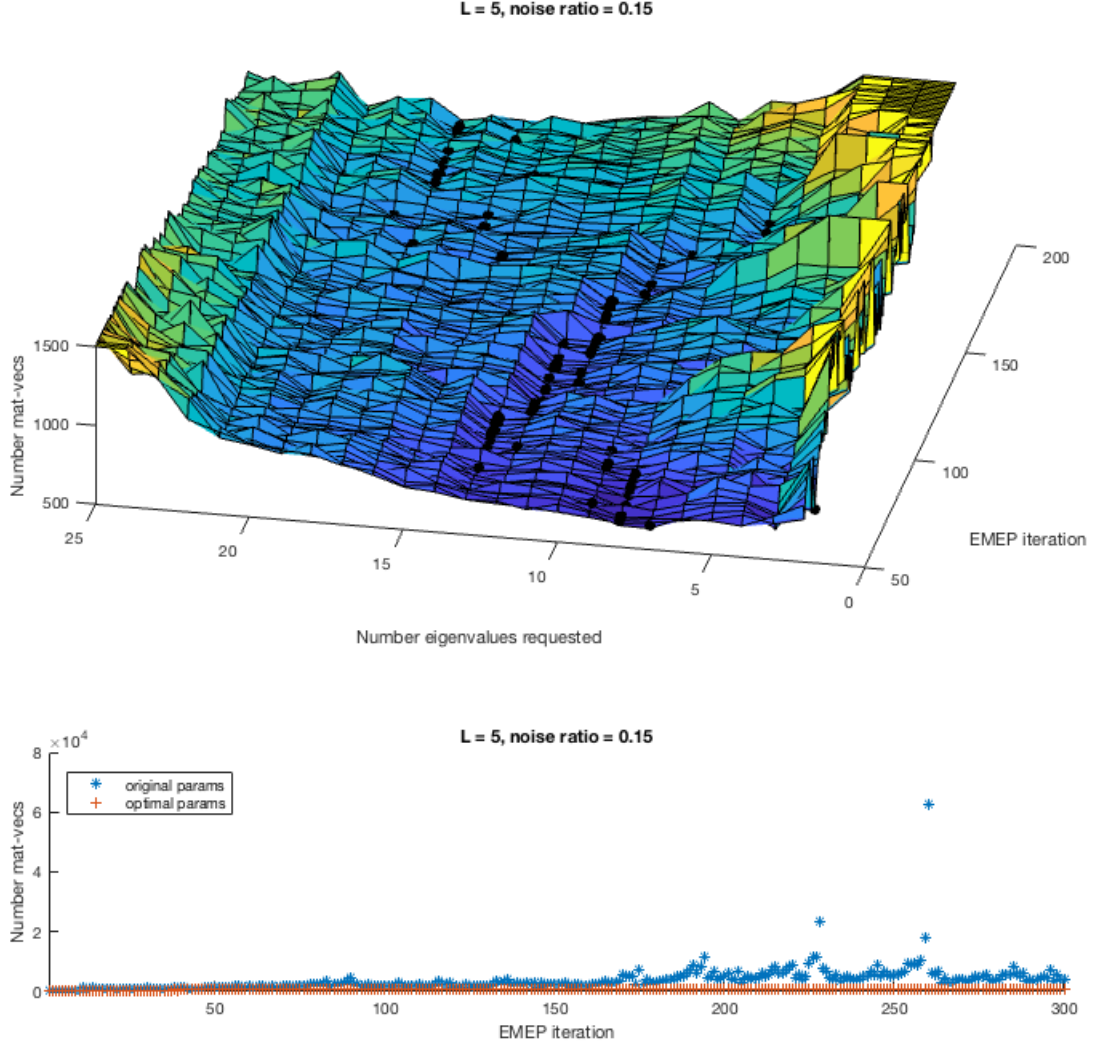


Figure 7.1: Number of matrix-vector products for an EMEP (6.2.2) with various IRAM parameters. Top: Number of matrix-vector products (capped at 1,500 for better viewing) for various EMEP iterates and number of requested eigenvalues  $r$ . Arnoldi decomposition size is set to  $m = 40$  and black dots indicate the *optimal* parameter  $r$  with the minimum number of products for each EMEP iterate. Bottom: Plot of IRAM results for the EMEP with optimal parameters from top plot and fixed parameters  $r = 2$  and  $m = 20$ . The EMEP is from a PLGD model with Gaussian noise (5.2.3) with noise ratio  $\epsilon_{\text{rel}} = 0.15$ , oversampling rate  $L = 5$ , and original signal from Figure 2.2 resized to  $64 \times 64$  pixels.

We now examine Figure 7.1 to develop an adaptive strategy for choosing the number

of requested eigenvalues  $r$  for the sequence of EMEP iterates. The top plot in Figure 7.1 shows that the optimal choice of parameter  $r$  typically changes only slightly between EMEP iterates. However, the optimal parameter  $r$  can increase quickly for later EMEP iterates (as we see around iterate 150 in the top plot in Figure 7.1). Based on these observations, we may develop a strategy for choosing a sequence of parameters  $r_1, r_2, \dots, r_{maxit}$  using two basic heuristics. Assume that each pair of parameters  $r_{i-1}$  and  $r_i$  differ by at least one. First, we compare the two most recent choices for  $r$  and the resulting number of matrix-vector products. If these two choices for  $r$  decreased the number of matrix-vector products, we continue to shift the value of  $r$  in this direction by one unit; and otherwise we shift  $r$  in the opposite direction. Next, we compare the four most recent choices for  $r$  and the resulting number of matrix-vector products using linear interpolation. If these recent choices suggest the same shift as the first comparison, then we shift  $r$  in this direction by two units rather than one.

Now we will formally develop an adaptive strategy for choosing the number of requested eigenvalues  $r_k$  for each EMEP (6.2.2) iterate  $k$ . First, we select a fixed Arnoldi decomposition size  $m$  and initialize  $r_1 = r_{min}$  and  $r_2 = r_1 + 1$  (where  $r_{min} = 2$ , and  $r_{max} = \min\{30, m - 5\}$ , and we set the default value  $m = 40$ ). At each step  $k \geq 2$  we update  $r_{k+1} = r_k + \delta$ , where  $\delta \in \{-2, -1, 1, 2\}$  is a shift based on the number of requested eigenvalues  $r_k, r_{k-1}, \dots$  and number of matrix-vector products  $t_k, t_{k-1}, \dots$  for the previous eigenvalues problems. The shift  $\delta$  is computed as follows. First, we determine a *2-step shift value*  $\delta_2 \in \{-1, 1\}$  based on  $r_{k-1}, r_k$  and  $t_{k-1}, t_k$ . If  $r_k > r_{k-1}$  and  $t_k < t_{k-1}$  then the number of matrix-vector products in the EMEP (6.2.2) decreased as the number of requested eigenvalues was increased, suggesting we should shift  $r_k$  by  $\delta_2 = 1$ . By the same reasoning for the other three inequality cases, we define the *2-step shift value* as

$$\delta_2 = \text{sign}(r_k - r_{k-1}) \cdot \text{sign}(t_{k-1} - t_k), \quad (7.2.1)$$

where  $\text{sign}(0)$  is defined as 1. Next, if  $k \geq 4$  then we compute a linear interpolation of the past four requested eigenvalue numbers and matrix-vector products by solving

$$\min_{\alpha, \beta} \|y - \alpha e - \beta x\|, \quad (7.2.2)$$

where  $y$  is the vector of matrix-vector product values  $t_{k-3}, t_{k-2}, t_{k-1}, t_k$ ,  $x$  is the vector of the number of requested eigenvalues  $r_{k-3}, r_{k-2}, r_{k-1}, r_k$ , and  $e = [1; 1; 1; 1]$ . If the solution to (7.2.2) has  $\beta > 0$  then the past four eigenvalue problems suggest that  $t_i$  increases with  $r_i$ , and thus we should decrease  $r_k$ . Thus we have the *4-step shift value*

$$\delta_4 = -\text{sign}(\beta), \quad (7.2.3)$$

where  $\beta$  is determined by (7.2.2). If  $\delta_2 = \delta_4$ , then the 2-step (7.2.1) and 4-step equations (7.2.3) both suggest we should shift in the direction of  $\delta_2$ , and we select the shift  $\delta = 2\delta_2$ . If  $\delta_2 \neq \delta_4$  then we rely on the 2-step equation (7.2.1) and select the shift value  $\delta = \delta_2$ .

Finally, if  $r_k = r_{min}$  then we set  $\delta = 1$  and if  $r_k = r_{max}$  then we set  $\delta = -1$ . Altogether, these steps lead to the *adaptive parameter method* for the EMEP (6.2.2).

---

**Algorithm 8** Adaptive parameter method for the EMEP (6.2.2)

---

**Input:** Sequence of matrices  $\{A_k\}_{k=1}^{maxit}$  from the EMEP (6.2.2), Arnoldi decomposition (6.3.4) size  $m$  (default parameter  $m = 40$ ).

**Output:** EMEP (6.2.2) solution eigenpairs  $\{(\lambda_1^{(k)}, v_1^{(k)})\}_{k=1}^{maxit}$  and  $\{(\lambda_2^{(k)}, v_2^{(k)})\}_{k=1}^{maxit}$ .

- 1: *Initialize:*  $r_{min} = 2$ ,  $r_{max} = \min\{30, m - 5\}$ ,  $r_1 = r_{min}$ ,  $t_0 = -1$ ,  $k = 1$ .
- 2: **while**  $k \leq maxit$  **do**
- 3:     *Algorithm 7:* Perform IRAM with matrix  $A_k$ , number of requested algebraically largest eigenvalues  $r_k$ , and maximum Arnoldi decomposition size  $m$ . Return eigenpairs  $(\lambda_1^{(k)}, v_1^{(k)})$ ,  $(\lambda_2^{(k)}, v_2^{(k)})$  and number of matrix-vector products  $t_k$ .
- 4:     **if**  $r_k = r_{min}$  **then**
- 5:          $r_{k+1} = r_k + 1$
- 6:     **else if**  $r_k = r_{max}$  **then**
- 7:          $r_{k+1} = r_k - 1$
- 8:     **else if**  $k < 4$  **then**
- 9:         Compute 2-step shift value  $\delta_2$  from (7.2.1) and set  $\delta = \delta_2$
- 10:         $r_{k+1} = r_k + \delta$
- 11:     **else**
- 12:        Compute 2-step shift value  $\delta_2$  from (7.2.1) and 4-step shift value  $\delta_4$  from (7.2.3)
- 13:        **if**  $\delta_2 = \delta_4$  **then**
- 14:          Set  $\delta = 2\delta_2$
- 15:        **else**
- 16:          Set  $\delta = \delta_2$
- 17:        **end if**
- 18:         $r_{k+1} = \min\{\max\{r_k + \delta, r_{min}\}, r_{max}\}$
- 19:     **end if**
- 20:      $k = k + 1$
- 21: **end while**
- 22: *Return:*  $\{(\lambda_1^{(k)}, v_1^{(k)})\}_{k=1}^{maxit}$  and  $\{(\lambda_2^{(k)}, v_2^{(k)})\}_{k=1}^{maxit}$ .

---

Note that the only parameter in Algorithm 8 is the Arnoldi decomposition (6.3.4) size  $m$ , which determines the size of the basis  $Q_m \in \mathbb{C}^{n \times m}$  in the Arnoldi decomposition  $AQ_m = Q_m H_m + r_m e_m^*$ . The choice of  $m$  is a trade-off between computational efficiency and data storage constraints. We seek the smallest value  $m$  possible, since  $Q_m$  must be stored in random-access memory and each column of  $Q_m$  is the size of the desired signal  $\bar{x}$  in the phase retrieval problem (2.1.1). However, as we will see in Section 7.3,  $m$  must be sufficiently large for the shifted QR iteration (Algorithm 6) in the IRAM to handle the EMEP (6.2.2) efficiently. For now we will let  $m = 40$  to examine the behavior of Algorithm



8.

We close this section by showing that Algorithm 8 selects a sequence  $r_1, r_2, \dots, r_{\max it}$  which varies significantly and generally tracks the optimal choice of parameter  $r$  for each EMEP (6.2.2) iterate.

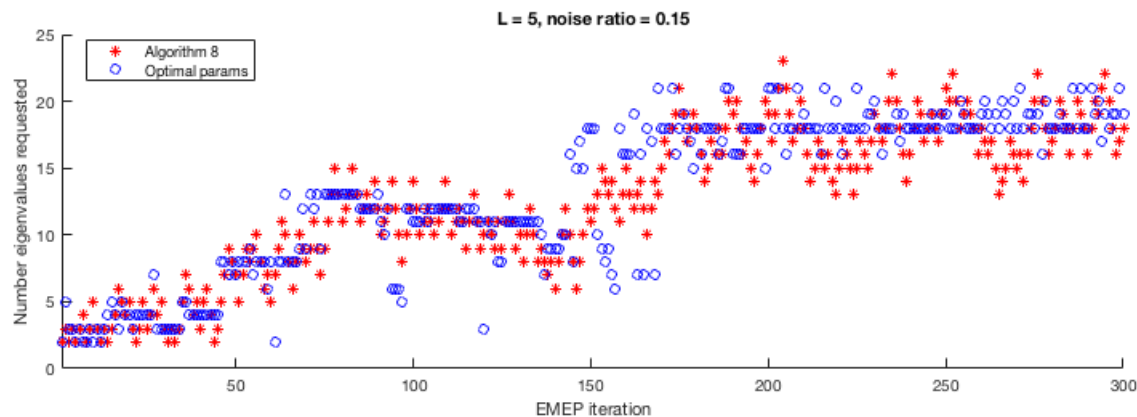


Figure 7.2: Plot comparing Algorithm 8 with the optimal choice of parameter  $r$  (number of requested eigenvalues in the IRAM). The EMEP (6.2.2) is from Figure 7.1 and Arnoldi decomposition (6.3.4) size is  $m = 40$ .

Figure 7.2 shows that the parameter values  $r_k$  chosen by Algorithm 8 effectively track the optimal parameters for the IRAM. For the majority of iterates, the value  $r_k$  is within two units from the optimal parameter value. As we will see in Section 7.3, these changes in  $r_k$  are related to the evolving spectrum of the EMEP (6.2.2).

### 7.3 Tracking the EMEP spectrum with the adaptive parameter method

In this section we explore the connection between the clustering of the algebraically largest eigenvalues in the EMEP (6.2.2) and the number of requested eigenvalues  $r_k$  as chosen by Algorithm 8. As the algebraically largest eigenvalues in the EMEP begin to cluster, we will see that the optimal value for  $r$  (corresponding to the minimum number of matrix-vector products for the IRAM) increases such that  $\lambda_{r+1}$  is not clustered with  $\lambda_1, \lambda_2, \dots, \lambda_r$ . Additionally, we will see that clustering of the algebraically largest eigenvalues can occur quickly in some EMEPs, causing the optimal value for  $r$  to increase quickly as well. The experiments in this section demonstrate that the value  $r_k$  chosen by Algorithm 8 properly tracks the optimal choice of  $r$  for EMEPs which evolve slowly or quickly. We also show

that the Arnoldi decomposition (6.3.4) size  $m = 40$  is an appropriate default parameter for Algorithm 8. We then close this section with a brief discussion of the possible cause of the optimal value for  $r$  to shift as the algebraically largest eigenvalues in the EMEP cluster.

This section focuses on two EMEPs (6.2.2) for which the sequence of parameters  $r_1, r_2, \dots, r_{maxit}$  chosen by Algorithm 8 varies greatly. Figure 7.3 depicts the results of Algorithm 8 for these EMEPs.

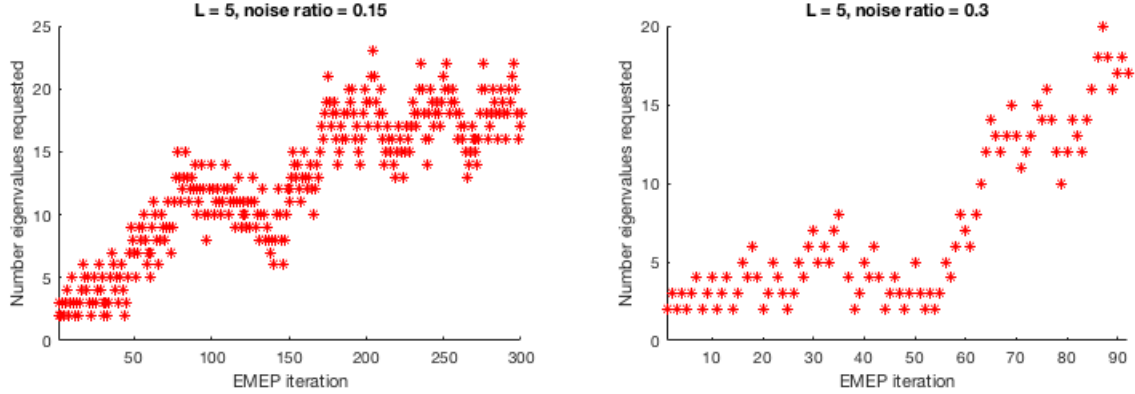


Figure 7.3: Number of requested eigenvalues  $r$  chosen by Algorithm 8 for two EMEPs (6.2.2). The EMEPs are from PLGD models with Gaussian noise (5.2.3), with noise ratio  $\epsilon_{\text{rel}} = 0.15, 0.30$ , oversampling rate  $L = 5$ , and original signal from Figure 2.2 resized to  $64 \times 64$  pixels.

Both EMEPs (6.2.2) in Figure 7.3 caused Algorithm 8 to increase the number of requested eigenvalues  $r$  from early to later iterates. Yet the rate at which  $r$  increased differs greatly for these two EMEPs. In particular, the experiment in Figure 7.3 with  $L = 5$  and  $\epsilon_{\text{rel}} = 0.30$  shows a rapid change in  $r_k$  around the iterate  $k = 60$ . On closer observation of this EMEP (Figure 7.3, right plot), we see that Algorithm 8 increased  $r_k$  by two units for several iterates  $55 \leq k \leq 70$ . These two-unit increases were caused by the 2-step shift value (7.2.1)  $\delta_2 = 1$  and 4-step shift value (7.2.3)  $\delta_4 = 1$  being equal, and suggest that the number of matrix-vector products was decreasing consistently as Algorithm 8 increased the number of requested eigenvalues  $r$ .

We will now examine the spectrum of the EMEPs (6.2.2) in Figure 7.3 to demonstrate that Algorithm 8 properly responds to two key tendencies of EMEPs. First, as the algebraically largest eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_{m-1}, \lambda_m$  of the EMEP begin to cluster, the number of requested eigenvalues  $r < m$  should be large enough such that the pair  $\{\lambda_r, \lambda_{r+1}\}$  will have sufficient separation, thus promoting convergence of the IRAM (Algorithm 7). And second, the rate at which this clustering occurs varies for different EMEPs, and thus  $r$  should increase at a corresponding rate. The number of matrix-vector products and eigenvalue difference  $\lambda_{r+1} - \lambda_r$  for the two EMEPs from Figure 7.3 are depicted in Figures 7.4

and 7.5, respectively. Note that these figures only depict about half of the EMEP iterates, so we may focus on the region where  $r_k$  varies greatest.

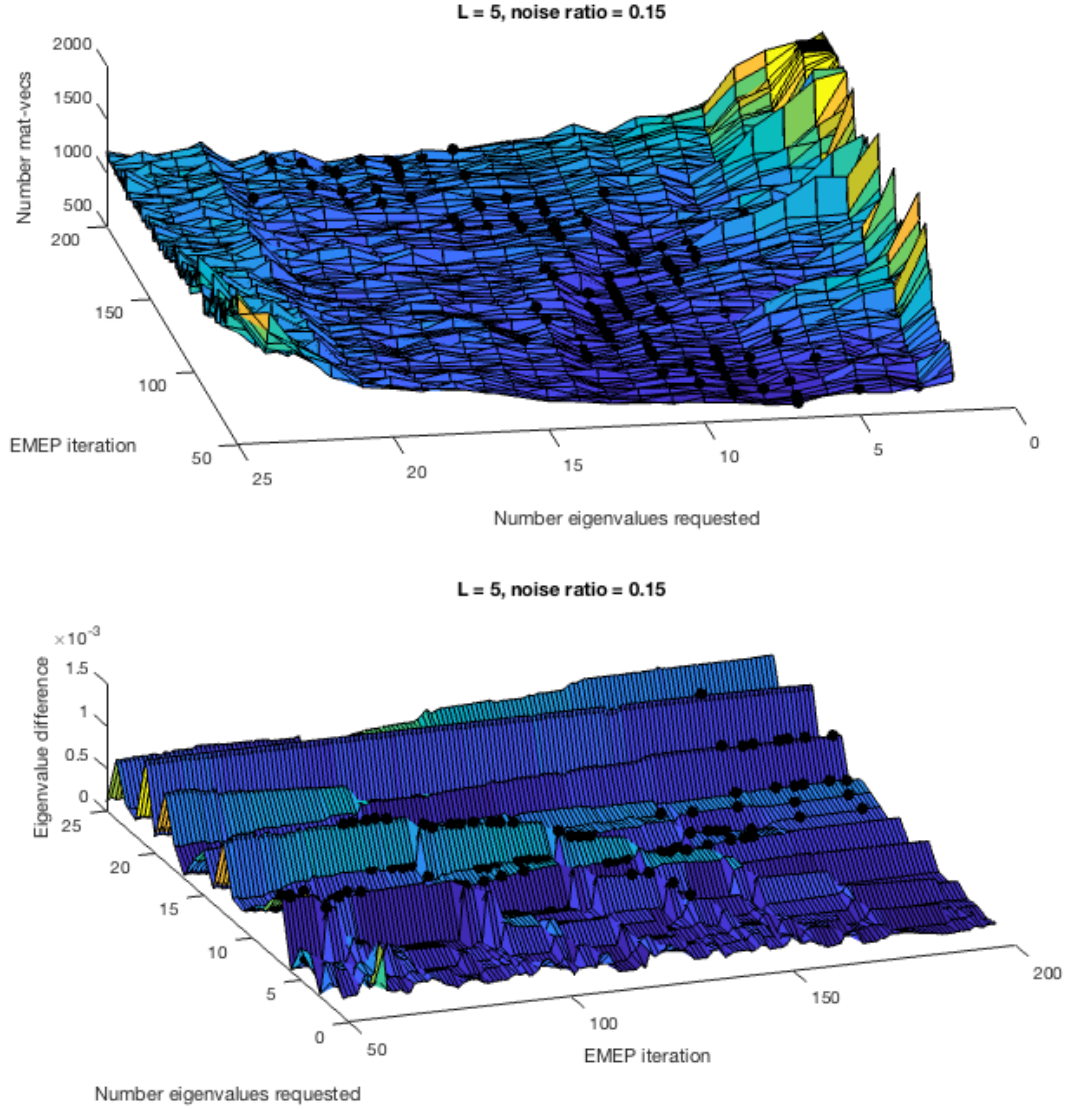


Figure 7.4: Behavior of Algorithm 8 for the experiment from Figure 7.3 with  $L = 5$ ,  $\epsilon_{\text{rel}} = 0.15$ . Top: Number of matrix-vector products (capped at 2,000 for better viewing) for each EMEP (6.2.2) iterate  $k$  and number of requested eigenvalues  $r$ . Bottom: eigenvalue differences  $\lambda_r - \lambda_{r+1}$  for each EMEP (6.2.2) iterate  $k$  and number of requested eigenvalues  $r$ . Black dots in both plots indicate the value  $r_k$  chosen by Algorithm 8.

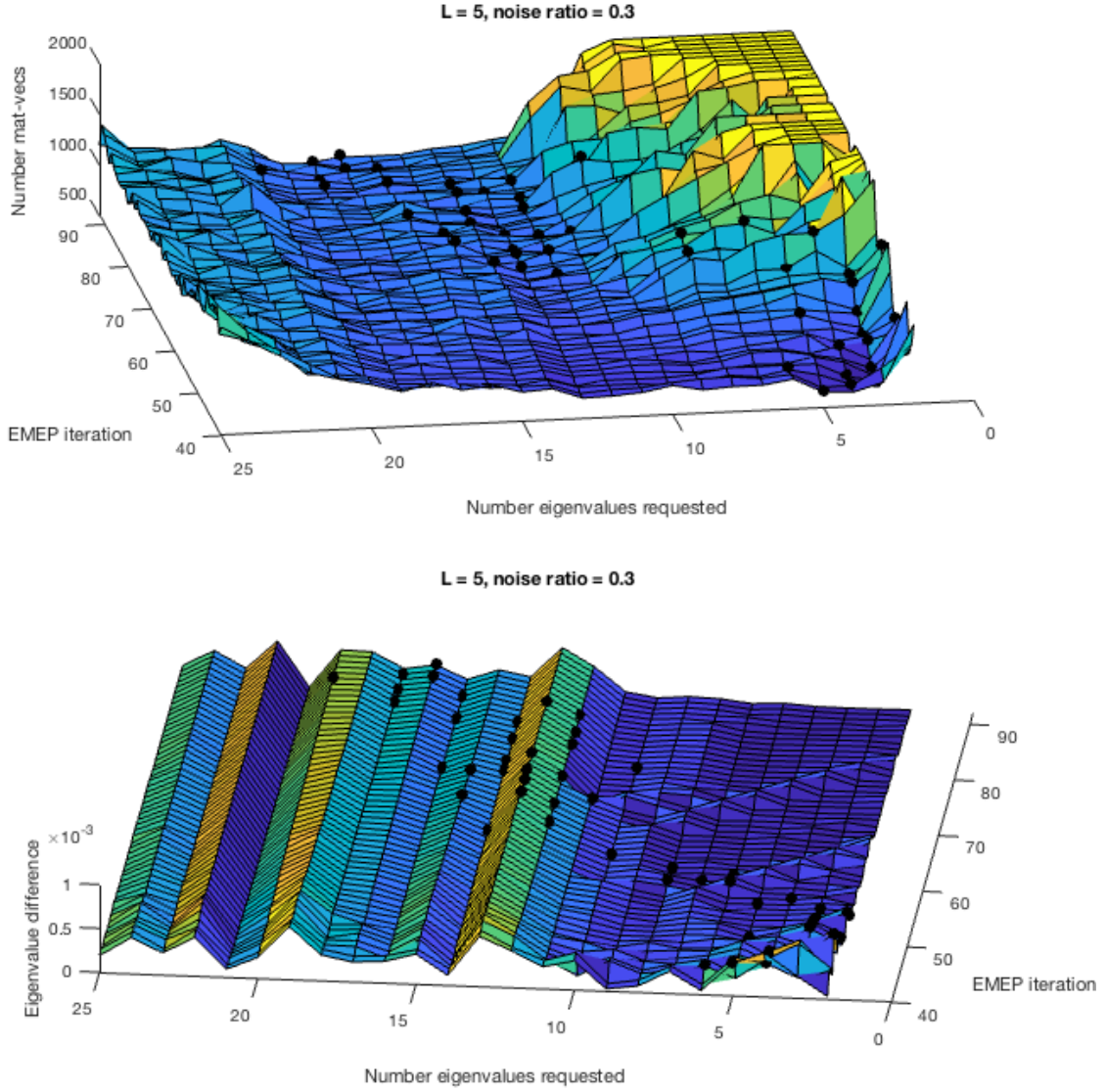


Figure 7.5: Behavior of Algorithm 8 for the experiment from Figure 7.3 with  $L = 5$ ,  $\epsilon_{\text{rel}} = 0.30$ . Top: Number of matrix-vector products (capped at 2,000 for better viewing) for each EMEP (6.2.2) iterate  $k$  and number of requested eigenvalues  $r$ . Bottom: eigenvalue differences  $\lambda_r - \lambda_{r+1}$  for each EMEP (6.2.2) iterate  $k$  and number of requested eigenvalues  $r$ . Black dots in both plots indicate the value  $r_k$  chosen by Algorithm 8.

In Figure 7.4, the top plot shows that  $r_k \approx 10$  offers the minimum number of matrix-vector products for EMEP iterates  $50 \leq k \leq 125$ . Correspondingly, the bottom plot in Figure 7.4 shows a “ridge” of eigenvalue differences  $\lambda_{r+1} - \lambda_r \approx 1 \times 10^{-3}$  around  $8 \leq r \leq 12$  for iterates  $50 \leq k \leq 125$ , and Algorithm 8 properly selects  $8 \leq r_k \leq 12$  for these iterates. Around iterate  $k = 150$ , this “ridge” of eigenvalue differences begins to flatten, creating a region of clustered algebraically largest eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_{10}$  for iterates  $k \geq 150$ . Likewise, the desired value  $r$  for the minimum number of matrix-vector products shifts to  $15 \leq r_k \leq 20$  for iterates  $k \geq 150$ . Yet the shifting of the desired value  $r$  in Figure 7.4 is gradual, and Algorithm 8 gradually increases  $r_k$  (as see in Figure 7.3, left plot) from  $r_k \approx 10$  at iterate  $k = 150$  to  $r_k \approx 18$  at iterate  $k = 200$ .

The bottom plot in Figure 7.5 also depicts less distinct “hills” of eigenvalue differences  $\lambda_{r+1} - \lambda_r \approx 0.5 \times 10^{-3}$  around  $2 \leq r \leq 8$  for iterates  $40 \leq k \leq 60$ . These “hills” of eigenvalue differences quickly flatten around iterates  $50 \leq k \leq 70$ . Likewise, the top plot in Figure 7.5 shows that the number of matrix-vector products dramatically increases for  $2 \leq r \leq 8$  around iterates  $50 \leq k \leq 70$ . In response to this change in the number of matrix-vector products, Algorithm 8 also increased  $r_k$  quickly around iterates  $55 \leq k \leq 70$ . For iterates  $k \geq 60$ , the top plot in Figure 7.5 shows that  $r = 12$  is the smallest parameter value necessary to have a minimal number of matrix-vector products. Correspondingly, the bottom plot in Figure 7.5 shows that for iterates  $k \geq 60$ , the pair  $\{\lambda_{12}, \lambda_{13}\}$  are the first pair of algebraically largest eigenvalues with greater separation than the preceding pairs. As desired, Algorithm 8 maintains  $r_k \geq 12$  for most of the iterates  $k \geq 70$ .

Next we show that the default parameter setting  $m = 40$  in Algorithm 8 is sufficiently large to minimize the number of matrix-vector products in the EMEPs (6.2.2) from Figure 7.3. To demonstrate that  $m = 40$  is appropriate, we will examine a more difficult eigenvalue problem (i.e., later iterate) from each of these EMEPs. Figure 7.6 depicts the number of matrix-vector products for various IRAM parameters  $r$  and  $m$  for an iterate from each EMEP.

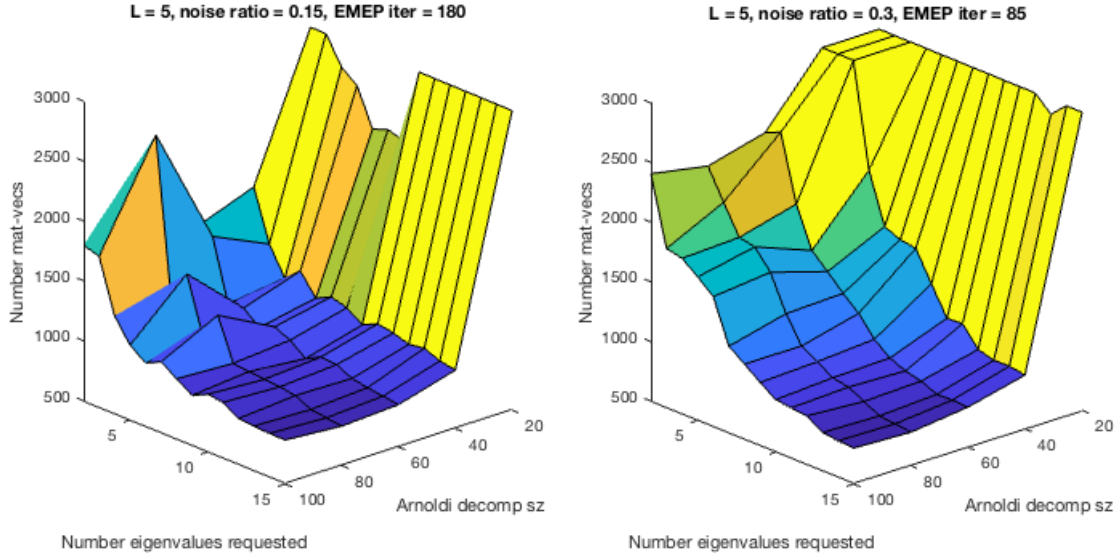


Figure 7.6: Number of matrix-vector products (capped at 3,000 for better viewing) for individual EMEP (6.2.2) iterates with varying number of requested eigenvalues  $r$  and Arnoldi decomposition (6.3.4) size  $m$ . Left: EMEP iterate 180 from Figure 7.4. Right: EMEP iterate 85 from Figure 7.5.

Figure 7.6 suggests that we should not select IRAM parameters below  $r = 9$  and  $m = 40$  for the EMEP iterate in the left plot, nor should we select parameters below  $r = 12$  and  $m = 40$  for the EMEP iterate in the right plot. Yet there is no significant benefit to selecting larger parameter values. Since the EMEP iterates in Figure 7.6 represent later, more difficult eigenvalue problems, this figure suggests that  $m = 40$  is sufficiently large for all iterates. As we saw in Figures 7.4 and 7.5, the desired number of requested eigenvalues  $r$  varies for earlier and later EMEP iterates. Thus Algorithm 8 has a fixed Arnoldi decomposition (6.3.4) size  $m = 40$  and adaptively adjusts  $r$  to minimize the number of matrix-vector products for each EMEP iterate.

We close this section by discussing a potential explanation for the IRAM (Algorithm 7) to perform poorly when the number of requested eigenvalues  $r$  is not sufficiently large for a given EMEP (6.2.2) iterate. As discussed in Section 6.3, the IRAM is based on the following two algorithms. Given a Hermitian matrix  $A \in \mathbb{C}^{n \times n}$ , the  $m$ -step Arnoldi iteration (Algorithm 5) generates a set of Ritz pairs  $(\theta_1, u_1), (\theta_2, u_2), \dots, (\theta_m, u_m)$  for  $A$  with respect to  $\mathcal{K}_m(A, q_1)$ . Next, the  $p$ -step shifted QR iteration (Algorithm 6) restarts the Arnoldi decomposition (6.3.4) by attempting to damp the unwanted part of the spectrum using the Ritz values  $\{\theta_{r+1}, \theta_{r+2}, \dots, \theta_m\}$  (where  $m = r + p$ ).

Assume we have an EMEP matrix iterate  $A_k$  with some number  $s$  of clustered algebraically largest eigenvalues,  $\lambda_1 \approx \lambda_2 \approx \dots \approx \lambda_s$ . And say we select the number of



requested eigenvalues  $r < s$ . When the IRAM builds an Arnoldi decomposition (6.3.4), the  $s$  largest Ritz values of  $A_k$  with respect to  $\mathcal{K}_m(A, q_1)$  may include values  $\theta_{r+1}, \theta_{r+2}, \dots, \theta_s$  which are close approximations to the desired eigenvalues. Thus when the shift values  $\mu_1 = \theta_{r+1}, \mu_2 = \theta_{r+2}, \dots, \mu_s = \theta_{r+s}, \dots, \mu_p = \theta_m$  are passed to the  $p$ -step shifted QR iteration, the implicit polynomial filter (6.3.12) will include values  $(\mu_1, \mu_2, \dots, \mu_s)$  which damp the desired part of the spectrum.

## 7.4 Numerical results for the adaptive parameter method

This section demonstrates the efficiency of the adaptive parameter method (Algorithm 8) for solving the EMEP (6.2.2). We begin by demonstrating that Algorithm 8 is more efficient than the original IRAM (Algorithm 7) parameter settings for a variety of EMEPs (6.2.2). Next, we show that Algorithm 8 is nearly optimal as a method for choosing the ideal number of requested eigenvalues  $r_k$  corresponding to the minimum number of matrix-vector products necessary for each EMEP (6.2.2) iteration. Finally, we demonstrate that the Arnoldi decomposition (6.3.4) size parameter  $m = 40$  for Algorithm 8 strikes a proper balance between increasing computational efficiency and minimizing data storage.

We begin this section by examining the performance of Algorithm 8 for various phase retrieval problems. Figure 7.7 depicts a set of experiments with randomly generated signals.

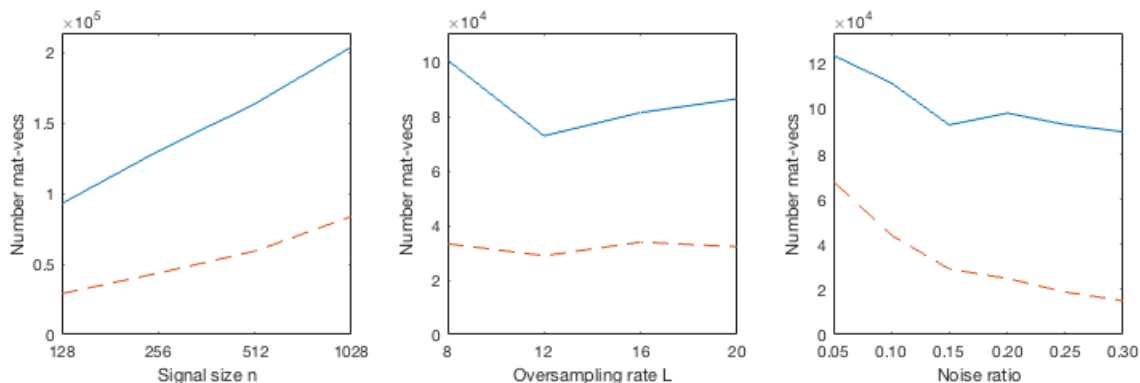


Figure 7.7: Performance results for Algorithm 8 (dashed line) and the original IRAM parameters  $m = 20$ ,  $r = 2$  (solid line). All problems are PLGD models with Gaussian noise (5.2.3) and signals are complex with standard Gaussian distribution (1.2.17). Each result is the mean of 10 experiments. Left: Varying signal size  $n$ , with fixed noise ratio  $\epsilon_{\text{rel}} = 0.15$  and oversampling scaled logarithmically with  $n$  (i.e.,  $L = 10, 12, 12, 14$ ) as indicated in Theorem 2.3.2. Middle: Varying oversampling rate  $L$ , with fixed signal size  $n = 128$  and noise ratio  $\epsilon_{\text{rel}} = 0.15$ . Right: Varying noise ratio  $\epsilon_{\text{rel}}$ , with fixed signal size  $n = 128$  and oversampling rate  $L = 10$ .



Figure 7.7 demonstrates that Algorithm 8 requires fewer matrix-vector products than the original IRAM parameters for the EMEP (6.2.2) for a wide range of phase retrieval problems. The left and middle plots in Figure 7.7 suggest that Algorithm 8 requires about 60% fewer matrix-vector products regardless of signal size  $n$  or oversampling rate  $L$ . Additionally, the right plot in Figure 7.7 suggests Algorithm 8 may reduce matrix-vector products by 80% or more for problems with significant noise.

Next, we demonstrate that Algorithm 8 is nearly optimal in the sense of choosing the number of requested eigenvalues  $r_k$  which minimizes the number of matrix-vector products for each EMEP (6.2.2) iterate  $k$ . Table 7.1 indicates the number of matrix-vector products for solving the six EMEPs with Algorithm 8, with the original IRAM parameters, and with the minimum possible number of matrix-vector products if each value  $r_k$  was chosen such that  $2 \leq r_k \leq 30$  and  $r_k$  corresponds to the minimum number of matrix-vector products for the IRAM to evaluate EMEP iterate  $k$ .

$L$	$\epsilon_{\text{rel}}$	EMEP its	Original $r = 2, m = 20$	Optimal $2 \leq r \leq 30$ $m = 40$		Algorithm 8 $m = 40$	
5	0.05	300	406,308	179,807	56%	198,070	51%
5	0.15	300	1,099,045	242,003	78%	258,385	76%
5	0.30	92	444,697	58,780	87%	69,510	84%
10	0.05	153	80,453	61,948	23%	68,709	15%
10	0.15	108	88,317	51,311	42%	57,231	35%
10	0.30	54	72,486	23,217	68%	25,809	64%

Table 7.1: Total number of matrix-vector products and percent decrease from the original IRAM parameters for various EMEPs (6.2.2). Parameter  $r$  is the number of requested eigenvalues in the IRAM (Algorithm 7) and  $m$  is the Arnoldi decomposition size (6.3.4). EMEPs come from PLGD models with Gaussian noise (5.2.3) with original signal from Figure 2.2 resized to  $64 \times 64$  pixels.

Table 7.1 demonstrates that Algorithm 8 decreases the number of matrix-vector products of each EMEP (6.2.2) from the original IRAM parameters by a percentage comparable to that of the optimal choice for parameter  $r$ . Notably, Algorithm 8 is particularly effective at decreasing the number of matrix-vector products when there is a large relative difference between the number of matrix-vector products for the original IRAM parameters and the optimal parameters. To further explore this performance behavior, Figure 7.8 depicts the two EMEPs from Table 7.1 with the largest and smallest relative difference in matrix-vector products (those with  $L = 5$ ,  $\epsilon_{\text{rel}} = 0.15$ , and  $L = 10$ ,  $\epsilon_{\text{rel}} = 0.05$ , respectively).

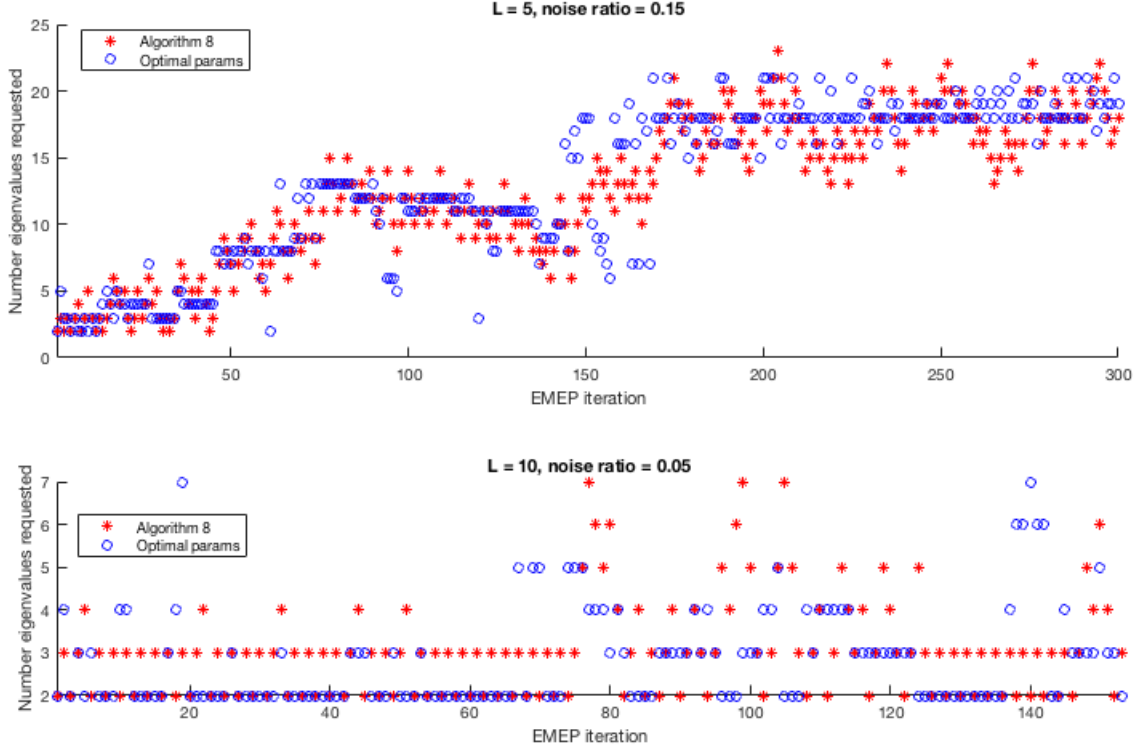


Figure 7.8: Number of requested eigenvalues  $r$  in the IRAM (Algorithm 7) for two EMEPs (6.2.2) from Table 7.1.

In both EMEPs (6.2.2) depicted in Figure 7.8, the number of requested eigenvalues  $r$  chosen by Algorithm 8 is usually within 1-3 units from the optimal parameter value. The bottom plot in Figure 7.8 suggests that Algorithm 8 required relatively more matrix-vector products than the optimal value of  $r$  because Algorithm 8 always changes the value of  $r$  by one or two units, thus shifting away from the optimal value  $r = 2$  for many EMEP iterates. As we saw in Figure 7.5, the optimal parameter  $r$  may shift rapidly for some EMEPs, and thus Algorithm 8 always changes  $r$  by one or two units to continue gathering performance information about the EMEP.

We now provide further justification for selecting  $m = 40$  as the default Arnoldi decomposition (6.3.4) size for Algorithm 8. Table 7.2 depicts the total number of matrix-vector products required to solve the EMEPs (6.2.2) from Table 7.1 with Algorithm 8 with various parameter values  $m$ .

n = 4,096			Original $r = 2, m = 20$	Adaptive parameter method (Algorithm 8)				
$L$	$\epsilon_{\text{rel}}$	EMEP its		$m = 20$	$m = 40$	$m = 60$	$m = 80$	$m = 100$
5	0.05	300	406,308	358,195	198,070	189,401	192,042	201,270
5	0.15	300	1,099,045	806,412	258,385	224,048	214,118	215,392
5	0.30	92	444,697	175,669	69,510	56,193	55,146	54,987
10	0.05	153	80,453	77,768	68,709	64,300	68,602	73,754
10	0.15	108	88,317	65,833	57,231	53,261	54,388	55,308
10	0.30	54	72,486	28,799	25,809	24,699	25,113	25,491

Table 7.2: Total number of matrix-vector products for various EMEPs (6.2.2). Parameter  $r$  is the number of requested eigenvalues in the IRAM (Algorithm 7) and  $m$  is the Arnoldi decomposition (6.3.4) size. EMEPs come from PLGD models with Gaussian noise (5.2.3) with original signal from Figure 2.2 resized to  $64 \times 64$  pixels.

Table 7.2 demonstrates that Algorithm 8 reduces the number of matrix-vector products from those of the original IRAM parameters for all experiments considered. Yet this cost reduction varies significantly depending on the choice of Arnoldi decomposition (6.3.4) size  $m$ . We seek a default setting for the parameter  $m$  which is sufficiently large to yield the benefits of Algorithm 8, yet sufficiently small as not to burden random-access memory constraints. To select an appropriate default value for  $m$ , we examine the two experiments from Table 7.2 with  $\epsilon_{\text{rel}} = 0.15, 0.30$  and  $L = 5$  which have the greatest original number of matrix-vector products, along with the greatest total decrease in cost when using Algorithm 8 with a sufficiently large parameter  $m$ . Figure 7.9 singles out these two experiments, depicting the number of matrix-vector products for each EMEP (6.2.2) iteration.

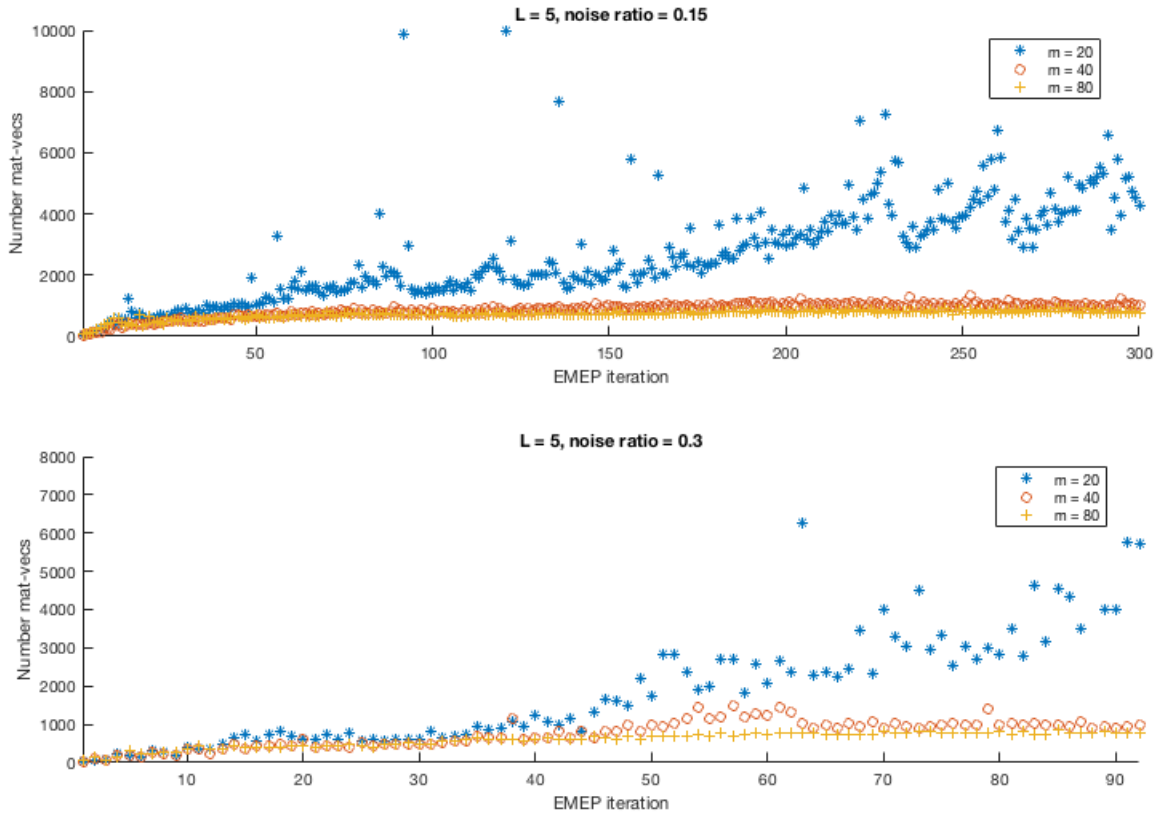


Figure 7.9: Number of matrix-vector products for each EMEP (6.2.2) iteration from two experiments in Figure 7.9 with various Arnoldi decomposition size  $m = 20, 40, 80$ .

Figure 7.9 demonstrates that the Arnoldi decomposition size of  $m = 20$  is not sufficiently large to allow Algorithm 8 to decrease the number of matrix-vector products. The dramatic matrix-vector product spikes for  $m = 20$  in Figure 7.9 resemble those first seen in Figure 6.1 when solving the EMEP (6.2.2) with the original IRAM parameters  $m = 20, r = 2$ . Yet when the Arnoldi decomposition size is increased to  $m = 40$ , these cost spikes effectively disappear. The change in number of matrix-vector products between  $m = 40$  and  $m = 80$  is minimal for each EMEP iterate. Thus the default parameter of  $m = 40$  for Algorithm 8 strikes the proper balance between efficiency and memory constraints.

As we have seen in this chapter, Algorithm 8 is an efficient strategy for using the IRAM (Algorithm 7) to handle the EMEP (6.2.2). We saw that changes in the number of requested eigenvalues, as chosen by Algorithm 8, correspond to clustering of the algebraically largest eigenvalues in the EMEP. Thus Algorithm 8 appears to promote convergence of the IRAM

by selecting an appropriate number of requested eigenvalues  $r$  for which the pair  $\lambda_r$  and  $\lambda_{r+1}$  will have sufficient separation. Algorithm 8 was shown to reduce the number of matrix-vector products required to handle EMEPs by 50 – 90% when the underlying phase retrieval models have minimal oversampling. Our experiments indicated that Algorithm 8 is nearly optimal as a method for selecting the number of requested eigenvalues which corresponds to the minimum possible number of matrix-vector products required for each EMEP iterate.

## Chapter 8

# Conclusion

Conclusion.

## Appendix A

# Proof of PhaseLift- $l_1$ gauge duality

In this appendix we show that the following pair of PhaseLift- $l_1$  models discussed in Section 2.4 are a primal and gauge dual pair

$$\begin{array}{ll}
 \min & \|\mathcal{A}(X) - b\|_1 \\
 \text{(PLP-}l_1\text{)} \quad \text{s.t.} & X \succeq 0
 \end{array}
 \qquad
 \begin{array}{ll}
 \min_y & \|y\|_\infty \\
 \text{(PLGD-}l_1\text{)} \quad \text{s.t.} & -\mathcal{A}^*y \succeq 0 \\
 & \langle b, y \rangle \geq 1.
 \end{array}
 \tag{A.0.1}$$

To determine the gauge dual of PLP- $l_1$ , first note that the objective function  $f(X) = \|\mathcal{A}(X) - b\|_1$  is not a gauge function, as  $f$  is not positively homogeneous. Using the substitution  $z = b - \mathcal{A}(X)$  and extending the linear operator  $\mathcal{A}$ , we may restate PLP- $l_1$  as

$$\begin{array}{ll}
 \min_{X,z} & \kappa(X, z) := \|z\|_1 + \delta_{\succeq 0}(X) \\
 \text{s.t.} & \overline{\mathcal{A}}(X, z) := \mathcal{A}(X) + z = b,
 \end{array}
 \tag{A.0.2}$$

which is now in the form of (3.2.21). The gauge functions  $\kappa_1(z) = \|z\|_1$  and  $\kappa_2(X) = \delta_{\succeq 0}(X)$  have polars  $\kappa_1^\circ(w) = \|w\|_\infty$  and  $\kappa_2^\circ(V) = \delta_{\succeq 0}(-V)$ . Thus, by Proposition 3.3.5,  $\kappa$  has the polar

$$\kappa^\circ(V, w) = \max\{\|w\|_\infty, \delta_{\succeq 0}(-V)\} = \|w\|_\infty + \delta_{\succeq 0}(-V).
 \tag{A.0.3}$$

Additionally, the adjoint of  $\overline{\mathcal{A}}$  is  $\overline{\mathcal{A}}^*y = (\mathcal{A}^*y, y)$ , giving

$$\kappa^\circ(\overline{\mathcal{A}}^*y) = \|y\|_\infty + \delta_{\succeq 0}(-\mathcal{A}^*y).
 \tag{A.0.4}$$

Passing  $\delta_{\succeq 0}(-\mathcal{A}^*y)$  into the constraint set, we see that PLGD- $l_1$  is the gauge dual of PLP- $l_1$ .

## Appendix B

# Further justification of residuals established in Chapter 5

1. The previous section demonstrated that the primal difference (5.2.7e) and dual variable difference (5.2.7i) effectively identify the point at which Algorithm 3 stagnates for PLGD models with Gaussian noise (5.2.3). We close this chapter with a justification for ignoring the remaining residuals (5.2.7). As we will see, each of these residuals is either unreliable or effectively a duplicate of another (computationally cheaper) residual.

One residual that may be eliminated is the primal relative error (5.2.7d). The original termination conditions for Algorithm 3 include the feasibility requirement (4.2.13)

$$\|\mathcal{A}(xx^*) - b\| \leq \epsilon + \text{tol}_{\text{feas}}(1 + \|b\|),$$

which is equivalent to the primal relative error (5.2.7d) using the relation  $\epsilon_{\text{rel}} = \epsilon/\|b\|$

$$\frac{\|\mathcal{A}(xx^*) - b\|}{\|b\|} \leq \epsilon_{\text{rel}} + \text{tol}_{\text{feas}} \left( \frac{1 + \|b\|}{\|b\|} \right). \quad (\text{B.0.1})$$

As indicated in Table 5.4, Algorithm 3 typically requires only a few iterations before a primal feasible signal  $x$  is found. Yet in certain cases this algorithm may never identify a feasible signal. Figure B.1 demonstrates plots the primal relative error (5.2.7d) for the particular model from Figure 5.3 which never attained primal feasibility.



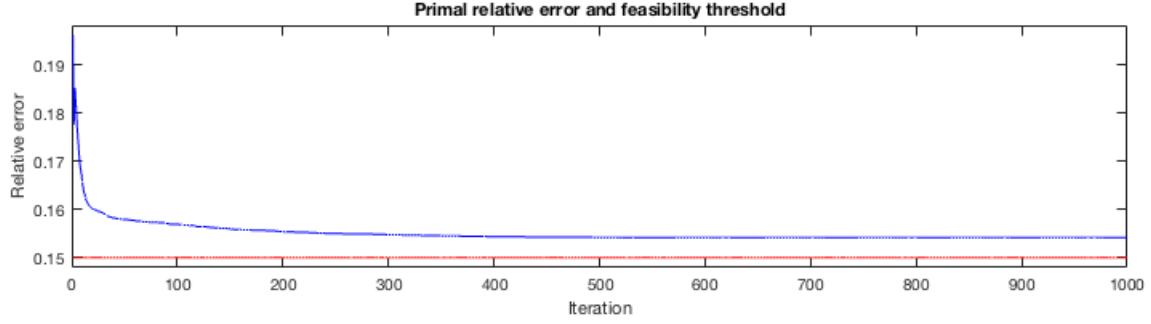


Figure B.1: Primal relative error (5.2.7d) (blue) and noise ratio threshold (red) for noise ratio  $\epsilon_{\text{rel}} = 0.15$  and oversampling rate  $L = 5$ . Iterate 1,000 remains infeasible with a primal relative error of 0.1541.

The experiment in Figure B.1 demonstrates a model where both the exact tolerance  $\epsilon_{\text{rel}} = 0.15$  and the relaxed tolerance  $\epsilon_{\text{rel}} + \text{tol}_{\text{feas}} \left( \frac{1+\|b\|}{\|b\|} \right) = 0.1502$  as suggested in (4.2.13) are unattainable by Algorithm 3. Nevertheless, this algorithm makes significant progress in recovering an approximate signal. Notably, wflow applied to the same model recovers a signal with a primal relative error (5.2.7d) of 0.3165. Thus we ignore the primal feasibility condition (4.2.13).

We may also disregard the duality gap (5.2.7f) as a termination condition. The experiments in Table 5.3 demonstrate that the duality gap tends to stagnate at different values for differing noise level and oversampling rate. Table B.1 indicates the final duality gap value for each experiment from Figure 5.3, further demonstrating this residual is not a reliable indicator of stagnation.

	$L = 5$	$L = 10$
$\epsilon_{\text{rel}} = 0.05$	8.24	8.10
$\epsilon_{\text{rel}} = 0.15$	39.23	29.41
$\epsilon_{\text{rel}} = 0.30$	38.93	59.73

Table B.1: Final values of duality gap (5.2.7f) after 1,000 iterations of Algorithm 3 with indicated noise ratios and oversampling rates.

The duality gap (5.2.7f) is not reliable because this measurement is dependent on the accuracy of both the dual objective value  $\lambda_1$  and the approximate signal norm  $\|xx^*\|_1$  as compared to their respective optimal values. As established in Section 5.3, Algorithm 3 cannot achieve optimality for most noisy phase retrieval models. Thus the complementarity condition  $\|xx^*\|_1 \cdot \lambda_1 = 1$  (Corollary 3.4.1, d) cannot be achieved. As a result, the duality gap (5.2.7f) stagnates unpredictably for varying

oversampling rates and noise levels and is not used as a termination condition for Algorithm 3.

Next we demonstrate that the duality gap difference (5.2.7g) is a redundant measurement. Denote the primal objective value as  $p = \|xx^*\|_1$  and its update with a hat. As Algorithm 3 stagnates,  $\hat{p}/p$  approaches 1. Additionally, note that typical optimal signals have fairly large norms (i.e.,  $\|\bar{x}\bar{x}^*\|_F = \mathcal{O}(10^3)$  or larger). If we assume  $\hat{p}/p \approx 1$  and  $\lambda_1 - 1/p \approx \lambda_1$  for later iterates then we have

$$\begin{aligned} \frac{|\gamma - \hat{\gamma}|}{\gamma} &= \frac{|(p\lambda_1 - 1) - (\hat{p}\hat{\lambda}_1 - 1)|}{p\lambda_1 - 1} \\ &= \frac{\left| \lambda_1 - \frac{\hat{p}}{p}\hat{\lambda}_1 \right|}{\lambda_1 - \frac{1}{p}} \\ &\approx \frac{|\lambda_1 - \hat{\lambda}_1|}{\lambda_1}. \end{aligned} \tag{B.0.2}$$

Thus the duality gap difference (5.2.7g) behaves similarly to the dual objective difference (5.2.7h) and may be disregarded.

Finally, the dual matrix difference (5.2.7j) is also a redundant measurement which may be disregarded. Note that the norm difference of dual matrices is bounded above by the dual variable difference (5.2.7i), since

$$\|A - \hat{A}\|_F = \|\mathcal{A}^*(y - \hat{y})\|_F \leq \|\mathcal{A}^*\| \|y - \hat{y}\|,$$

where we have the induced norm

$$\|\mathcal{A}^*\| = \sup_{\|w\|=1} \|\mathcal{A}^*w\|_F.$$

And computationally, the dual variable difference (5.2.7i) is computed with a vector norm, where as the dual matrix difference (5.2.7j) requires an additional eigenvalue computation (in this case the largest magnitude eigenvalue). Thus the dual matrix difference (5.2.7j) can be ignored due to its excessive computational cost and close relationship to the dual variable difference (5.2.7i).

Finally, we may also eliminate the dual objective difference (5.2.7h) as a candidate residual for termination. Figure B.2 depicts the behavior of this residual for the models in Figure 5.3. As with Figure 5.4, the vertical axis indicates specific tolerances and the horizontal axis indicates the first iterate at which Algorithm 3 would satisfy this tolerance.

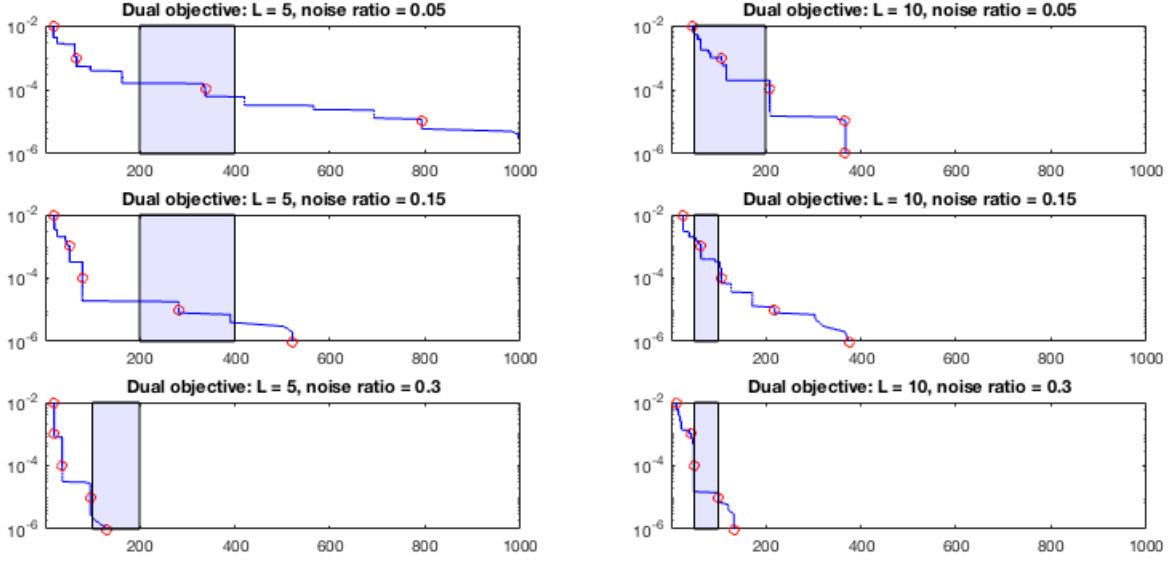


Figure B.2: Plots of dual objective difference values (5.2.7h) against the iterate at which Algorithm 3 first satisfies this tolerance for the models discussed in Figure 5.3. Red circles are placed at tolerances  $10^{-n}$ . The blue rectangles indicate the proposed intervals of termination from Table 5.5.

Figure B.2 demonstrates that the dual objective difference (5.2.7h) behaves erratically, decreasing too slowly for low-noise models and too quickly for high-noise models. To simplify our discussion, label the dual objective tolerance as  $\text{tol}_{\text{DO}}$ . If we set  $\text{tol}_{\text{DO}} = 10^{-5}$ , then Algorithm 3 will terminate far too late for the top-left model in Figure B.2. However, if we set  $\text{tol}_{\text{DO}} = 10^{-4}$  then the algorithm may terminate far too early for the middle-left and bottom-left models. Likewise, the models at right in Figure B.2 do not have a consistent tolerance value  $\text{tol}_{\text{DO}}$  which reliably selects for termination within the desired intervals. Thus the tolerance  $\text{tol}_{\text{DO}}$  is not a reliable indicator of stagnation of Algorithm 3.

# Bibliography

- [ABD<sup>+</sup>17] A. Y. Aravkin, J. V. Burke, D. Drusvyatskiy, M. P. Friedlander, and K. MacPhee, *Foundations of gauge and perspective duality*, arXiv preprint arXiv:1702.08649 (2017).
- [BB86] R. N. Bracewell and R. N. Bracewell, *The fourier transform and its applications*, third ed., vol. 31999, McGraw-Hill New York, 1986.
- [BB88] J. Barzilai and J. M. Borwein, *Two-point step size gradient methods*, IMA Journal of Numerical Analysis **8** (1988), no. 1, 141–148.
- [BCG11] S. R. Becker, E. J. Candès, and M. C. Grant, *Templates for convex cone problems with applications to sparse signal recovery*, Mathematical programming computation **3** (2011), no. 3, 165.
- [BDP<sup>+</sup>07] O. Bunk, A. Diaz, F. Pfeiffer, C. David, B. Schmitt, D. K. Satapathy, and J. F. van der Veen, *Diffraction imaging for periodic samples: retrieving one-dimensional concentration profiles across microfluidic channels*, Acta Crystallographica Section A: Foundations of Crystallography **63** (2007), no. 4, 306–314.
- [Ber16] D. P. Bertsekas, *Nonlinear programming*, third ed., Athena Scientific, 2016.
- [BR16] S. Bahmani and J. Romberg, *Phase retrieval meets statistical learning theory: A flexible convex relaxation*, arXiv preprint arXiv:1610.04210 (2016).
- [BTN01] A. Ben-Tal and A. Nemirovski, *Lectures on modern convex optimization: analysis, algorithms, and engineering applications*, vol. 2, Siam, 2001.
- [BV04] S. Boyd and L. Vandenberghe, *Convex optimization*, Cambridge university press, 2004.
- [CDS01] S. S. Chen, D. L. Donoho, and M. A. Saunders, *Atomic decomposition by basis pursuit*, SIAM review **43** (2001), no. 1, 129–159.

- [CESV13] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski, *Phase retrieval via matrix completion*, SIAM J. Imaging Sciences **6** (2013), no. 1, 199–225.
- [CG90] P. Comon and G. H. Golub, *Tracking a few extreme singular values and vectors in signal processing*, Proceedings of the IEEE **78** (1990), no. 8, 1327–1343.
- [CL14] E. J. Candès and X. Li, *Solving quadratic equations via phaselift when there are about as many equations as unknowns*, Foundations of Computational Mathematics **14** (2014), no. 5, 1017–1026.
- [CLM<sup>+</sup>16] T. T. Cai, X. Li, Z. Ma, et al., *Optimal rates of convergence for noisy sparse phase retrieval via thresholded wirtinger flow*, The Annals of Statistics **44** (2016), no. 5, 2221–2251.
- [CLS15] E. J. Candès, X. Li, and M. Soltanolkotabi, *Phase retrieval via wirtinger flow: Theory and algorithms*, IEEE Trans. Information Theory **61** (2015), no. 4, 1985–2007.
- [CMP10] A. Chai, M. Moscoso, and G. Papanicolaou, *Array imaging using intensity-only measurements*, Inverse Problems **27** (2010), no. 1, 015005.
- [CRT06] E. J. Candes, J. K. Romberg, and T. Tao, *Stable signal recovery from incomplete and inaccurate measurements*, Communications on pure and applied mathematics **59** (2006), no. 8, 1207–1223.
- [CSV13] E. J. Candes, T. Strohmer, and V. Voroninski, *Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming*, Communications on Pure and Applied Mathematics **66** (2013), no. 8, 1241–1274.
- [DM08] X. G. Doukopoulos and G. V. Moustakides, *Fast and stable subspace tracking*, IEEE Transactions on Signal Processing **56** (2008), no. 4, 1452–1465.
- [DMM<sup>+</sup>11] H. Duadi, O. Margalit, V. Mico, J. A. Rodrigo, T. Alieva, J. Garcia, and Z. Zalevsky, *Digital holography and phase retrieval*, Holography, Research and Technologies, InTech, 2011.
- [ELB17] V. Elser, T.-Y. Lan, and T. Bendory, *Benchmark problems for phase retrieval*, arXiv preprint arXiv:1706.00399 (2017).
- [EM12] Y. C. Eldar and S. Mendelson, *Phase retrieval: Stability and recovery guarantees*, CoRR **abs/1211.0872** (2012), 1211.0872.
- [FD87] C. Fienup and J. Dainty, *Phase retrieval and image reconstruction for astronomy*, Image Recovery: Theory and Application **231** (1987), 275.

- [Fie82] J. R. Fienup, *Phase retrieval algorithms: a comparison*, Appl. Opt. **21** (1982), no. 15, 2758–2769.
- [FM16] M. P. Friedlander and I. Macedo, *Low-rank spectral optimization via gauge duality*, SIAM J. Scientific Computing **38** (2016), no. 3.
- [FMP14] M. P. Friedlander, I. Macedo, and T. K. Pong, *Gauge optimization and duality*, SIAM Journal on Optimization **24** (2014), no. 4, 1999–2022.
- [Fre87] R. M. Freund, *Dual gauge programs, with applications to quadratic programming and the minimum-norm problem*, Math. Program. **38** (1987), no. 1, 47–67.
- [Gri91] R. D. Grigorieff, *A note on von neumann’s trace inequality*, (1991).
- [GS72] R. W. Gerchberg and W. O. Saxton, *A practical algorithm for the determination of phase from image and diffraction plane pictures*, Optik **35** (1972), 237–250.
- [GS18a] T. Goldstein and C. Studer, *Phasemax: Convex phase retrieval via basis pursuit*, IEEE Transactions on Information Theory (2018).
- [GS18b] M. Guizar-Sicairos, *Phase retrieval for x-ray coherent lensless imaging*, 2018.
- [GVL12] G. H. Golub and C. F. Van Loan, *Matrix computations*, 4 ed., JHU Press, 2012.
- [Har93] R. W. Harrison, *Phase problem in crystallography*, JOSA a **10** (1993), no. 5, 1046–1055.
- [JEH15] K. Jaganathan, Y. C. Eldar, and B. Hassibi, *Phase retrieval: An overview of recent developments*, CoRR **abs/1510.07713** (2015), 1–24, 1510.07713.
- [JSL17] X. Jiang, H.-C. So, and X. Liu, *Robust phase retrieval via admm with outliers*, arXiv preprint arXiv:1702.06157 (2017).
- [Kat17] V. Katkovnik, *Phase retrieval from noisy data based on sparse approximation of object phase and amplitude*, arXiv preprint arXiv:1709.01071 (2017).
- [LS84] A. Levi and H. Stark, *Image restoration by the method of generalized projections with application to restoration from magnitude*, IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP ’84, San Diego, California, USA, March 19-21, 1984, 1984, pp. 88–91.
- [LSY98] R. B. Lehoucq, D. C. Sorensen, and C. Yang, *Arpack users’ guide: solution of large-scale eigenvalue problems with implicitly restarted arnoldi methods*, vol. 6, Siam, 1998.

- [MCKS99] J. Miao, P. Charalambous, J. Kirz, and D. Sayre, *Extending the methodology of x-ray crystallography to allow imaging of micrometre-sized non-crystalline specimens*, Nature **400** (1999), no. 6742, 342.
- [Mil90] R. P. Millane, *Phase retrieval in crystallography and optics*, JOSA A **7** (1990), no. 3, 394–411.
- [MISE08] J. Miao, T. Ishikawa, Q. Shen, and T. Earnest, *Extending x-ray crystallography to allow the imaging of noncrystalline materials, cells, and single protein complexes*, Annu. Rev. Phys. Chem. **59** (2008), 387–410.
- [MWL<sup>+</sup>12] A. V. Martin, F. Wang, N. Loh, T. Ekeberg, F. R. Maia, M. Hantke, G. van der Schot, C. Y. Hampton, R. G. Sierra, A. Aquila, et al., *Noise-robust coherent diffractive imaging with a single diffraction pattern*, Optics Express **20** (2012), no. 15, 16650–16661.
- [Nat95] B. K. Natarajan, *Sparse approximate solutions to linear systems*, SIAM journal on computing **24** (1995), no. 2, 227–234.
- [NS12] T. Ngo and Y. Saad, *Scaled gradients on grassmann manifolds for matrix completion*, Advances in Neural Information Processing Systems, 2012, pp. 1412–1420.
- [NW06] J. Nocedal and S. J. Wright, *Numerical optimization 2nd*, Springer, 2006.
- [PB<sup>+</sup>14] N. Parikh, S. Boyd, et al., *Proximal algorithms*, Foundations and Trends® in Optimization **1** (2014), no. 3, 127–239.
- [RFP10] B. Recht, M. Fazel, and P. A. Parrilo, *Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization*, SIAM review **52** (2010), no. 3, 471–501.
- [Roc70] R. T. Rockafellar, *Convex analysis*, Princeton university press, 1970.
- [RS80] M. Reed and B. Simon, *Functional analysis, vol. i*, Academic Press, San Diego, 1980.
- [RXC<sup>+</sup>13] J. A. Rodriguez, R. Xu, C.-C. Chen, Y. Zou, and J. Miao, *Oversampling smoothness: an effective algorithm for phase retrieval of noisy diffraction intensities*, Journal of applied crystallography **46** (2013), no. 2, 312–318.
- [Saa11] Y. Saad, *Numerical methods for large eigenvalue problems: revised edition*, vol. 66, Siam, 2011.

- [SBE14] Y. Shechtman, A. Beck, and Y. C. Eldar, *Gespar: Efficient phase retrieval of sparse signals*, IEEE transactions on signal processing **62** (2014), no. 4, 928–938.
- [SBFM09] M. Schmidt, E. Berg, M. Friedlander, and K. Murphy, *Optimizing costly functions with simple constraints: A limited-memory projected quasi-newton algorithm*, Artificial Intelligence and Statistics, 2009, pp. 456–463.
- [Sch05] M. Schmidt, *minfunc: unconstrained differentiable multivariate optimization in matlab*, <http://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>, 2005.
- [Sch08] ———, *minconf: projection methods for optimization with simple constraints in matlab*, <http://www.cs.ubc.ca/~schmidtm/Software/minConf.html>, 2008.
- [SCS10] Y. Saad, J. R. Chelikowsky, and S. M. Shontz, *Numerical methods for electronic structure calculations of materials*, SIAM review **52** (2010), no. 1, 3–54.
- [SEC<sup>+</sup>15] Y. Shechtman, Y. C. Eldar, O. Cohen, H. N. Chapman, J. Miao, and M. Segev, *Phase retrieval with application to optical imaging: A contemporary overview*, IEEE Signal Process. Mag. **32** (2015), no. 3, 87–109.
- [SESS11] Y. Shechtman, Y. C. Eldar, A. Szameit, and M. Segev, *Sparsity based sub-wavelength imaging with partially incoherent light via quadratic compressed sensing*, CoRR **abs/1104.4406** (2011), 1–16, 1104.4406.
- [Sor92] D. C. Sorensen, *Implicit application of polynomial filters in ak-step arnoldi method*, Siam journal on matrix analysis and applications **13** (1992), no. 1, 357–385.
- [Sor97] ———, *Implicitly restarted arnoldi/lanczos methods for large scale eigenvalue calculations*, (1997), 119–165.
- [SQW16] J. Sun, Q. Qu, and J. Wright, *A geometric analysis of phase retrieval*, Information Theory (ISIT), 2016 IEEE International Symposium on, IEEE, 2016, pp. 2379–2383.
- [Ste92] G. W. Stewart, *An updating algorithm for subspace tracking*, IEEE Transactions on Signal Processing **40** (1992), no. 6, 1535–1541.
- [TTT99] K.-C. Toh, M. J. Todd, and R. H. Tütüncü, *Sdpt3a matlab software package for semidefinite programming, version 1.3*, Optimization methods and software **11** (1999), no. 1-4, 545–581.
- [Wal63] A. Walther, *The question of phase retrieval in optics*, Optica Acta: International Journal of Optics **10** (1963), no. 1, 41–49.



- [Yan95] B. Yang, *Projection approximation subspace tracking*, IEEE Transactions on Signal processing **43** (1995), no. 1, 95–107.
- [ZH04] H. Zhang and W. W. Hager, *A nonmonotone line search technique and its application to unconstrained optimization*, SIAM journal on Optimization **14** (2004), no. 4, 1043–1056.
- [ZYLX17] H. Zhang, S. You, Z. Lin, and C. Xu, *Fast compressive phase retrieval under bounded noise.*, AAAI, 2017, pp. 2884–2890.