

TITLE OF YOUR THESIS

A Dissertation
Presented to
The Academic Faculty

By

George P. Burdell

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
School of Fiction

Georgia Institute of Technology

January 1927

Copyright © George P. Burdell 1927

TITLE OF YOUR THESIS

Approved by:

Dr. Burdell, Advisor
School of Myths
Georgia Institute of Technology

Dr. Two
School of Mechanical Engineering
Georgia Institute of Technology

Dr. Three
School of Electrical Engineering
Georgia Institute of Technology

Dr. Four
School of Computer Science
Georgia Institute of Technology

Dr. Five
School of Public Policy
Georgia Institute of Technology

Dr. Six
School of Nuclear Engineering
Georgia Institute of Technology

Date Approved: January 11, 2000

A great quote to start the thesis

George P. Burdell

A great dedication goes here.

ACKNOWLEDGEMENTS

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

TABLE OF CONTENTS

| | |
|---|----|
| Acknowledgments | v |
| List of Tables | ix |
| List of Figures | x |
| Chapter 1: Introduction | 1 |
| 1.1 Dictionaries and Dictionary Learning | 1 |
| 1.1.1 Convolutional Dictionaries | 1 |
| 1.2 Convolutional Neural Networks | 1 |
| 1.3 Multi-Layer Dictionaries | 1 |
| 1.4 Contributions and Organization of Dissertation | 1 |
| Chapter 2: Learning Dictionaries for Multi-Channel Signals | 3 |
| 2.1 Introduction | 3 |
| 2.2 Dictionary Types | 3 |
| 2.3 Pursuit and Sparse Coding | 4 |
| 2.4 ADMM | 5 |
| 2.5 Applying ADMM to the Sparse Coding Problem | 7 |
| 2.5.1 Exploiting Dictionary Structure for the Inverse Problem | 9 |

| | | |
|---|--|-----------|
| 2.6 | Sparse Coding for Multi-Channel Signals: Alternatives to My Novel Approach | 11 |
| 2.7 | A Novel Approach to Sparse Coding: ADMM with Low-Rank Updates . . . | 14 |
| 2.7.1 | Low-Rank Updates | 14 |
| 2.7.2 | Handling Dictionary Normalization | 18 |
| 2.7.3 | Dictionary Updates | 21 |
| 2.8 | Conclusion | 21 |
| Chapter 3: Learning Multi-Layer Dictionaries | | 23 |
| 3.1 | Introduction | 23 |
| 3.2 | Literature Review | 23 |
| 3.3 | Multi-Layer ADMM with Low-Rank Updates | 24 |
| 3.3.1 | Coefficients Update Equation | 26 |
| 3.3.2 | Proximal Updates | 27 |
| 3.3.3 | Dual Updates | 30 |
| 3.4 | Summary | 30 |
| Chapter 4: JPEG Artifact Removal | | 31 |
| 4.1 | Introduction | 31 |
| 4.2 | JPEG Algorithm | 31 |
| 4.3 | Literature Review | 32 |
| 4.4 | Modelling Compressed JPEG Images | 32 |
| 4.5 | Handling Quantization | 36 |
| 4.6 | Experiments | 36 |

| | | |
|--|--|-----------|
| 4.6.1 | Experiment Setup | 36 |
| 4.6.2 | Results | 36 |
| 4.7 | Conclusion | 36 |
| Chapter 5: Practical Considerations Concerning Tensorflow | | 37 |
| 5.1 | Boundary Handling | 37 |
| 5.2 | Removing Low-Frequency Signal Content | 37 |
| 5.2.1 | JPEG Artifact Removal | 37 |
| 5.3 | Tensorflow and Keras | 37 |
| 5.3.1 | Why Not Use Gradient Tape and TensorFlow-1-Style Code? | 37 |
| 5.3.2 | Shared Weights Between Layers | 38 |
| 5.3.3 | Custom Partial Gradients | 38 |
| 5.3.4 | Updating TensorFlow Variables After Applying Gradients | 39 |
| 5.3.5 | The Perils of Using Built-In Functions for Complex Tensors and Arrays | 41 |
| Appendix A: Experimental Equipment | | 43 |
| Appendix B: Data Processing | | 44 |
| References | | 47 |
| Vita | | 48 |

LIST OF TABLES

| | | |
|-----|-----------------------------------|---|
| 1.1 | This is an example Table. | 1 |
|-----|-----------------------------------|---|

LIST OF FIGURES

| | | |
|-----|---|----|
| 1.1 | This is an example Figure. | 2 |
| 2.1 | This is another example Figure, rotated to landscape orientation. | 22 |

SUMMARY

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

CHAPTER 1

INTRODUCTION

Dictionaries and Dictionary Learning

Convolutional Dictionaries

Convolutional Neural Networks

Multi-Layer Dictionaries

Contributions and Organization of Dissertation

Table 1.1: This is an example Table.

| x | f(x) | g(x) |
|---|------|------|
| 1 | 6 | 4 |
| 2 | 6 | 3 |
| 3 | 6 | 2 |
| 4 | 6 | 2 |

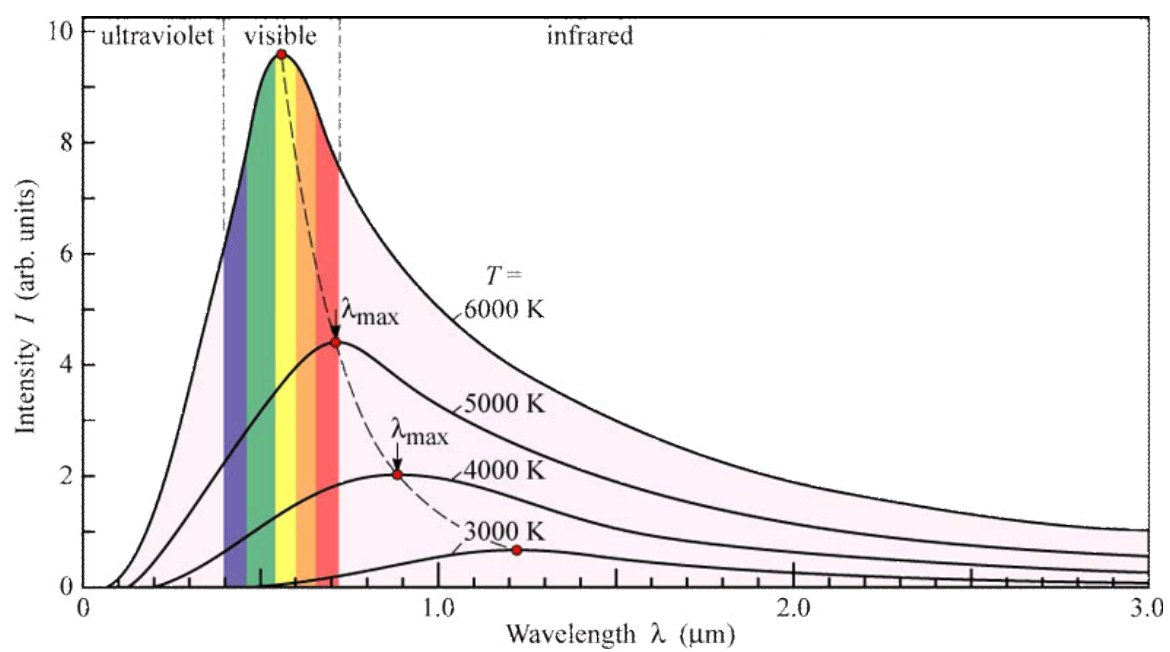


Figure 1.1: This is an example Figure.

CHAPTER 2

LEARNING DICTIONARIES FOR MULTI-CHANNEL SIGNALS

Introduction

When using a multi-layer dictionary model, the coefficients corresponding to a dictionary from one layer become the "signal" for the subsequent layer. The number of channels for this "signal" is the number of dictionary filters from the previous layer. Much of the literature on learning convolutional dictionaries is tailored to applications with signals that only have a small number of channels. This chapter presents a novel method for learning convolutional dictionaries from and for multi-channel signals.

Dictionary Types

There are many ways to construct a convolutional sparse representation of a multi-channel signal, but broadly the distinctions reduce down to if and how signal channels share dictionaries and coefficients, and if and how those non-shared entities interact across channels.

It is common in many applications for dictionary models to share dictionaries across channels, which requires the use multi-channel coefficients. If such models were used in a multi-layer dictionary model, the tensor rank would increase with each subsequent layer.

For this work, I focus instead on the multi-channel dictionary with shared coefficients. This structure matches that of convolutional neural networks, and the number of channels for a subsequent dictionary is the number of filters for the dictionary from the previous layer.

Pursuit and Sparse Coding

The dictionary model decomposes the signal s_i into a dictionary D (which generalizes to other signals) and the coefficients x_i (which are specific to the signal s_i):

$$s_i \approx Dx_i \quad (2.1)$$

(Here the subscript i specifies a particular signal and its corresponding coefficients.) A pursuit algorithm finds the coefficients x_i corresponding to a particular signal s_i for known dictionary D . If the number of dictionary atoms (columns) is larger than the dimension of the signal, then the number of unknowns is larger than the number of equations, and many solutions for x_i represent s_i equally well (at least in an L2 sense). Researchers and practitioners commonly either impose a sparsity constraint on the coefficients or add a coefficient L1 penalty to the objective function, which removes the ambiguity from the problem construction. When such a penalty or constraint is added, pursuit is sometimes called sparse coding. With the added coefficient L1 penalty, the pursuit optimization problem looks like this:

$$x_i = \arg \min_x \frac{1}{2} \|s_i - Dx\|_2^2 + \lambda \|x\|_1 \quad (2.2)$$

where λ is a hyperparameter greater than zero controlling how much the L1 norm of the coefficients is penalized. Researchers have proposed many ways to solve this problem. If the dictionary is convolutional and the number of channels is low, a standard approach is to use the Alternating direction Method of Multipliers (ADMM) algorithm.

ADMM

ADMM is a convex-optimization algorithm used to solve the optimization problem:

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{y}}{\text{minimize}} f(\mathbf{x}) + g(\mathbf{y}) \\ & \text{subject to } \mathbf{Ax} + \mathbf{By} + \mathbf{c} = \mathbf{0} \end{aligned} \quad (2.3)$$

where f and g are convex functions [1]. (I will address how to put the sparse coding problem in this form in the next section.)

The ADMM algorithm makes use of the augmented Lagrangian, a particular expression that has a saddle point at the solution to the constrained optimization problem:

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = f(\mathbf{x}) + g(\mathbf{y}) + \mathbf{u}^H(\mathbf{Ax} + \mathbf{By} + \mathbf{c}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{By} + \mathbf{c}\|_2^2 \quad (2.4)$$

where ρ is a hyperparameter greater than zero and \mathbf{u} is the dual variable for the constraints.

At the saddle-point solution, the augmented Lagrangian is at a minimum in respect to \mathbf{x} and \mathbf{y} , but at a maximum in respect to \mathbf{u} .

The ADMM algorithm is an iterative search for the saddle point of the augmented Lagrangian. Each iteration consists of a primal update for \mathbf{x} , a primal update for \mathbf{y} , and a dual update for \mathbf{u} :

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \mathbf{y}^{(t)}, \mathbf{u}^{(t)}) \quad (2.5)$$

$$\mathbf{y}^{(t+1)} = \arg \min_{\mathbf{y}} \mathcal{L}_\rho(\mathbf{x}^{(t+1)}, \mathbf{y}, \mathbf{u}^{(t)}) \quad (2.6)$$

$$\mathbf{u}^{(t+1)} = \mathbf{u}^{(t)} + \rho(\mathbf{Ax}^{(t+1)} + \mathbf{By}^{(t+1)} + \mathbf{c}) \quad (2.7)$$

The primal updates serve to move towards the minimum of the augmented Lagrangian in respect to \mathbf{x} and \mathbf{y} with \mathbf{u} fixed, and the dual update fixes \mathbf{x} and \mathbf{y} , and performs gradient ascent on \mathbf{u} with stepsize ρ . Under very mild assumptions, this process converges to the saddle point of the augmented Lagrangian, which matches the solution to the constrained optimization problem.

There are two common variations of the ADMM algorithm that this work will make use of. The first is the scaled form, which comes from completing the square for the augmented lagrangian function:

$$L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = f(\mathbf{x}) + g(\mathbf{y}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{By} + \mathbf{c} + \frac{\mathbf{u}}{\rho}\|_2^2 - \frac{1}{2\rho} \|\mathbf{u}\|_2^2 \quad (2.8)$$

The term $-\frac{1}{2\rho} \|\mathbf{u}\|_2^2$ can be ignored for the primal updates because it has no dependence on the primal variables. For this reason, it is sometimes more convenient to keep track of $\frac{\mathbf{u}}{\rho}$ instead of \mathbf{u} , since that is the form that appears in the augmented Lagrangian after completing the square.

$$\frac{\mathbf{u}^{(t+1)}}{\rho} = \frac{\mathbf{u}^{(t)}}{\rho} + \mathbf{Ax}^{(t+1)} + \mathbf{By}^{(t+1)} + \mathbf{c} \quad (2.9)$$

This form is known as scaled ADMM.

Another common variation of ADMM updates the dual variable more frequently.

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{y}^{(t)}, \mathbf{u}^{(t)}) \quad (2.10)$$

$$\mathbf{u}^{(t+\frac{1}{2})} = \mathbf{u}^{(t)} + (\alpha - 1)\rho(\mathbf{Ax}^{(t+1)} + \mathbf{By}^{(t)} + \mathbf{c}) \quad (2.11)$$

$$\mathbf{y}^{(t+1)} = \arg \min_{\mathbf{y}} L_{\rho}(\mathbf{x}^{(t+1)}, \mathbf{y}, \mathbf{u}^{(t+\frac{1}{2})}) \quad (2.12)$$

$$\mathbf{u}^{(t+1)} = \mathbf{u}^{(t+\frac{1}{2})} + \rho(\mathbf{Ax}^{(t+1)} + \mathbf{By}^{(t+1)} + \mathbf{c}) \quad (2.13)$$

When $\alpha > 1$, this is known as over-relaxation, and if $\alpha < 1$, this is known as under-relaxation.¹ α is always chosen to be greater than zero. In some applications, researchers have found using over-relaxation converges faster than without over-relaxation [2], but optimal choice of α is problem-dependent [3].

Applying ADMM to the Sparse Coding Problem

Recall from section 2.3, equation 2.2 for sparse coding.

$$\mathbf{x}_i = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{s}_i - \mathbf{Dx}\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (2.14)$$

This can be rewritten to match the ADMM form from equation 2.3:

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{y}}{\text{minimize}} \frac{1}{2} \|\mathbf{s}_i - \mathbf{Dx}\|_2^2 + \lambda \|\mathbf{y}\|_1 \\ & \text{subject to } \mathbf{y} - \mathbf{x} = \mathbf{0} \end{aligned} \quad (2.15)$$

Given sufficient iterations, \mathbf{x} and \mathbf{y} will both be close to the optimal, but they may not be equal. Either can be used as an approximate solution to the sparse coding problem.

Computing the augmented Lagrangian of convex optimization problem in expression

¹I have elected to notate over/under relaxation differently than standard, but the α is the same, and the notations are mathematically equivalent. The standard notation does not use the first dual update, and instead includes another variable $\mathbf{h}^{(t+1)} = \mathbf{Ax}^{(t+1)} - (1 - \alpha)(\mathbf{Ax}^{(t+1)} + \mathbf{By}^{(t)} + \mathbf{c})$ and substitutes $\mathbf{h}^{(t+1)}$ for $\mathbf{Ax}^{(t+1)}$ in the dual-update equation and the second primal-update equation. While more familiar to readers who have dealt with ADMM before, this standard notation complicates ADMM with an extra variable and obscures how the dual update and second primal update relate to the augmented Lagrangian.

2.15 yields the following equation:

$$L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = \frac{1}{2} \|\mathbf{s}_i - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{y}\|_1 + \frac{\rho}{2} \|\mathbf{y} - \mathbf{x} + \frac{\mathbf{u}}{\rho}\|_2^2 - \frac{1}{2\rho} \|\mathbf{u}\|_2^2 \quad (2.16)$$

Starting with the \mathbf{x} -update:

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{y}^{(t)}, \mathbf{u}^{(t)}) \quad (2.17)$$

Since the minimum is desired, setting the gradient to zero will produce the solution.

$$\nabla_{\mathbf{x}^{(t+1)}} L_\rho(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}, \mathbf{u}^{(t)}) = \mathbf{0} \quad (2.18)$$

$$\mathbf{0} = \mathbf{D}^T \mathbf{D} \mathbf{x}^{(t+1)} - \mathbf{D}^T \mathbf{s}_i + \rho \mathbf{x}^{(t+1)} - \rho(\mathbf{y}^{(t)} + \frac{\mathbf{u}^{(t)}}{\rho}) \quad (2.19)$$

$$(\rho \mathbf{I} + \mathbf{D}^T \mathbf{D}) \mathbf{x}^{(t+1)} = \mathbf{D}^T \mathbf{s}_i + \rho(\mathbf{y}^{(t)} + \frac{\mathbf{u}^{(t)}}{\rho}) \quad (2.20)$$

$$\mathbf{x}^{(t+1)} = (\rho \mathbf{I} + \mathbf{D}^T \mathbf{D})^{-1} (\mathbf{D}^T \mathbf{s}_i + \rho(\mathbf{y}^{(t)} + \frac{\mathbf{u}^{(t)}}{\rho})) \quad (2.21)$$

In subsection 2.5.1, there is a discussion of the implications of this update equation, how to compute it for cases in which the signal has a low number of channels, and the challenges it poses for signals with many channels.

If using over-relaxation², there is a dual update:

$$\frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho} = \frac{\mathbf{u}^{(t)}}{\rho} + (\alpha - 1)(\mathbf{y}^{(t)} - \mathbf{x}^{(t+1)}) \quad (2.22)$$

²or under-relaxation

Moving on to the \mathbf{y} -update:

$$\mathbf{y}^{(t+1)} = \arg \min_{\mathbf{y}} L_{\rho}(\mathbf{x}^{(t+1)}, \mathbf{y}, \mathbf{u}^{(t+\frac{1}{2})}) \quad (2.23)$$

Excluding the terms that don't include \mathbf{y} , I have

$$\mathbf{y}^{(t+1)} = \arg \min_{\mathbf{y}} \lambda \|\mathbf{y}\|_1 + \frac{\rho}{2} \|\mathbf{y} - \mathbf{x}^{(t+1)} + \frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho}\|_2^2 \quad (2.24)$$

This is a well-known problem, whose solution is

$$\mathbf{y}^{(t+1)} = S_{\frac{\lambda}{\rho}}(\mathbf{x}^{(t+1)} - \frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho}) \quad (2.25)$$

where S is the shrinkage operator:

$$S_b(x) = \begin{cases} x - b & x > b \\ 0 & -b < x < b \\ x + b & x < -b \end{cases} \quad (2.26)$$

In the case of a vector, matrix, or tensor input, the shrinkage operator is applied element by element.

Finally, the last update equation for the dual variable:

$$\frac{\mathbf{u}^{(t+1)}}{\rho} = \frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho} + \mathbf{y}^{(t+1)} - \mathbf{x}^{(t+1)} \quad (2.27)$$

Exploiting Dictionary Structure for the Inverse Problem

Returning to the \mathbf{x} update:

$$\mathbf{x}^{(t+1)} = (\rho \mathbf{I} + \mathbf{D}^T \mathbf{D})^{-1} \left(\mathbf{D}^T \mathbf{s}_i + \rho(\mathbf{y}^{(t)} + \frac{\mathbf{u}^{(t)}}{\rho}) \right) \quad (2.28)$$

For problems using a dictionary with convolutional structure, this inverse for the convolutional sparse coding problem is very structured. Exploiting this structure is important for efficient computation, because the matrix $\rho\mathbf{I} + \mathbf{D}^T\mathbf{D}$ is a large matrix.

Writing \mathbf{D} in a block structure, I have

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_{1,1}, & \dots, & \mathbf{D}_{1,M} \\ \vdots & \ddots & \vdots \\ \mathbf{D}_{C,1}, & \dots, & \mathbf{D}_{C,M} \end{bmatrix} \quad (2.29)$$

where $\mathbf{D}_{c,m}$ is a toplitz matrix capturing channel c of the m th filter of the dictionary.

Toplitz matrices are diagonalizable with Fourier eigenvectors:

$$\mathbf{D} = \begin{bmatrix} \mathcal{F}^{-1}\hat{\mathbf{D}}_{1,1}\mathcal{F}, & \dots, & \mathcal{F}^{-1}\hat{\mathbf{D}}_{1,M}\mathcal{F} \\ \vdots & \ddots & \vdots \\ \mathcal{F}^{-1}\hat{\mathbf{D}}_{C,1}\mathcal{F}, & \dots, & \mathcal{F}^{-1}\hat{\mathbf{D}}_{C,M}\mathcal{F} \end{bmatrix} \quad (2.30)$$

where $\hat{\mathbf{D}}_{c,m}$ is a diagonal matrix whose elements are the discrete Fourier transform (FFT) of channel c of the m th dictionary filter.

This sparsely banded structure is a useful form in analyzing the structure of the inverse problem:

$$(\rho\mathbf{I} + \mathbf{D}^T\mathbf{D})^{-1} = \mathcal{F}^{-1}(\rho\mathbf{I} + \hat{\mathbf{D}}^H\hat{\mathbf{D}})^{-1}\mathcal{F} \quad (2.31)$$

where

$$\hat{\mathbf{D}} = \begin{bmatrix} \hat{\mathbf{D}}_{1,1}, & \dots, & \hat{\mathbf{D}}_{1,M} \\ \vdots & \ddots & \vdots \\ \hat{\mathbf{D}}_{C,1}, & \dots, & \hat{\mathbf{D}}_{C,M} \end{bmatrix} \quad (2.32)$$

and in a slight abuse of notation, \mathcal{F} computes the FFT separately on the coefficients for each filter. In [4], Bristow et al. observe the matrix $\rho\mathbf{I} + \hat{\mathbf{D}}^H\hat{\mathbf{D}}$ is sparsely banded, so the inverse can be broken down into much smaller inverse problems, and one only needs to compute the inverse of an $M \times M$ matrix for every element in the signal. ($\rho\mathbf{I} + \hat{\mathbf{D}}^H\hat{\mathbf{D}}$ is

an $M \times M$ block matrix, whose blocks are diagonal. Each submatrix collects one element from the diagonal of each of the blocks.)

Furthermore, the maximum rank of these submatrices is C , so if C is small, these inverses can be computed even more efficiently using the Woodbury matrix identity [5] [6] [7].

According to the Woodbury matrix identity [8], for any invertible matrix U and any matrix V :

$$(U + V^H V)^{-1} = U^{-1} - U^{-1} V^H (I + V U^{-1} V^H)^{-1} V U^{-1} \quad (2.33)$$

So,

$$(\rho I + \hat{D}^H \hat{D})^{-1} = \frac{1}{\rho} I - \frac{1}{\rho} \hat{D}^H (\rho I + \hat{D} \hat{D}^H)^{-1} \hat{D} \quad (2.34)$$

This means that instead of computing the inverse of an $M \times M$ matrix for every pixel in the image, one could instead choose to compute the inverse of a $C \times C$ matrix for each pixel in the image.³

Sparse Coding for Multi-Channel Signals: Alternatives to My Novel Approach

In applying ADMM to the convolutional sparse coding problem, [5] [6] [7] exploit the low-rank structure of the inverse problem in the x update for efficient computation. Unfortunately, this relies on the number of channels being small. Broadly, there are two main approaches to avoid or simplify this challenging inverse problem: either construct a variant of the ADMM algorithm that simplifies the inverse problem, or use a proximal gradient approach that avoids it altogether.

In [9][10], the authors use the ADMM algorithm for sparse coding. They observe that if the dictionary is a tight frame, that is, $DD^T = I$, then the inverse can be simplified without

³Generally, Cholesky or LDLT decomposition would be preferable to explicitly computing the inverse, and the efficiency gains due to the Woodbury matrix identity are relevant regardless of the chosen representation.

using the frequency representation.

$$(\rho \mathbf{I} + \mathbf{D}^T \mathbf{D})^{-1} = \frac{1}{\rho} \mathbf{I} - \frac{1}{\rho(\rho + 1)} \mathbf{D}^T \mathbf{D} \quad (2.35)$$

This produces the \mathbf{x} update equation

$$\mathbf{x}^{(t+1)} = \frac{1}{\rho + 1} \mathbf{D}^T \mathbf{s} + \left(\mathbf{I} - \frac{1}{\rho + 1} \mathbf{D}^T \mathbf{D} \right) \left(\mathbf{z}^{(t)} - \frac{\gamma^{(t)}}{\rho} \right) \quad (2.36)$$

In their work, they use the equations built on the assumption that the dictionary is a tight frame, but develop no mechanism to ensure that their assumption is accurate. Thus, ultimately $\frac{1}{\rho} \mathbf{I} - \frac{1}{\rho(\rho+1)} \mathbf{D}^T \mathbf{D}$ merely serves as an approximation to $(\rho \mathbf{I} + \mathbf{D}^T \mathbf{D})^{-1}$. Empirically, they observe the algorithm converges, but the dictionaries they learn are not tight frames, so the solution they converge to is not optimal⁴.

Other works avoid the ADMM algorithm entirely.

The iterative shrinkage thresholding algorithm (ISTA) is an iterative algorithm that minimizes the sum of two convex functions f and g . f is required to be smooth. It is helpful for f to be easily differentiable and g to have a simple proximal operator.

$$\text{prox}_g(\boldsymbol{\mu}) = \arg \min_{\boldsymbol{\nu}} \frac{1}{2} \|\boldsymbol{\nu} - \boldsymbol{\mu}\|_2^2 + g(\boldsymbol{\nu}) \quad (2.37)$$

Then, ISTA has the following update equation, where the constant L controls step size.

$$\mathbf{x}^{(t+1)} = \text{prox}_g \left(\mathbf{x}^{(t)} - \frac{1}{L} \nabla_{\mathbf{x}} f(\mathbf{x}^{(t)}) \right) \quad (2.38)$$

FISTA is similar to ISTA, but adds momentum [11].

$$\mathbf{z}^{(t+1)} = \text{prox}_g \left(\mathbf{x}^{(t)} - \frac{1}{L} \nabla_{\mathbf{x}} f(\mathbf{x}^{(t)}) \right) \quad (2.39)$$

⁴The solution does not minimize the sparse coding objective function.

$$r^{(t+1)} = \frac{1}{2} \left(1 + \sqrt{1 + 4(r^{(t)})^2} \right) \quad (2.40)$$

$$\mathbf{x}^{(t+1)} = \mathbf{z}^{(t+1)} + \frac{r^{(t)} - 1}{r^{(t+1)}} (\mathbf{z}^{(t+1)} - \mathbf{x}^{(t)}) \quad (2.41)$$

Applying FISTA to the sparse coding problem, $\frac{1}{2} \|\mathbf{s} - \mathbf{D}\mathbf{x}\|_2^2$ is straightforward to differentiate and $\lambda \|\mathbf{x}\|_1$ has a simple proximal operator.

$$\nabla_{\mathbf{x}} \left(\frac{1}{2} \|\mathbf{s} - \mathbf{D}\mathbf{x}\|_2^2 \right) = \mathbf{D}^T \mathbf{D}\mathbf{x} - \mathbf{D}^T \mathbf{s} \quad (2.42)$$

$$\text{prox}_{\lambda \|\cdot\|_1}(\cdot) = \mathbf{S}_\lambda \quad (2.43)$$

So, the FISTA equations for convolutional basis pursuit are the following:

$$\mathbf{z}^{(t+1)} = \mathbf{S}_\lambda \left(\mathbf{x}^{(t)} - \frac{1}{L} \mathbf{D}^T (\mathbf{D}\mathbf{x}^{(t)} - \mathbf{s}) \right) \quad (2.44)$$

$$r^{(t+1)} = \frac{1}{2} \left(1 + \sqrt{1 + 4(r^{(t)})^2} \right) \quad (2.45)$$

$$\mathbf{x}^{(t+1)} = \mathbf{z}^{(t+1)} + \frac{r^{(t)} - 1}{r^{(t+1)}} (\mathbf{z}^{(t+1)} - \mathbf{x}^{(t)}) \quad (2.46)$$

In [7], Wohlberg compares FISTA to ADMM on a sparse coding task and finds FISTA converges much slower than ADMM. However, the comparison is made on signals with few channels, so ADMM is able to exploit the structure of \mathbf{D} for efficient \mathbf{x} updates.

In a recent work [12], Chodosh and Lucey use updates prox-linear updates using more general convex solver methods detailed in [13].

The updates come from the formula:

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} \left(\nabla f(\bar{\mathbf{x}}^{(t)}) \right)^T (\mathbf{x} - \bar{\mathbf{x}}^{(t)}) + \frac{L}{2} \|\mathbf{x} - \bar{\mathbf{x}}^{(t)}\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (2.47)$$

where $\bar{\mathbf{x}}^{(t)} = \mathbf{x}^{(t)} + \omega_k(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})$ and ω_t is a momentum factor.

This yields the update equation⁵:

$$\mathbf{x}^{(t+1)} = \mathbf{S}_{\frac{\lambda}{L}} \left(\bar{\mathbf{x}}^{(t)} + \mathbf{D}^T(\mathbf{s} - \mathbf{D}\bar{\mathbf{x}}) \right) \quad (2.48)$$

While neither Chodosh and Lucey nor the work they cite mentions FISTA, the resemblance is very close. There are two distinctions:

1. Momentum is computed slightly differently: to match FISTA, the prox-linear updates would need to use $\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t-1)}$ for momentum. Instead, they use $\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}$.
2. The prox-linear approach scales the momentum steps differently.

Given these similarities, it is likely the performance between the two methods is similar.

A Novel Approach to Sparse Coding: ADMM with Low-Rank Updates

In this section, I present a novel approach to sparse coding for signals with a large number of channels. The approach uses the ADMM algorithm described in section 2.4 and will share many similarities to the standard ADMM sparse coding approach described in section 2.5 for signals with few channels.

Low-Rank Updates

Under many circumstances, inverse representations can be updated efficiently, provided the update adheres to a low-rank structure. Recall the frequency representation of the

⁵In their paper, they add a non-negativity constraint and allow different λ for the coefficients of each filter (and possibly spatially varied as well). They also are constructing the equations specifically for a multi-layer network. I simplified their equations to illustrate how their approach relates to the FISTA algorithm.

convolutional dictionary:

$$\hat{\mathbf{D}} = \begin{bmatrix} \hat{\mathbf{D}}_{1,1}, & \dots, & \hat{\mathbf{D}}_{1,M} \\ \vdots & \ddots & \vdots \\ \hat{\mathbf{D}}_{C,1}, & \dots, & \hat{\mathbf{D}}_{C,M} \end{bmatrix} \quad (2.49)$$

where $\hat{\mathbf{D}}_{c,m}$ is diagonal for all c and m . Let $\hat{\mathbf{D}}_{c,m}[\hat{k}]$ be the \hat{k} th element of the diagonal and let

$$\hat{\mathbf{D}}[\hat{k}] = \begin{bmatrix} \hat{\mathbf{D}}_{1,1}[\hat{k}], & \dots, & \hat{\mathbf{D}}_{1,M}[\hat{k}] \\ \vdots & \ddots & \vdots \\ \hat{\mathbf{D}}_{C,1}[\hat{k}], & \dots, & \hat{\mathbf{D}}_{C,M}[\hat{k}] \end{bmatrix} \quad (2.50)$$

Then $\hat{\mathbf{D}}[\hat{k}]$ is a $C \times M$ matrix collecting the \hat{k} th frequency of all channels and filters of \mathbf{D} .

Thus, $(\rho\mathbf{I} + \hat{\mathbf{D}}^H \hat{\mathbf{D}})^{-1}$ really consists of \hat{k} separate inverse problems: $(\rho\mathbf{I} + \hat{\mathbf{D}}^H[\hat{k}] \hat{\mathbf{D}}[\hat{k}])^{-1}$.

Consider the update equation.

$$\hat{\mathbf{D}}[\hat{k}]^{(n+1)} = \hat{\mathbf{D}}[\hat{k}]^{(n)} + \mathbf{U}\mathbf{V}[\hat{k}]^H \quad (2.51)$$

where \mathbf{U} is an orthogonal matrix of size $C \times L$ and $\mathbf{V}[\hat{k}]$ is an orthogonal matrix of size $M \times L$.⁶

Then,

$$(\rho\mathbf{I} + (\hat{\mathbf{D}}^{(n+1)})^H[\hat{k}] \hat{\mathbf{D}}^{(n+1)}[\hat{k}])^{-1} = \rho\mathbf{I} + (\mathbf{D}^{(n)}[\hat{k}] + \mathbf{U}\mathbf{V}^H[\hat{k}])^H (\mathbf{D}^{(n)}[\hat{k}] + \mathbf{U}\mathbf{V}^H[\hat{k}]) \quad (2.52)$$

⁶ \mathbf{U} and \mathbf{V} are also iteration specific, but to notate that would over-clutter the equations. For efficient, low-rank updates to the inverse representation, I could allow both \mathbf{U} and \mathbf{V} to vary in respect to frequency \hat{k} (instead of just \mathbf{V}). However, I also need to limit the spatial support of the dictionary (so that the filter size is small), and preventing \mathbf{U} from varying across frequency is part of a means to satisfy that constraint.

For brevity and simplicity, I will drop the notation indexing the frequency \hat{k} and selecting the iteration n for matrix $\hat{\mathbf{D}}^{(n)}[\hat{k}]$ and simply use $\hat{\mathbf{D}}$ instead. However, the reader should keep in mind the $\hat{\mathbf{D}}$ here is a dense $C \times M$ matrix capturing the component of the dictionary \mathbf{D} for implicit frequency \hat{k} , not the sparsely banded $\hat{\mathbf{D}}$ of size $KC \times M$ from earlier in this section.

$$(\rho \mathbf{I} + \hat{\mathbf{D}}^H \hat{\mathbf{D}})^{-1} = \rho \mathbf{I} + \mathbf{D}^H \mathbf{D} + \mathbf{V} \mathbf{U}^H \mathbf{U} \mathbf{V}^H + \mathbf{V} \mathbf{U}^H \mathbf{D} + \mathbf{D}^H \mathbf{U} \mathbf{V}^H \quad (2.53)$$

Given that \mathbf{V} and \mathbf{U} are orthogonal matrices, $\mathbf{V} \mathbf{U}^H \mathbf{U} \mathbf{V}^H$ can easily be broken into L rank-one Hermitian updates.

$$\mathbf{V} \mathbf{U}^H \mathbf{U} \mathbf{V}^H = \sum_{\ell=1}^L \mathbf{u}_{\ell}^H \mathbf{u}_{\ell} \mathbf{v}_{\ell} \mathbf{v}_{\ell}^H \quad (2.54)$$

Similarly, $\mathbf{V} \mathbf{U}^H \mathbf{D} + \mathbf{D}^H \mathbf{U} \mathbf{V}^H$ can be broken into L Hermitian, rank-two updates:

$$\mathbf{V} \mathbf{U}^H \mathbf{D} + \mathbf{D}^H \mathbf{U} \mathbf{V}^H = \sum_{\ell=1}^L \mathbf{v}_{\ell} \mathbf{u}_{\ell}^H \mathbf{D} + \mathbf{D}^H \mathbf{u}_{\ell} \mathbf{v}_{\ell}^H \quad (2.55)$$

Inverse representations can be efficiently updated if the update is Hermitian and rank one. The details of such updates are discussed in the appendix.

The Hermitian rank-two update consists of two rank-one terms, but the terms are not Hermitian, complicating the update process. However, this can be resolved through eigen-decomposition.

$$\mathbf{v}_{\ell} \mathbf{u}_{\ell}^H \hat{\mathbf{D}} + \hat{\mathbf{D}} \mathbf{u}_{\ell} \mathbf{v}_{\ell}^H = \begin{bmatrix} \mathbf{v}_{\ell} & \hat{\mathbf{D}}^H \mathbf{u}_{\ell} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{D}}^H \mathbf{u}_{\ell} & \mathbf{v}_{\ell} \end{bmatrix}^H \quad (2.56)$$

While matrix products are not commutative, some of the eigenvalues of matrix products are commutative.

Furthermore, for general matrices \mathbf{A} and \mathbf{B} the eigenvectors of \mathbf{AB} and \mathbf{BA} are re-

lated:

$$\mathbf{B}\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \implies \mathbf{A}\mathbf{B}\mathbf{A}\mathbf{x} = \lambda\mathbf{A}\mathbf{x} \quad (2.57)$$

where λ is the eigenvalue and \mathbf{x} is a vector.⁷

So, if \mathbf{x} is an eigenvector of $\mathbf{B}\mathbf{A}$, $\mathbf{A}\mathbf{x}$ is an eigenvector of $\mathbf{A}\mathbf{B}$.

$$\begin{bmatrix} \hat{\mathbf{D}}^H \mathbf{u}_\ell & \mathbf{v}_\ell \end{bmatrix}^H \begin{bmatrix} \mathbf{v}_\ell & \hat{\mathbf{D}}^H \mathbf{u}_\ell \end{bmatrix} = \begin{bmatrix} \mathbf{u}_\ell^H \hat{\mathbf{D}} \mathbf{v}_\ell & \mathbf{u}_\ell^H \hat{\mathbf{D}} \hat{\mathbf{D}}^H \mathbf{u}_\ell \\ \mathbf{v}_\ell^H \mathbf{v}_\ell & \mathbf{v}_\ell^H \hat{\mathbf{D}}^H \mathbf{u}_\ell \end{bmatrix} \quad (2.58)$$

The eigenvalues and corresponding eigenvectors of a 2×2 matrix can be computed using the quadratic formula. Assuming that the 2×2 matrix has 2 distinct eigenvalues, the expressions for these are below.⁸

$$\text{eigval} \left(\begin{bmatrix} a & b \\ c & a^* \end{bmatrix} \right) = \text{real}(a) \pm \sqrt{bc - (\text{imag}(a))^2} \quad (2.59)$$

$$\text{eigvec} \left(\begin{bmatrix} a & b \\ c & a^* \end{bmatrix} \right) = \begin{bmatrix} b \\ -j \text{imag}(a) \pm \sqrt{bc - (\text{imag}(a))^2} \end{bmatrix} \quad (2.60)$$

For the sake of brevity, I will drop the subscripts for \mathbf{u} and \mathbf{v} .

Letting $\eta_u = \|\hat{\mathbf{D}}^H \mathbf{u}\|_2^2$, $\eta_v = \|\mathbf{v}\|_2^2$, and $\eta_{u,v} = \mathbf{u}^H \hat{\mathbf{D}} \mathbf{v}$:

$$\text{eigval} \left(\begin{bmatrix} \eta_{u,v} & \eta_u \\ \eta_v & \eta_{u,v}^* \end{bmatrix} \right) = \text{real}(\eta_{u,v}) \pm \sqrt{\eta_v \eta_u - (\text{imag}(\eta_{u,v}))^2} \quad (2.61)$$

⁷I appologize for the reuse of certain variables here. Please do not confuse this \mathbf{x} for the sparse coding coefficients, λ for the L1 penalty factor applied to the coefficients, or \mathbf{A} and \mathbf{B} for the matrices in the ADMM constraints.

⁸It is not guarenteed the 2×2 matrix will have 2 distinct eigenvalues. In the practical considerations section, I consider those cases.

$$\text{eigvec} \left(\begin{bmatrix} \eta_{\mathbf{u},\mathbf{v}} & \eta_{\mathbf{u}} \\ \eta_{\mathbf{v}} & \eta_{\mathbf{u},\mathbf{v}}^* \end{bmatrix} \right) = \begin{bmatrix} \eta_{\mathbf{u}} \\ -j \text{imag}(\eta_{\mathbf{u},\mathbf{v}}) \pm \sqrt{\eta_{\mathbf{v}}\eta_{\mathbf{u}} - (\text{imag}(\eta_{\mathbf{u},\mathbf{v}}))^2} \end{bmatrix} \quad (2.62)$$

Therefore,

$$\text{eigvec}(\mathbf{v}\mathbf{u}^H \hat{\mathbf{D}} + \hat{\mathbf{D}}^H \mathbf{u}\mathbf{v}^H) = \eta_{\mathbf{u}}\mathbf{v} + \left(-j \text{imag}(\eta_{\mathbf{u},\mathbf{v}}) \pm \sqrt{\eta_{\mathbf{v}}\eta_{\mathbf{u}} - (\text{imag}(\eta_{\mathbf{u},\mathbf{v}}))^2} \right) \hat{\mathbf{D}}^H \mathbf{u} \quad (2.63)$$

$$\text{eigval}(\mathbf{v}\mathbf{u}^H \hat{\mathbf{D}} + \hat{\mathbf{D}}^H \mathbf{u}\mathbf{v}^H) = \text{real}(\eta_{\mathbf{u},\mathbf{v}}) \pm \sqrt{\eta_{\mathbf{v}}\eta_{\mathbf{u}} - (\text{imag}(\eta_{\mathbf{u},\mathbf{v}}))^2} \quad (2.64)$$

This decomposition splits the Hermitian rank-two update into two Hermitian rank-one updates that can be used to update the inverse representation for $\rho\mathbf{I} + \hat{\mathbf{D}}^H[\hat{k}]\hat{\mathbf{D}}[\hat{k}]$.

Recall once again, the update under consideration:

$$\hat{\mathbf{D}}^{(n+1)}[\hat{k}] = \hat{\mathbf{D}}^{(n)}[\hat{k}] + \mathbf{U}\mathbf{V}^H[\hat{k}] \quad (2.65)$$

This update must be of rank L at every frequency. Furthermore, the dictionary filter is spatially limited to its filter size. This second constraint is met if $\mathbf{V}^H[\hat{k}]$ is similarly spatially limited.

Handling Dictionary Normalization

Consider the optimization problem:

$$\begin{aligned} & \min_{\mathbf{x}, \mathbf{y}, \mathbf{z}} \frac{1}{2} \|\mathbf{s} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{y}\|_1 \\ & \text{subject to } \mathbf{R}^{-1}\mathbf{y} - \mathbf{R}^{-1}\mathbf{x} = 0 \end{aligned} \quad (2.66)$$

where \mathbf{R} is a diagonal matrix with scaled identity blocks:

$$\mathbf{R} = \begin{bmatrix} r_1 \mathbf{I} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & r_2 \mathbf{I} & & \vdots \\ \vdots & & \ddots & \\ \mathbf{0} & \dots & & r_M \mathbf{I} \end{bmatrix} \quad (2.67)$$

This optimization problem has the augmented Lagrangian function:

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = \frac{1}{2} \|\mathbf{s} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{y}\|_1 + \mathbf{u}^H \mathbf{R}^{-1}(\mathbf{y} - \mathbf{x}) + \frac{\rho}{2} \|\mathbf{R}^{-1}(\mathbf{y} - \mathbf{x})\|_2^2 \quad (2.68)$$

$$\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = -\mathbf{R}^{-1}\mathbf{u} - \mathbf{D}^H \mathbf{s} + \mathbf{D}^T \mathbf{D}\mathbf{x} + \rho \mathbf{R}^{-2}\mathbf{x} - \rho \mathbf{R}^{-2}\mathbf{y} \quad (2.69)$$

For $\mathbf{x}, \mathbf{y}, \mathbf{u}$ such that $\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = 0$:

$$(\rho \mathbf{R}^{-2} + \mathbf{D}^T \mathbf{D})\mathbf{x} = \rho \mathbf{R}^{-2}\mathbf{y} + \mathbf{R}^{-1}\mathbf{u} + \mathbf{D}^T \mathbf{s} \quad (2.70)$$

$$\mathbf{R}^{-1} (\rho \mathbf{I} + (\mathbf{D}\mathbf{R})^T (\mathbf{D}\mathbf{R})) \mathbf{R}^{-1}\mathbf{x} = \rho \mathbf{R}^{-2}\mathbf{y} + \mathbf{R}^{-1}\mathbf{u} + \mathbf{D}^T \mathbf{s} \quad (2.71)$$

$$(\rho \mathbf{I} + (\mathbf{D}\mathbf{R})^T (\mathbf{D}\mathbf{R})) \mathbf{R}^{-1}\mathbf{x} = \rho \mathbf{R}^{-1}\mathbf{z} + \mathbf{u} + (\mathbf{D}\mathbf{R})^T \mathbf{s} \quad (2.72)$$

$$\mathbf{R}^{-1}\mathbf{x} = (\rho \mathbf{I} + (\mathbf{D}\mathbf{R})^H (\mathbf{D}\mathbf{R}))^{-1} (\rho \mathbf{R}^{-1}\mathbf{y} + \mathbf{u} + (\mathbf{D}\mathbf{R})^T \mathbf{s}) \quad (2.73)$$

So,

$$\min_{\mathbf{x}} L_{\rho}(\mathbf{x}, \mathbf{y}, \mathbf{u}) = \mathbf{R} \left(\rho \mathbf{I} + (\mathbf{D}\mathbf{R})^T (\mathbf{D}\mathbf{R}) \right)^{-1} \left(\rho \mathbf{R}^{-1} \mathbf{y} + \mathbf{u} + (\mathbf{D}\mathbf{R})^T \mathbf{s} \right) \quad (2.74)$$

However, taking a similar approach to that of scaled ADMM, it will be simpler to track $\mathbf{R}^{-1}\mathbf{x}$ instead of \mathbf{x} directly.

$$\mathbf{R}^{-1}\mathbf{x}^{(t+1)} = \left(\rho \mathbf{I} + (\mathbf{D}\mathbf{R})^T (\mathbf{D}\mathbf{R}) \right)^{-1} \left((\mathbf{D}\mathbf{R})^T \mathbf{s} + \rho \left(\mathbf{R}^{-1}\mathbf{y}^{(t)} + \frac{\mathbf{u}^{(t)}}{\rho} \right) \right) \quad (2.75)$$

Moving on to the \mathbf{y} update,

$$\min_{\mathbf{y}} L_{\rho}(\mathbf{x}, \mathbf{y}, \mathbf{u}) = S_{\frac{\lambda \mathbf{R}^2}{\rho}} \left(\mathbf{x} - \frac{\mathbf{R}\mathbf{u}}{\rho} \right) \quad (2.76)$$

$$\mathbf{R}^{-1}\mathbf{y}^{(t+1)} = S_{\frac{\lambda \mathbf{R}}{\rho}} \left(\mathbf{R}^{-1}\mathbf{x}^{(t+1)} - \frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho} \right) \quad (2.77)$$

Finally, the dual updates are

$$\frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho} = \frac{\mathbf{u}^{(t)}}{\rho} + (\alpha - 1)(\mathbf{R}^{-1}\mathbf{y}^{(t)} - \mathbf{R}^{-1}\mathbf{x}^{(t+1)}) \quad (2.78)$$

$$\frac{\mathbf{u}^{(t+1)}}{\rho} = \frac{\mathbf{u}^{(t+\frac{1}{2})}}{\rho} + \mathbf{R}^{-1}\mathbf{y}^{(t+1)} - \mathbf{R}^{-1}\mathbf{x}^{(t+1)} \quad (2.79)$$

Thus, with this modification to the sparse coding optimization problem, the inverse representation used in the \mathbf{x} updates can be updated efficiently (given that the dictionary updates adhere to a particular low-rank structure), and normalization can be handed through a normalization factor \mathbf{R}^{-1} .

Dictionary Updates

A desired dictionary update can be computed through stochastic gradient descent or other dictionary-update scheme. However, the actual update will only approximate the desired one, obtained through singular-value decomposition.

Conclusion

In this chapter, I have derived a novel sparse coding algorithm for signals with a large number of channels. One of the steps in the iterative algorithm involves solving an inverse problem, but the optimization is constructed such that the representation of the inverse can be updated efficiently.

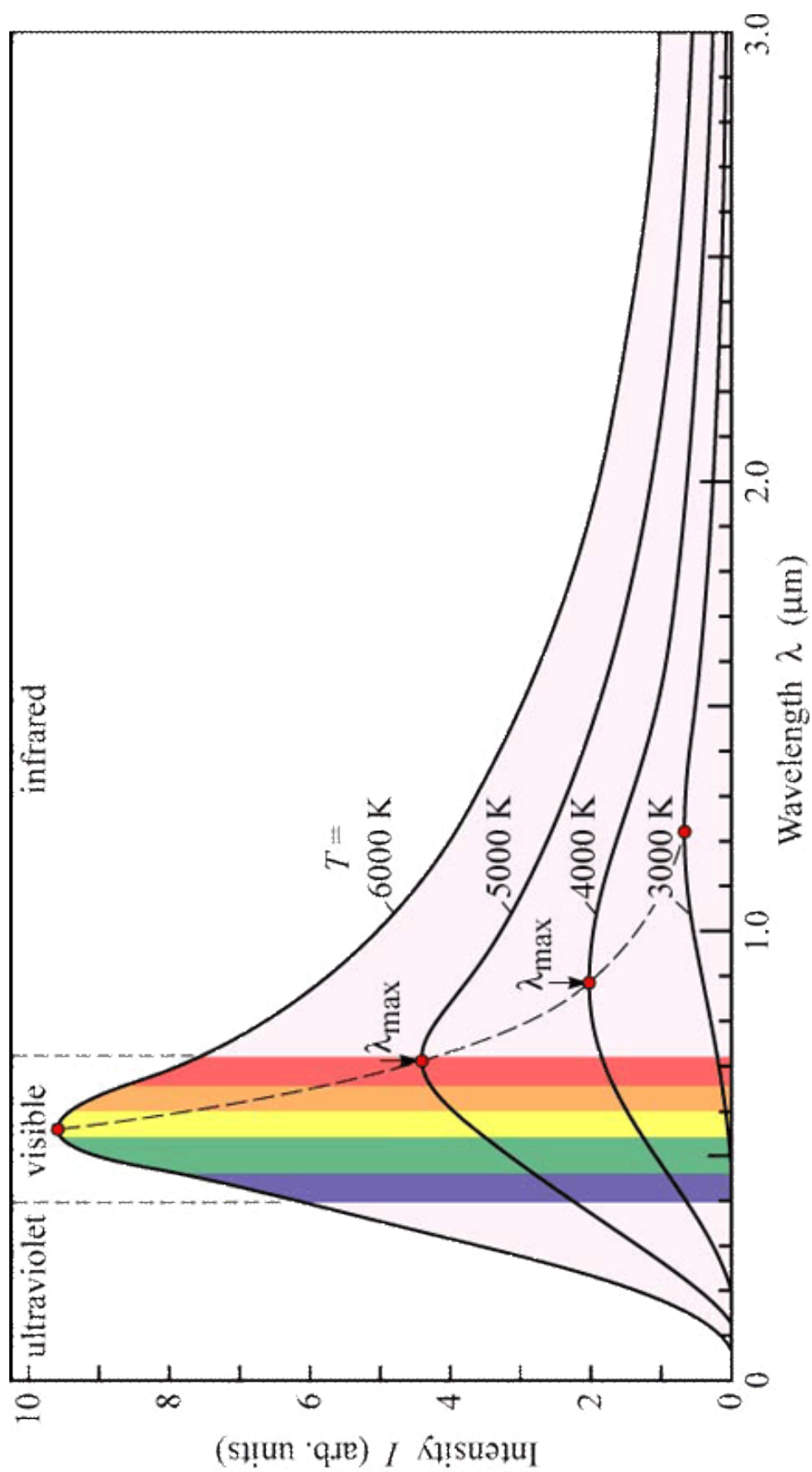


Figure 2.1: This is another example Figure, rotated to landscape orientation.

CHAPTER 3

LEARNING MULTI-LAYER DICTIONARIES

Introduction

A multi-layer dictionary model is composed of multiple dictionaries; the model treats the dictionary coefficients of a previous layer as the signal for the subsequent layer. This model dates back to Zeiler’s Deconvolutional Neural Networks [14] and can be thought of as a deep autoencoder [15, Chapter 14][16]. Some researchers have interpreted convolutional neural networks as multi-layer dictionary models, the convolution and its corresponding rectified linear units serving as a crude pursuit algorithm [17]. In this chapter, I explain how to apply the novel dictionary learning algorithm from the prior chapter to the multi-layer dictionary learning problem.

Literature Review

In 2010, Zeiler et al. proposed a multi-layer dictionary model termed a deconvolutional network. The learning process for dictionary filters is entirely unsupervised, and they learn their filters layer-by-layer. Their algorithm is greedy in the sense that there is no feedback from subsequent layers to influence the learning process on the previous layer. This approach was tested both on the task of removing added gaussian noise to images, and also as a feature extraction method for object recognition on the Caltech-101 dataset [18]. While this research drew a lot of attention at the time, as the success of alternative models like convolutional neural networks grew, the popularity of deconvolutional networks decreased.

Multi-layer dictionaries also appear in Bayesian models, going by names such as hierarchical convolutional factor analysis [19][20] and deep deconvolutional learning [21]. These networks use probabilistic models to prune network architecture and provide interpretable

dictionaries. Inference can be slow.

In more recent work, [10] and [9] use ADMM for pursuit on a multi-layer dictionary model. Their pursuit algorithm attempts to solve the minimization problem:

$$\underset{\mathbf{x}}{\text{minimize}} = \sum_{\ell=1}^L \frac{\mu_{\ell}}{2} \|\mathbf{x}_{\ell-1} - \mathbf{D}_{\ell} \mathbf{x}_{\ell}\|_2^2 + \lambda_{\ell} \|\mathbf{x}_{\ell}\|_1 \quad (3.1)$$

where $\mathbf{x}_0 = \mathbf{s}$ is the signal. They convert this to a constrained optimization for the ADMM algorithm.

$$\underset{\mathbf{x}, \mathbf{z}}{\min} \sum_{\ell=1}^L \frac{\mu_{\ell}}{2} \|\mathbf{z}_{\ell-1} - \mathbf{D}_{\ell} \mathbf{x}_{\ell}\|_2^2 + \lambda_{\ell} \|\mathbf{z}_{\ell}\|_1 \quad (3.2)$$

$$\text{subject to } \mathbf{z}_{\ell} - \mathbf{x}_{\ell} = 0$$

where $\mathbf{z}_0 = \mathbf{s}$ is the signal. The \mathbf{x} updates involve solving an inverse problem. They use a tight-frame assumption to approximate the inverse.

Finally, in [12], Chodosh and Lucey use a similar model to [10] and [9], but replace the ADMM approach with FISTA-like linear-proximal iterative steps.

Multi-Layer ADMM with Low-Rank Updates

This chapter demonstrates how to apply the novel sparse coding method for multi-channel signals to a multi-layer dictionary pursuit problem.

To start off, it is helpful to write a multi-layer dictionary optimization problem. To keep things compact, let $\mathbf{x}_0 = \mathbf{s}$.

I use the same multi-layer dictionary model as [10] and [9], but I use different ADMM constraints:

$$\underset{\mathbf{x}}{\text{minimize}} = \sum_{\ell=1}^L \frac{\mu_{\ell}}{2} \|\mathbf{x}_{\ell-1} - \mathbf{D}_{\ell} \mathbf{x}_{\ell}\|_2^2 + \lambda_{\ell} \|\mathbf{x}_{\ell}\|_1 \quad (3.3)$$

Applying the ADMM algorithm, I add a secondary primal variable \mathbf{z}_{ℓ} for each layer ℓ and constrain it to be equal to \mathbf{x}_{ℓ} . As in the previous chapter, I scale this constraint so that

the inverse representation can be updated efficiently without losing the normalized quality of the dictionary. This replaces the tight-frame assumption used in [10] and [9]. Again, keeping things compact, let $\mathbf{z}_0 = \mathbf{s}$.

$$\min_{\mathbf{x}, \mathbf{z}} \sum_{\ell=1}^L \frac{\mu_\ell}{2} \|\mathbf{z}_{\ell-1} - \mathbf{D}_\ell \mathbf{x}_\ell\|_2^2 + \lambda_\ell \|\mathbf{z}_\ell\|_1 \quad (3.4)$$

$$\text{subject to } \sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \mathbf{z}_\ell - \sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \mathbf{x}_\ell = 0$$

where $\mathbf{z}_0 = \mathbf{s}$ is not a primal variable, but instead the signal itself.

This optimization problem has the augmented Lagrangian function:

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\gamma}) = f(\mathbf{x}, \mathbf{z}) + \sum_{\ell=1}^L \frac{\rho}{2} \|\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1}(\mathbf{z}_\ell - \mathbf{x}_\ell) + \frac{\boldsymbol{\gamma}_\ell}{\rho}\|_2^2 - \frac{1}{2\rho} \|\boldsymbol{\gamma}_\ell\|_2^2 \quad (3.5)$$

where

$$f(\mathbf{x}, \mathbf{z}) = \sum_{\ell=1}^L \frac{\mu_\ell}{2} \|\mathbf{z}_{\ell-1} - \mathbf{D}_\ell \mathbf{x}_\ell\|_2^2 + \lambda_\ell \|\mathbf{z}_\ell\|_1 \quad (3.6)$$

Recall that the ADMM algorithm (with relaxation) iteratively alternates between primal and dual updates, using four update equations:

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{z}^{(t)}, \boldsymbol{\gamma}^{(t)}) \quad (3.7)$$

$$\frac{\boldsymbol{\gamma}^{(t+\frac{1}{2})}}{\rho} = \frac{\boldsymbol{\gamma}^{(t)}}{\rho} + (\alpha - 1)(\mathbf{A}\mathbf{x}^{(t+1)} + \mathbf{B}\mathbf{z}^{(t)} + \mathbf{c}) \quad (3.8)$$

$$\mathbf{z}^{(t+1)} = \arg \min_{\mathbf{z}} \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{z}, \boldsymbol{\gamma}^{(t+\frac{1}{2})}) \quad (3.9)$$

$$\frac{\boldsymbol{\gamma}^{(t+1)}}{\rho} = \frac{\boldsymbol{\gamma}^{(t+\frac{1}{2})}}{\rho} + \mathbf{A}\mathbf{x}^{(t+1)} + \mathbf{B}\mathbf{z}^{(t+1)} + \mathbf{c} \quad (3.10)$$

where $\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} + \mathbf{c} = \mathbf{0}$ are the affine constraints.

The first of these updates is the \mathbf{x} update, which updates the coefficients.

Coefficients Update Equation

The coefficients update comes from equation 3.7, which can be derived through setting the gradient of the Lagrangian equal to zero and solving for \mathbf{x} .

$$\nabla_{\mathbf{x}_\ell} f(\mathbf{x}, \mathbf{z}) = \mu_\ell \mathbf{D}_\ell^T \mathbf{D}_\ell \mathbf{x}_\ell - \mu_\ell \mathbf{D}_\ell^T \mathbf{z}_{\ell-1} \quad (3.11)$$

$$\nabla_{\mathbf{x}_\ell} \frac{1}{2} \|\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} (\mathbf{z}_\ell - \mathbf{x}_\ell) + \frac{\gamma_\ell}{\rho}\|_2^2 = \mu_\ell \mathbf{R}_\ell^{-2} \mathbf{x} - \mu_\ell \mathbf{R}_\ell^{-2} \mathbf{z}_\ell - \frac{\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \gamma_\ell}{\rho} \quad (3.12)$$

Therefore,

$$\nabla_{\mathbf{x}_\ell} \mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \gamma) = \mu_\ell \mathbf{D}_\ell^T \mathbf{D}_\ell \mathbf{x}_\ell - \mu_\ell \mathbf{D}_\ell^T \mathbf{z}_{\ell-1} + \rho \left(\mu_\ell \mathbf{R}_\ell^{-2} \mathbf{x} - \mu_\ell \mathbf{R}_\ell^{-2} \mathbf{z}_\ell - \frac{\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \gamma_\ell}{\rho} \right) \quad (3.13)$$

For $\mathbf{x}, \mathbf{z}, \gamma$, such that $\nabla_{\mathbf{x}_\ell} \mathcal{L}_\rho(\mathbf{x}_1, \dots, \mathbf{x}_L, \mathbf{z}_1, \dots, \mathbf{z}_L, \gamma_1, \dots, \gamma_L) = 0$:

$$\mu_\ell (\rho \mathbf{R}_\ell^{-2} + \mathbf{D}_\ell^T \mathbf{D}_\ell) \mathbf{x}_\ell = \mu_\ell \mathbf{D}_\ell^T \mathbf{z}_{\ell-1} + \rho \mu_\ell \mathbf{R}_\ell^{-2} \mathbf{z}_\ell + \sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \gamma_\ell \quad (3.14)$$

$$(\rho \mathbf{R}_\ell^{-2} + \mathbf{D}_\ell^T \mathbf{D}_\ell) \mathbf{x}_\ell = \mathbf{D}_\ell^T \mathbf{z}_{\ell-1} + \rho \mathbf{R}_\ell^{-2} \mathbf{z}_\ell + \frac{\mathbf{R}_\ell^{-1} \gamma_\ell}{\sqrt{\mu_\ell}} \quad (3.15)$$

$$\mathbf{x}_\ell = (\rho \mathbf{R}_\ell^{-2} + \mathbf{D}_\ell^T \mathbf{D}_\ell)^{-1} \left(\mathbf{D}_\ell^T \mathbf{z}_{\ell-1} + \rho \mathbf{R}_\ell^{-2} \mathbf{z}_\ell + \frac{\mathbf{R}_\ell^{-1} \gamma_\ell}{\sqrt{\mu_\ell}} \right) \quad (3.16)$$

This solution is the \mathbf{x} update for the ADMM algorithm, but a couple extra steps can put

it into a form that is easier to use.

$$\mathbf{x}_\ell = \mathbf{R}_\ell (\rho \mathbf{I} + (\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{D}_\ell \mathbf{R}_\ell)^{-1} \mathbf{R}_\ell \left(\mathbf{D}_\ell^T \mathbf{z}_{\ell-1} + \rho \mathbf{R}_\ell^{-2} \mathbf{z}_\ell + \frac{\mathbf{R}_\ell^{-1} \gamma_\ell}{\sqrt{\mu_\ell}} \right) \quad (3.17)$$

$$\mathbf{R}_\ell^{-1} \mathbf{x}_\ell = (\rho \mathbf{I} + (\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{D}_\ell \mathbf{R}_\ell)^{-1} \left((\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{z}_{\ell-1} + \rho \mathbf{R}_\ell^{-1} \mathbf{z}_\ell + \frac{\gamma_\ell}{\sqrt{\mu_\ell}} \right) \quad (3.18)$$

So, therefore the update equation for $\mathbf{R}_\ell^{-1} \mathbf{x}_\ell$ is the following:

$$\mathbf{R}_\ell^{-1} \mathbf{x}_\ell^{(t+1)} = (\rho \mathbf{I} + (\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{D}_\ell \mathbf{R}_\ell)^{-1} \left((\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{z}_{\ell-1}^{(t)} + \rho \left(\mathbf{R}_\ell^{-1} \mathbf{z}_\ell^{(t)} + \frac{\gamma_\ell^{(t)}}{\rho \sqrt{\mu_\ell}} \right) \right) \quad (3.19)$$

Before moving onto another update equation, there are a few useful things to note here. The form of the inverse matrix identically matches the form from the last chapter, so the inverse representation can be updated efficiently if the updates have a low-rank structure. Furthermore, $\mathbf{D}_\ell \mathbf{R}_\ell$ is the unnormalized dictionary that is updated through low-rank updates. The normalized dictionary does not need to be explicitly calculated at all. It would be easy to isolate \mathbf{x}_ℓ , but it will be simpler to keep track of $\mathbf{R}_\ell^{-1} \mathbf{x}_\ell$ instead, similar to how $\frac{\mathbf{u}}{\rho}$ is tracked instead of \mathbf{u} in the scaled ADMM algorithm.

Proximal Updates

The second set of primal updates comes from equation 3.9, repeated here for convenience:

$$\mathbf{z}^{(t+1)} = \arg \min_{\mathbf{z}} \mathcal{L} \left(\mathbf{x}^{(t+1)}, \mathbf{z}, \gamma^{(t+\frac{1}{2})} \right) \quad (3.20)$$

where, as before,

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\gamma}) = f(\mathbf{x}, \mathbf{z}) + \sum_{\ell=1}^L \frac{\rho}{2} \|\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1}(\mathbf{z}_\ell - \mathbf{x}_\ell) + \frac{\boldsymbol{\gamma}_\ell}{\rho}\|_2^2 - \frac{1}{2\rho} \|\boldsymbol{\gamma}_\ell\|_2^2 \quad (3.21)$$

and

$$f(\mathbf{x}, \mathbf{z}) = \sum_{\ell=1}^L \frac{\mu_\ell}{2} \|\mathbf{z}_{\ell-1} - \mathbf{D}_\ell \mathbf{x}_\ell\|_2^2 + \lambda_\ell \|\mathbf{z}_\ell\|_1 \quad (3.22)$$

For $\mathbf{x}, \mathbf{z}, \boldsymbol{\gamma}$, such that $\nabla_{\mathbf{z}_\ell} \mathcal{L}_\rho(\mathbf{x}_1, \dots, \mathbf{x}_L, \mathbf{z}_0, \dots, \mathbf{z}_L, \boldsymbol{\gamma}_0, \dots, \boldsymbol{\gamma}_L) = 0$:

$$\nabla_{\mathbf{z}} \frac{\mu_{\ell+1}}{2} \|\mathbf{z}_\ell - \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1}\|_2^2 + \lambda_\ell \|\mathbf{z}_\ell\|_1 + \frac{\rho}{2} \|\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1}(\mathbf{z}_\ell - \mathbf{x}_\ell) + \frac{\boldsymbol{\gamma}_\ell}{\rho}\|_2^2 = 0 \quad (3.23)$$

Something important to note here is that each element of \mathbf{z}_ℓ can be treated independently, that is:

$$\nabla_{\mathbf{z}_\ell[i]} \frac{\mu_{\ell+1}}{2} (\mathbf{z}_\ell[i] - (\mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1})[i])^2 + \lambda_\ell |\mathbf{z}_\ell[i]| + \frac{\rho}{2} (\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1}[i](\mathbf{z}_\ell[i] - \mathbf{x}_\ell[i]) + \frac{\boldsymbol{\gamma}_\ell[i]}{\rho})^2 = 0 \quad (3.24)$$

$$\nabla_{\mathbf{z}_\ell[i]} \frac{\mu_{\ell+1}}{2} (\mathbf{z}_\ell[i] - (\mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1})[i])^2 + \lambda_\ell |\mathbf{z}_\ell[i]| + \frac{\rho \mu_\ell}{2 \mathbf{R}_\ell^2[i]} (\mathbf{z}_\ell[i] - \mathbf{x}_\ell[i] + \frac{\mathbf{R}_\ell[i] \boldsymbol{\gamma}_\ell[i]}{\rho \sqrt{\mu_\ell}})^2 = 0 \quad (3.25)$$

For the sake of brevity, I will now drop the indexing:

$$\nabla_{\mathbf{z}_\ell} \frac{\mu_{\ell+1}}{2} (\mathbf{z}_\ell^2 - 2(\mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1}) \mathbf{z}_\ell) + \lambda_\ell |\mathbf{z}_\ell| + \frac{\rho \mu_\ell}{2 \mathbf{R}_\ell^2} (\mathbf{z}_\ell^2 - 2 \mathbf{x}_\ell \mathbf{z}_\ell + \frac{2 \mathbf{R}_\ell \boldsymbol{\gamma}_\ell \mathbf{z}_\ell}{\rho \sqrt{\mu_\ell}}) = 0 \quad (3.26)$$

$$\nabla_{\mathbf{z}_\ell} \frac{1}{2} (\mu_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2}) \mathbf{z}_\ell^2 - \mu_{\ell+1} \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1} \mathbf{z}_\ell - \rho \mu_\ell \mathbf{R}_\ell^{-2} \mathbf{x}_\ell \mathbf{z}_\ell + \sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \gamma_\ell \mathbf{z}_\ell + \lambda_\ell |\mathbf{z}_\ell| = 0 \quad (3.27)$$

$$\nabla_{\mathbf{z}_\ell} \frac{1}{2} \mathbf{z}_\ell^2 - \frac{\mu_{\ell+1} \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2} \mathbf{x}_\ell - \sqrt{\mu_\ell} \mathbf{R}_\ell^{-1} \gamma_\ell}{\mu_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2}} \mathbf{z}_\ell + \frac{\lambda_\ell}{\mu_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2}} |\mathbf{z}_\ell| = 0 \quad (3.28)$$

$$\mathbf{z}_\ell = \mathbf{S} \frac{\lambda_\ell}{\mu_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2}} \left(\frac{\mu_{\ell+1} \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2} (\mathbf{x}_\ell - \frac{\mathbf{R}_\ell \gamma_\ell}{\rho \sqrt{\mu_\ell}})}{\mu_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2}} \right) \quad (3.29)$$

$$\mathbf{z}_\ell = \frac{1}{\mu_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-2}} \mathbf{S}_{\lambda_\ell} \left(\mu_{\ell+1} \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1} + \rho \mu_\ell \mathbf{R}_\ell^{-1} \left(\mathbf{R}_\ell^{-1} \mathbf{x}_\ell - \frac{\gamma_\ell}{\rho \sqrt{\mu_\ell}} \right) \right) \quad (3.30)$$

$$\mathbf{z}_\ell^{(t+1)} = (\rho \mu_\ell \mathbf{I} + \mu_{\ell+1} \mathbf{R}_\ell^2)^{-1} \mathbf{R}_\ell^2 \mathbf{S}_{\lambda_\ell} \left(\mu_{\ell+1} \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1}^{(t+1)} + \rho \mu_\ell \mathbf{R}_\ell^{-1} \left(\mathbf{R}_\ell^{-1} \mathbf{x}_\ell^{(t+1)} - \frac{\gamma_\ell^{(t+\frac{1}{2})}}{\rho \sqrt{\mu_\ell}} \right) \right) \quad (3.31)$$

Note there is a dependence on $\mathbf{R}_{\ell+1}^{-1} \mathbf{x}_{\ell+1}$. The last layer will have to be considered separately. Using the same procedure, the update for \mathbf{z}_L can be derived. Given how similar the derivations are to those used for the other \mathbf{z} layers, I will skip to the result.

$$\mathbf{z}_L = \mathbf{R}_L \mathbf{S}_{\frac{\lambda_L \mathbf{R}_L}{\rho \mu_L}} \left(\mathbf{R}_L^{-1} \mathbf{x}_L - \frac{\gamma_L}{\rho \sqrt{\mu_L}} \right) \quad (3.32)$$

$$\mathbf{z}_L^{(t+1)} = \mathbf{R}_L \mathbf{S}_{\frac{\lambda_L \mathbf{R}_L}{\rho \mu_L}} \left(\mathbf{R}_L^{-1} \mathbf{x}_L^{(t+1)} - \frac{\gamma_L^{(t+\frac{1}{2})}}{\rho \sqrt{\mu_L}} \right) \quad (3.33)$$

Dual Updates

Rather than tracking γ_ℓ or $\frac{\gamma_\ell}{\rho}$ explicitly, it will be easier to track $\frac{\gamma_\ell}{\rho\sqrt{\mu_\ell}}$. The update equations are very straightforward.

$$\frac{\gamma_\ell^{(t+\frac{1}{2})}}{\rho\sqrt{\mu_\ell}} = \frac{\gamma_\ell^{(t)}}{\rho\sqrt{\mu_\ell}} + (\alpha - 1)(\mathbf{R}_\ell^{-1}\mathbf{z}_\ell^{(t)} - \mathbf{R}^{-1}\mathbf{z}_\ell^{(t+1)}) \quad (3.34)$$

$$\frac{\gamma_\ell^{(t+1)}}{\rho\sqrt{\mu_\ell}} = \frac{\gamma_\ell^{(t+\frac{1}{2})}}{\rho\sqrt{\mu_\ell}} + \mathbf{R}_\ell^{-1}\mathbf{z}_\ell^{(t+1)} - \mathbf{R}^{-1}\mathbf{z}_\ell^{(t+1)} \quad (3.35)$$

Summary

In this chapter, I have applied the novel sparse coding algorithm from the previous chapter to a multi-layer dictionary model. If the dictionaries are updated with low-rank updates, the inverse representation necessary for the \mathbf{x} updates in the algorithm can be updated efficiently. This approach offers an alternative to direct proximal methods such as FISTA or mathematically suspect inverse approximations like the tight-frame assumption.

CHAPTER 4

JPEG ARTIFACT REMOVAL

Introduction

Despite the existence of better compression algorithms, use of the JPEG compression algorithm is ubiquitous: it is the most commonly used image compression algorithm. Overzealous JPEG compression can produce visible distortions, and image restoration from these distortions is a challenging problem. There are two aspects of JPEG compression which make the restoration process more challenging than simpler restoration problems like deblurring or removing salt-and-pepper noise: JPEG's block-based approach is not spatially invariant, and the quantization is nonlinear. This chapter describes a novel approach to address the challenges of JPEG image restoration using the ADMM-based convolutional sparse coding for a multi-layer dictionary model.

JPEG Algorithm

The JPEG compression process begins with an RGB image input, and consists of five steps. The first is a color transformation, transitioning from RGB to YUV. Then, the U and V color channels are downsampled. The DCT for each 8×8 block is computed (separately for each channel). The DCT coefficients are then quantized using a quantization matrix determined by a user-chosen JPEG quality factor. Finally, these quantized coefficients are reordered and encoded using a lossless variable length coding process.

The standard reconstruction process reverses the lossless encoding, computes the IDCT of the blocks, upsamples the color channels, and reverses the color transform.

Literature Review

Modelling Compressed JPEG Images

Some researchers have observed convolutional dictionary models struggle with large smooth components of signals, likely due to the fact that shifted versions of smooth filters have high coherence.

For this reason, it is often a good idea to subtract a smoothed version \mathbf{s}_{smth} of the signal, and only apply the dictionary model to the residual $\mathbf{s}_{\text{rough}}$.

$$\mathbf{s}_{\text{clean}} = \mathbf{s}_{\text{smth}} + \mathbf{s}_{\text{rough}} \quad (4.1)$$

$$\mathbf{s}_{\text{rough}} \approx \mathbf{D}_1 \mathbf{x}_1 \quad (4.2)$$

When restoring an image after JPEG compression, the original image $\mathbf{s}_{\text{clean}}$ is not known. Instead, the compressed image \mathbf{s} is observed.

$$\mathbf{s} = \mathbf{QW} \mathbf{s}_{\text{clean}} \quad (4.3)$$

$$\mathbf{s} \approx \mathbf{QW}(\mathbf{s}_{\text{smth}} + \mathbf{D}_1 \mathbf{x}_1) \quad (4.4)$$

where \mathbf{W} maps the signal to 8×8 block frequency coefficients (from the cosine transform), and \mathbf{Q} quantizes them.

A means of estimating \mathbf{s}_{smth} from JPEG-compressed image \mathbf{s} is discussed in the Practical Considerations chapter.

From this idea, I construct the pursuit problem:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \frac{\mu_1}{2} \|\mathbf{s} - \mathbf{QW}(\mathbf{D}_1 \mathbf{x}_1 + \mathbf{s}_{\text{smth}})\|_2^2 + \sum_{\ell=2}^L \frac{\mu_\ell}{2} \|\mathbf{x}_{\ell-1} - \mathbf{D}_\ell \mathbf{x}_\ell\|_2^2 + \sum_{\ell=1}^L \|\mathbf{b}_\ell \cdot \mathbf{x}_\ell\|_1 \\ & \text{subject to } \mathbf{x}_\ell > \mathbf{0} \end{aligned} \tag{4.5}$$

with $\mathbf{b}_\ell \geq \mathbf{0}$.

My approach to solve this problem uses the ADMM algorithm, where $\mathbf{x}_1, \dots, \mathbf{x}_L$ are the first set of primal variables, $\mathbf{v}, \mathbf{z}_1, \dots, \mathbf{z}_L$ are the second set of primal variables, and $\boldsymbol{\eta}$ is the dual variable corresponding to the \mathbf{v} constraint and $\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_L$ are the dual variables corresponding to constraints on $\mathbf{z}_1, \dots, \mathbf{z}_L$. Here is the corresponding optimization problem:

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{v}, \mathbf{z}}{\text{minimize}} \frac{\mu_1}{2} \|\mathbf{v} - \mathbf{D}_1 \mathbf{x}_1 - \mathbf{s}_{\text{smth}}\|_2^2 + \sum_{\ell=2}^L \frac{\mu_\ell}{2} \|\mathbf{z}_{\ell-1} - \mathbf{D}_\ell \mathbf{x}_\ell\|_2^2 + \sum_{\ell=1}^L \|\mathbf{b}_\ell \cdot \mathbf{z}_\ell\|_1 \\ & \text{subject to } \mathbf{z}_\ell \geq \mathbf{0} \end{aligned} \tag{4.6}$$

$$\sqrt{\mu} \mathbf{R}_\ell^{-1} (\mathbf{z}_\ell - \mathbf{x}_\ell) = \mathbf{0}$$

$$\mathbf{QW}(\mathbf{v}) - \mathbf{s} = \mathbf{0}$$

The constraint $\mathbf{QW}(\mathbf{v}) - \mathbf{s} = \mathbf{0}$ is not convex because of the quantization. A linear approximation of the quantization operation (as seen in [12]) would produce an affine constraint. However, I have a different novel means of handling the quantization operator in the constraint, described in the next section. For now, I will focus on the other variable updates.

Setting $\mathbf{z}_0 = \mathbf{v} - \mathbf{s}_{\text{smth}}$, the \mathbf{x} update is identical to the \mathbf{x} update in the last chapter.

$$\mathbf{R}_\ell^{-1} \mathbf{x}_\ell^{(t+1)} = (\rho \mathbf{I} + (\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{D}_\ell \mathbf{R}_\ell)^{-1} \left((\mathbf{D}_\ell \mathbf{R}_\ell)^T \mathbf{z}_{\ell-1}^{(t)} + \rho \left(\mathbf{R}_\ell^{-1} \mathbf{z}_\ell^{(t)} + \frac{\boldsymbol{\gamma}_\ell^{(t)}}{\rho \sqrt{\mu_\ell}} \right) \right) \quad (4.7)$$

For the \mathbf{z} update, it is helpful to introduce a common convex-optimization trick. Consider the following function:

$$\mathbb{1}_{\mathbf{z} \geq \mathbf{0}} = \begin{cases} 0 & \mathbf{z} \geq \mathbf{0} \\ +\infty & \text{otherwise} \end{cases} \quad (4.8)$$

The function is convex, and if it is added to the objective function, it implicitly enforces the constraint $\mathbf{z} \geq \mathbf{0}$.

So,

$$\begin{aligned} \mathcal{L}_\rho(\mathbf{x}, \mathbf{v}, \mathbf{z}, \boldsymbol{\eta}, \boldsymbol{\gamma}) = & f(\mathbf{x}, \mathbf{v}, \mathbf{z}) + \frac{\rho}{2} \|\mathbf{QW}(\mathbf{v}) - \mathbf{s} + \frac{\boldsymbol{\eta}}{\rho}\|_2^2 - \frac{1}{2\rho} \|\boldsymbol{\eta}\|_2^2 + \\ & \sum_{\ell=1}^L \frac{\rho}{2} \|\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1}(\mathbf{z}_\ell - \mathbf{x}_\ell) + \frac{\boldsymbol{\gamma}_\ell}{\rho}\|_2^2 - \frac{1}{2\rho} \|\boldsymbol{\gamma}_\ell\|_2^2 \end{aligned} \quad (4.9)$$

where

$$f(\mathbf{x}, \mathbf{v}, \mathbf{z}) = \frac{\mu_1}{2} \|\mathbf{v} - \mathbf{D}_1 \mathbf{x}_1 - \mathbf{s}_{\text{smth}}\|_2^2 + \sum_{\ell=2}^L \frac{\mu_\ell}{2} \|\mathbf{z}_{\ell-1} - \mathbf{D}_\ell \mathbf{x}_\ell\|_2^2 + \sum_{\ell=1}^L \|\mathbf{b}_\ell \cdot \mathbf{z}_\ell\|_1 + \mathbb{1}_{\mathbf{z} \geq \mathbf{0}} \quad (4.10)$$

This means the \mathbf{z} update satisfies this equation:

$$\mathbf{z}_\ell = \arg \min_{\mathbf{z}} \frac{\mu}{2} \|\mathbf{z} - \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1}\|_2^2 + \|\mathbf{b}_\ell \cdot \mathbf{z}\|_1 + \frac{\rho}{2} \|\sqrt{\mu_\ell} \mathbf{R}_\ell^{-1}(\mathbf{z} - \mathbf{x}_\ell) + \frac{\boldsymbol{\gamma}_\ell}{\rho}\|_2^2 + \mathbb{1}_{\mathbf{z} \geq \mathbf{0}} \quad (4.11)$$

Like before, each element of \mathbf{z} can be treated independently for the \mathbf{z} update.

$$\mathbf{z}_\ell[i] = \arg \min_z \frac{\mu}{2} (z - (\mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1})[i])^2 + \mathbf{b}_\ell[i] |z| + \frac{\rho}{2} \|\sqrt{\mu} \mathbf{R}_\ell^{-1} [i] (z - \mathbf{x}_\ell[i]) + \frac{\gamma_\ell[i]}{\rho}\|_2^2 + \mathbb{1}_{z \geq 0} \quad (4.12)$$

$$\mathbf{z}_{\text{est}} = \arg \min_z \frac{\mu}{2} \|\mathbf{z} - \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1}\|_2^2 + \|\mathbf{b}_\ell \cdot \mathbf{z}\|_1 + \frac{\rho}{2} \|\sqrt{\mu} \mathbf{R}_\ell^{-1} (\mathbf{z} - \mathbf{x}_\ell) + \frac{\gamma_\ell}{\rho}\|_2^2 \quad (4.13)$$

$$\mathbf{z}_\ell[i] = \begin{cases} \mathbf{z}_{\text{est}}[i] & \mathbf{z}_{\text{est}}[i] \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (4.14)$$

A non-negativity constraint converts the shrinkage operator to a biased rectified linear unit.

$$\mathbf{z}_\ell^{(t+1)} = (\mu_{\ell+1} \mathbf{R}_\ell^2 + \rho \mu_\ell \mathbf{I})^{-1} \mathbf{R}_\ell^2 \text{ReLU} \left(\mu_{\ell+1} \mathbf{D}_{\ell+1} \mathbf{x}_{\ell+1}^{(t+1)} + \rho \mu_\ell \mathbf{R}_\ell^{-1} \left(\mathbf{R}_\ell^{-1} \mathbf{x}_\ell^{(t+1)} - \frac{\gamma_\ell^{(t+\frac{1}{2})}}{\rho \sqrt{\mu_\ell}} \right) - \mathbf{b}_\ell \right) \quad (4.15)$$

$$\mathbf{z}_L^{(t+1)} = \mathbf{R}_L \text{RELU} \left(\mathbf{R}_L^{-1} \mathbf{x}_L^{(t+1)} - \frac{\gamma_L^{(t+\frac{1}{2})}}{\rho \sqrt{\mu_L}} - \frac{\mathbf{R}_L \mathbf{b}_L}{\rho \mu_L} \right) \quad (4.16)$$

Handling Quantization

Experiments

Experiment Setup

Results

Conclusion

CHAPTER 5

PRACTICAL CONSIDERATIONS CONCERNING TENSORFLOW

Boundary Handling

Removing Low-Frequency Signal Content

JPEG Artifact Removal

Tensorflow and Keras

Most of the computations for my research rely on TensorFlow version 2.3.1 [22], a Python library for machine learning specializing in building models with differentiable, parameterizable composite functions and learning model parameters using gradient descent or other gradient-based optimization methods. TensorFlow is a common platform for researchers and developers working on artificial neural networks, and there are many tutorials and examples freely available online, so I will not replicate that work here. This chapter section the reader already has some familiarity with TensorFlow and Keras [23] (a high-level library inside TensorFlow). The goal of this section is to provide the reader with the tools and workarounds to be able to replicate my work without resorting to hacking things together with gradient tape and/or TensorFlow-1-style code.

Why Not Use Gradient Tape and TensorFlow-1-Style Code?

Keras offers a high-level environment. Code written in Keras's framework is easier to integrate with other work. Gradient tape is great for hacking something together or debugging, but promotes styles of coding that are less readable, less maintainable, and less portable. Keras also has a lower learning curve than the broader TensorFlow library.

Shared Weights Between Layers

Trainable TensorFlow variables declared outside of any Keras layer will not be automatically added to a Keras model's list of trainable variables. In most cases, this limitation is not a problem; it is intuitive to declare a layer's weights inside that layer. However, sometimes the same variable is needed in multiple distinct layers. To be include a variable in the model's trainable variables, it is sufficient to declare the variable in one layer and pass the variable (or the layer it was initialized in) as an input argument to the `__init__` function of the other layers that share that variable. This will work even if the Keras model does not use the layer that declared the variable.¹

Custom Partial Gradients

TensorFlow offers a well-documented means of replacing TensorFlow's gradient computations of an operation with specified custom gradient computations. However, if the operation involves multiple tensors that are inputs or trainable variables, the standard approach replaces all the gradients with custom gradients. If TensorFlow's gradient computations are sufficient for some tensors but not others, a workaround is necessary. This workaround is best explained by example.

Suppose the operation is the following:

$$z = f(x, y)$$

for which the standard TensorFlow gradient computations of f are desired in respect to x , but the custom gradient computations desired in respect to y are specified in function $g(\nabla_z \mathcal{L})$. This can be rewritten as the following:

¹One could instead declare the variable outside any layers, pass it into the `__init__` functions of all the variables that depend on it, and then manually add the variable to the model's list of trainable variables, but I do not recommend this approach. The resulting code will be less readable and much less maintainable.

```

@tf.custom_gradient
def h(z,y):
    def grad_fun(grad):
        return (tf.identity(grad),g(grad))
    return z,grad_fun
z = f(x,tf.stop_gradient(y))
z = h(z,y)

```

The function h does nothing on the forward pass, but in the backward pass computes the custom gradient in respect to y as intended.

Updating TensorFlow Variables After Applying Gradients

To update TensorFlow Variables after applying gradients, it is necessary to track which variables are affected and what their corresponding update functions are. To accomplish this, I store the update functions in a Python dictionary using variable names as the dictionary keys. This Python dictionary needs to be widely accessible so that layers can add update functions when they are initialized; a simple way to do this is to make the update function Python dictionary a class attribute. The keys need to be unique, but TensorFlow variable names can conflict. It is easy to avoid this problem by checking for conflicts before adding a new update function.

```

class PostProcess:
    update = {}
    def add_update(varName,update_fun):
        assert varName not in PostProcess.update
        PostProcess.update[varName] = update_fun

```

In the standard Keras training paradigm, models are trained using the fit function, a method in the Keras model object. The fit function calls the function `train_step`, where gradients are applied. To update TensorFlow Variables after gradients are applied, `train_step`

is the function to modify. The only change that needs to be made is adding a function call to all update functions that correspond to the model's list of trainable variables.

```
class Model_subclass(tf.keras.Model):  
    def train_step(self, data):  
        trainStepOutputs =  
            tf.keras.Model.train_step(self, data)  
        update_ops = []  
        for tv in self.trainable_variables:  
            if tv.name in PostProcess.update:  
                PostProcess.update[tv.name]()  
        return trainStepOutputs
```

Changes to Tensorflow variables in the update function must use the assign command (or its variants: assign_add, assign_sub, ect). Otherwise, TensorFlow will detect that computations lie outside of its computational graph and throw an error. Note that using the assign command on Python variables that are not TensorFlow variables will produce some very cryptic error messages, so be sure to use the assign command correctly. If the value change of one TensorFlow variable depends on the value of another TensorFlow variable value pre-update, it may be necessary to use the TensorFlow control_dependencies command to get TensorFlow to track that dependency. TensorFlow has a useful tool called TensorBoard that helps visualize TensorFlow's dependencies, but a workaround is required to use TensorBoard on update functions that are called after applying gradients. To use TensorBoard to visualize dependencies in an update function, temporarily call the update function in the layer's call method, use TensorBoard to verify all necessary dependencies are being tracked, then remove the update function call from the layer's call method.

The Perils of Using Built-In Functions for Complex Tensors and Arrays

The TensorFlow Probability version 0.11.1 [24] is an extension of TensorFlow mostly used for probabilistic models. The library contains a Cholesky update function, but the function does not properly handle complex inputs. To compute Cholesky updates for complex inputs, users should either write their own implementation or use my code (included in supplementary material). Similarly, the Randomized SVD algorithm in the Python scikit-learn library does not properly handle complex inputs.

Errors like these are fairly common, so when dealing with complex data, researchers and practitioners should carefully verify that the function libraries they rely on are properly handling complex numbers.

Appendices

APPENDIX A

EXPERIMENTAL EQUIPMENT

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

APPENDIX B

DATA PROCESSING

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

REFERENCES

- [1] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *et al.*, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [2] J. Eckstein, “Parallel alternating direction multiplier decomposition of convex programs,” *Journal of Optimization Theory and Applications*, vol. 80, no. 1, pp. 39–62, 1994.
- [3] R. Nishihara, L. Lessard, B. Recht, A. Packard, and M. Jordan, “A general analysis of the convergence of admm,” in *International Conference on Machine Learning*, PMLR, 2015, pp. 343–352.
- [4] H. Bristow, A. Eriksson, and S. Lucey, “Fast convolutional sparse coding,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 391–398.
- [5] M. Šorel and F. Šroubek, “Fast convolutional sparse coding using matrix inversion lemma,” *Digital Signal Processing*, vol. 55, pp. 44–51, 2016.
- [6] F. Heide, W. Heidrich, and G. Wetzstein, “Fast and flexible convolutional sparse coding,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5135–5143.
- [7] B. Wohlberg, “Efficient algorithms for convolutional sparse representations,” *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 301–315, 2015.
- [8] H. V. Henderson and S. R. Searle, “On deriving the inverse of a sum of matrices,” *Siam Review*, vol. 23, no. 1, pp. 53–60, 1981.
- [9] N. Chodosh, C. Wang, and S. Lucey, “Deep convolutional compressed sensing for lidar depth completion,” in *Asian Conference on Computer Vision*, Springer, 2018, pp. 499–513.
- [10] C. Murdock, M. Chang, and S. Lucey, “Deep component analysis via alternating direction neural networks,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 820–836.
- [11] A. Beck and M. Teboulle, “A fast iterative shrinkage-thresholding algorithm for linear inverse problems,” *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183–202, 2009.

- [12] N. Chodosh and S. Lucey, “When to use convolutional neural networks for inverse problems,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8226–8235.
- [13] Y. Xu and W. Yin, “A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion,” *SIAM Journal on imaging sciences*, vol. 6, no. 3, pp. 1758–1789, 2013.
- [14] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus, “Deconvolutional networks,” in *2010 IEEE Computer Society Conference on computer vision and pattern recognition*, IEEE, 2010, pp. 2528–2535.
- [15] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [16] A. Rangamani, A. Mukherjee, A. Basu, A. Arora, T. Ganapathi, S. Chin, and T. D. Tran, “Sparse coding and autoencoders,” in *2018 IEEE International Symposium on Information Theory (ISIT)*, IEEE, 2018, pp. 36–40.
- [17] V. Pappas, Y. Romano, and M. Elad, “Convolutional neural networks analyzed via convolutional sparse coding,” *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 2887–2938, 2017.
- [18] L. Fei-Fei, R. Fergus, and P. Perona, “Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories,” in *2004 conference on computer vision and pattern recognition workshop*, IEEE, 2004, pp. 178–178.
- [19] B. Chen, G. Polatkan, G. Sapiro, L. Carin, and D. B. Dunson, “The hierarchical beta process for convolutional factor analysis and deep learning,” in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, 2011, pp. 361–368.
- [20] B. Chen, G. Polatkan, G. Sapiro, D. Blei, D. Dunson, and L. Carin, “Deep learning with hierarchical convolutional factor analysis,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1887–1901, 2013.
- [21] Y. Pu, X. Yuan, and L. Carin, “Generative deep deconvolutional learning,” *ArXiv preprint arXiv:1412.6039*, 2014.
- [22] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden,

M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, *Tensorflow: Large-scale machine learning on heterogeneous systems*, Software available from tensorflow.org, 2015.

[23] F. Chollet *et al.*, *Keras*, <https://keras.io>, 2015.

[24] J. V. Dillon, I. Langmore, D. Tran, E. Brevdo, S. Vasudevan, D. Moore, B. Patton, A. Alemi, M. Hoffman, and R. A. Saurous, “Tensorflow distributions,” *ArXiv preprint arXiv:1711.10604*, 2017.

VITA

Vita may be provided by doctoral students only. The length of the vita is preferably one page. It may include the place of birth and should be written in third person. This vita is similar to the author biography found on book jackets.