

Chapter 8

Convex Optimization for Structured Signal Recovery

In the previous part, we showed rigorously that under fairly broad conditions, many important classes of structured signals can be recovered via computationally trackable optimization problems, such as ℓ^1 minimization for recovering sparse signals and nuclear norm minimization for recovering low-rank matrices. As we will see in Part III of this book, many of these structures are essential for modeling many high-dimensional data that arise in real-world applications. Hence, from a practical perspective, it is important that we develop efficient and scalable algorithms for these classes of optimization problems. We take on this task in this and the next couple of chapters.

In this chapter, we will mainly focus on the *convex approach* for structured signal recovery. There are two compelling reasons why we are in favor of studying the convex approach first. Firstly, the previous chapters have established very precise conditions under which convex programs give exactly correct solutions to the recovery problems, whereas conditions for other approaches are less strong or less understood. Secondly, as we will see through this chapter, the class of convex programs we are dealing with have very interesting and unique properties that lend themselves to much more fast and scalable solutions than generic convex programs. Although we will primarily use ℓ^1 norm or nuclear norm minimization as working examples, the techniques that we introduce are rather general. They can be easily applied to much larger class of convex programs that share the same properties and structures.

In the next two chapters, we will study two alternative approaches to develop algorithms for structured signal recovery. One approach is to rely

on simple and straightforward greedy strategies. This class of techniques are important because in many practical applications one may need to trade the range of working conditions of algorithms for better scalability or simplicity for implementation. Another approach is to explore stronger, not necessarily convex, regularizers for the desired structures such as sparsity. In practice, these regularizers are often observed to enjoy a wider range of working conditions than convex regularizers.

This chapter (or book) is not intended to give a comprehensive introduction to (convex) optimization¹. Instead, the chapters will focus mainly on showing how one can exploit special structures of the problems so as to develop significantly more efficient and scalable algorithms than generic convex optimization methods. To make the book better self-contained, we will briefly survey related concepts and properties of convex functions as well as generic optimization methods in Appendices B-D.

8.1 Challenges and Opportunities

We have shown that under fairly broad conditions, many important classes of structured signals can be recovered by convex optimization – in particular, we analyzed in depth the ℓ^1 heuristic for recovering sparse signals and the nuclear norm heuristic for recovering low-rank matrices. As we will see later in the book, many of these structures are essential for modeling many high-dimensional data that arise in real-world applications. So, from a practical perspective, it is important to have efficient, scalable algorithms for these classes of optimization problems.

In this chapter, we will describe a few basic ideas which go a long way towards achieving this, by exploiting special properties of the particular convex optimization problems that arise in structured signal recovery. Our discussion will center around four model problems: basis pursuit (i.e., equality constrained ℓ^1 minimization), its penalized version (basis pursuit denoising), robust PCA, and its penalized version. We recap these four optimization problems below:

Recall the problem of recovering a sparse vector $\mathbf{x}_o \in \mathbb{R}^n$ from observations $\mathbf{y} = \mathbf{A}\mathbf{x}_o \in \mathbb{R}^m$ via convex program, also known as *basis pursuit* (BP):

$$\begin{aligned} & \text{minimize} && \|\mathbf{x}\|_1 \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y}. \end{aligned} \tag{8.1.1}$$

A variant of the problem considers the noisy case, in which the observations \mathbf{y} are contaminated by moderate Gaussian noise $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{z}$, also known

¹for which there are already excellent references such as [Boyd and Vandenberghe, 2004].

as the *Lasso*:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (8.1.2)$$

where λ is a scalar weight parameter.

In robust PCA, the goal is to recover a low-rank matrix \mathbf{L}_o from sparsely corrupted observations $\mathbf{Y} = \mathbf{L}_o + \mathbf{S}_o$. A natural approach suggested in earlier chapters is to solve the so called *principal component pursuit* (PCP) program:

$$\begin{aligned} & \text{minimize} && \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \\ & \text{subject to} && \mathbf{L} + \mathbf{S} = \mathbf{Y}, \end{aligned} \quad (8.1.3)$$

where $\lambda > 0$ is a scalar weight. Again, if the data are noisy we could also consider to solve a stable version of the PCP program:

$$\text{minimize} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2 + \lambda_1 \|\mathbf{L}\|_* + \lambda_2 \|\mathbf{S}\|_1, \quad (8.1.4)$$

to produce stable estimates $\hat{\mathbf{L}}$ and $\hat{\mathbf{S}}$.

Challenge of scale.

When the dimension of the problem is not so high, one could simply apply classical second-order convex optimization algorithms, such as the interior-point methods (see e.g. [Boyd and Vandenberghe, 2004]), to solve the above convex programs. These require very few iterations to converge to a highly accurate solution: $O(\log \log(1/\varepsilon))$, where ε is the target accuracy. However, for problems in n variables, each iteration requires the solution of a linear system of equations of size $n \times n$, incurring a per-iteration cost of $O(n^3)$. For applications in modern signal processing, the number of variables can be quite high – in the case of image processing, it is typically of the same magnitude as the number of pixels, easily in the range of millions. For such problems, the cost of a single iteration is prohibitively large. We will need to consider simple alternatives with cheaper iterations.

Example 8.1.1 (Solving large-scale BP problems via interior-point methods). *Just to motivate yourself from practical perspective, construct a simulation to plot the average run time of BP on CVX from 100 to 1000 dimensions. See how a generic convex optimization solver (that is mainly based on second-order, interior point method) scales for this class of optimization problems.*

Difficulty with non-smoothness.

The large scale of the problems forces us to consider simple, scalable algorithms that use only first-order information about the objective function. The prototypical first-order method is *gradient descent*. However, one technical difficulty arises though as the objective functions may contain non-smooth terms that are not differentiable. For instance, the ℓ^1 norm

$\|\mathbf{x}\|_1$ in the BPDN problem does not have a gradient in the normal sense. In such cases, the simplest solution is to employ a generic subgradient method, as we did in Chapter 2. Although the subgradient method has very simple and efficient iterations, its rate of convergence is very poor, typically $O(1/\sqrt{k})$.² That means it usually takes many (thousands of) iterations for the algorithm to converge to the optimal solution. In this chapter, we will show how to exploit some important properties of structured signal recovery. Such properties allow us to develop gradient descent algorithms as if the objective is smooth, the so called Proximal Gradient (PG) method. The same properties also allow us to utilize acceleration techniques that were designed for smooth functions and lead to much more scalable and fast-converging algorithms, with convergence rates much better than the generic situation.

Enforcing equality constraints.

To solve the basis pursuit problem (8.1.1), we need to ensure that the final solution \mathbf{x} satisfies the equality constraint $\mathbf{y} = \mathbf{A}\mathbf{x}$ exactly. A naive way to enforce the equality constraint is to incorporate it as a penalty term and minimize: $\min_{\mathbf{x}} \|\mathbf{x}\|_1 + \frac{\mu}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$. This is a similar optimization problem as the BPDN only that we need to solve a series of problems of this type for an increasing sequence of $\mu_i \rightarrow \infty$ so that the equality constraint is to be enforced in the end. However, as the weight μ_i increases, the corresponding BPDN problem becomes increasingly ill-conditioned and hence algorithms converge slower. In this chapter, we will see how to employ the *augmented Lagrangian multiplier* (ALM) technique to alleviate this difficulty.

Exploiting separable structures.

Very often the structured signal that we are recovering is a superposition of multiple structured terms. This is the case for the principal component pursuit (PCP) program that we have mentioned earlier:

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \quad \text{subj. to} \quad \mathbf{Y} = \mathbf{L} + \mathbf{S}. \quad (8.1.5)$$

As we will show in this chapter, such separable structures of the objective function can be naturally exploited through methods such as *alternating direction of multipliers method* (ADMM). We end up with more simple and efficient algorithms for solving this class of convex programs as ADMM converts the global optimization to several parallel subproblems of much smaller dimension.

²See [Nemirovski, 1995, Nemirovski, 2007] for a characterization of typical subgradient methods.

8.2 Proximal Gradient Methods

As one can see, the optimization problems we are dealing with have a common structure: they all try to minimize convex objective functions of the form:

$$F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x}), \quad (8.2.1)$$

where $f(\mathbf{x})$ is a smooth convex term and $g(\mathbf{x})$ is convex but non-smooth term. For instance, in the Lasso problem (8.1.2), we could set $f(\mathbf{x}) \doteq \frac{1}{2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ and $g(\mathbf{x}) \doteq \lambda\|\mathbf{x}\|_1$. We want to develop both scalable and fast-converging algorithms for this type of problems.

Since the composite objective function $F(\mathbf{x})$ is not differentiable, general gradient descent algorithms do not apply. The first recourse in this situation is to replace the gradient with a *subgradient*, yielding the simple subgradient descent iteration

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \mathbf{g}_k, \quad \mathbf{g}_k \in \partial F(\mathbf{x}_k). \quad (8.2.2)$$

The main disadvantage to this approach is its relatively poor convergence rate. In general, the convergence rate of subgradient descent for non-smooth objective functions is (see [Nesterov, 2003b]):

$$F(\mathbf{x}_k) - F(\mathbf{x}_\star) = O(1/\sqrt{k}). \quad (8.2.3)$$

The constants in the big- O notation depend on various properties of the problem. The important point is that for even a moderate target accuracy ε , we will have to set $k = \Omega(\varepsilon^{-2})$ very large.

Interlude: convergence rates of first order methods.

We can contrast the behavior of the subgradient method with the behavior of the simple *gradient descent* method for minimizing a smooth function. Consider, briefly, the simpler problem

$$\text{minimize } f(\mathbf{x}), \quad (8.2.4)$$

with f convex and differentiable. The gradient descent iteration for this problem is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \nabla f(\mathbf{x}_k). \quad (8.2.5)$$

This iteration comes from a first-order approximation to f at $\mathbf{x} = \mathbf{x}_k$:

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle. \quad (8.2.6)$$

Because f is convex, this first order approximation provides a global lower bound on f . Nevertheless, we expect this lower bound to be more accurate in a neighborhood of \mathbf{x} . The size of this neighborhood depends substantially on the properties of f . For example, if f is relatively smooth, and its gradient does not vary much from point to point, we might imagine that the first order approximation at \mathbf{x} would be accurate over a relatively large

region. To make this more formal, we say that a differentiable function $f(\mathbf{x})$ has L -Lipschitz gradient if

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|, \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n \quad (8.2.7)$$

for some $L > 0$. The quantity L is known as the *Lipschitz constant* of ∇f . When this condition holds, a bit of calculus shows that we can complement the linear lower bound (8.2.6) with a corresponding quadratic upper bound:

Lemma 8.2.1. *Suppose that f is differentiable, and ∇f is L -Lipschitz. Then for every \mathbf{x}, \mathbf{x}' ,*

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2. \quad (8.2.8)$$

Proof. We calculate:

$$f(\mathbf{x}') = f(\mathbf{x} + t(\mathbf{x}' - \mathbf{x}))|_{t=1} \quad (8.2.9)$$

$$= f(\mathbf{x}) + \int_{t=0}^1 \frac{d}{dt} f(\mathbf{x} + t(\mathbf{x}' - \mathbf{x})) dt \quad (8.2.10)$$

$$= f(\mathbf{x}) + \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{x}' - \mathbf{x})), \mathbf{x}' - \mathbf{x} \rangle dt \quad (8.2.11)$$

$$= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{x}' - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle dt \quad (8.2.12)$$

$$\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \int_{t=0}^1 tL \|\mathbf{x}' - \mathbf{x}\|_2^2 dt \quad (8.2.13)$$

$$= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2, \quad (8.2.14)$$

giving the claim. \square

Thus, when ∇f is Lipschitz, we have a matching quadratic upper bound

$$f(\mathbf{x}') \leq \hat{f}(\mathbf{x}', \mathbf{x}) \doteq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2, \quad (8.2.15)$$

and that the upper bound agrees with f at the point \mathbf{x} at which it is formed: $f(\mathbf{x}) = \hat{f}(\mathbf{x}, \mathbf{x})$. Suppose that we minimize this upper bound, with respect to \mathbf{x}' . A quick calculation shows that the minimizer has a very familiar form:

$$\arg \min_{\mathbf{x}'} \hat{f}(\mathbf{x}', \mathbf{x}) = \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}). \quad (8.2.16)$$

This is simply a gradient descent step, taken from \mathbf{x} , with a special choice of step size $\gamma = 1/L$. Moreover, because $\hat{f}(\mathbf{x}, \mathbf{x}) = f(\mathbf{x})$, this minimization does not increase the objective function: if $\mathbf{x}'_\star \in \arg \min_{\mathbf{x}'} \hat{f}(\mathbf{x}', \mathbf{x})$, then

$$f(\mathbf{x}'_\star) \leq \hat{f}(\mathbf{x}'_\star, \mathbf{x}) \leq \hat{f}(\mathbf{x}, \mathbf{x}) = f(\mathbf{x}). \quad (8.2.17)$$

Thus, if we apply the gradient descent method with step size $1/L$, we are guaranteed to produce a monotone sequence of function values $f(\mathbf{x}_k)$. Furthermore, we can show convergence³ to the optimal function value at a rate of $O(1/k)$:

$$f(\mathbf{x}_k) - f(\mathbf{x}_\star) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{2k} = O(1/k). \quad (8.2.18)$$

This is still not a particularly fast rate of convergence, but it is much better than the $O(1/\sqrt{k})$ rate of convergence experienced by the subgradient algorithm on nonsmooth functions.

From gradient to proximal gradient.

Can we draw inspiration from the gradient method to produce a more efficient algorithm for minimizing the composite function $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$, with f differentiable? Again, the gradient method does not directly apply, since F is nondifferentiable. Nevertheless, if the gradient ∇f of the smooth term is Lipschitz, we can still make a simpler upper bound to F , by upper bounding f by a quadratic and leaving the nonsmooth term g untouched:

$$\hat{F}(\mathbf{x}', \mathbf{x}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2 + g(\mathbf{x}'). \quad (8.2.19)$$

Since above repeatedly minimizing \hat{f} produced the gradient method, resulting in a better convergence rate, let us try minimizing \hat{F} :

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \hat{F}(\mathbf{x}, \mathbf{x}_k). \quad (8.2.20)$$

For commonly encountered g , this minimization often takes on a very simple form. Completing the square in (8.2.19), we obtain that

$$\hat{F}(\mathbf{x}, \mathbf{x}_k) = \frac{L}{2} \left\| \mathbf{x} - \left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \right) \right\|_2^2 + g(\mathbf{x}) + h(\mathbf{x}_k), \quad (8.2.21)$$

where $h(\mathbf{x}_k)$ is a term that depends only on \mathbf{x}_k .

Hence, the iteration (8.2.20) becomes

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \frac{L}{2} \left\| \mathbf{x} - \left(\mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \right) \right\|_2^2 + g(\mathbf{x}) \quad (8.2.22)$$

$$= \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{w}_k\|_2^2, \quad (8.2.23)$$

where for convenience we define $\mathbf{w}_k = \mathbf{x}_k - (1/L)\nabla f(\mathbf{x}_k)$. Thus, at each step of the iteration (8.2.20), we have to minimize the function g plus a separable quadratic $\frac{L}{2} \|\mathbf{x} - \mathbf{w}\|_2^2$. In a sense, this problem asks us to make g as small as possible, while not straying too far from the point \mathbf{w} . Because

³We will not prove (8.2.18) here, although we will obtain a more general result below which implies it.

$\|\cdot\|_2^2$ is strongly convex, this problem always has a solution, and the solution is unique. So, the sequence \mathbf{x}_k defined recursively by (8.2.20) is well-defined.

In fact, the operation of minimizing a convex function g plus a separable quadratic $\|\mathbf{x} - \mathbf{w}\|_2^2$ recurs so frequently in convex analysis and optimization that it has its own name. This is known as the *proximal operator* for the convex function $g(\mathbf{x})$:

Definition 8.2.2 (Proximal Operator). *The proximal operator of a convex function g is*

$$\text{prox}_g(\mathbf{w}) \doteq \arg \min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 \right\}. \quad (8.2.24)$$

In this language, iteration (8.2.20) can be written as

$$\mathbf{x}_{k+1} = \text{prox}_{g/L}(\mathbf{w}_k). \quad (8.2.25)$$

Fortunately, many of the convex functions (or norms) that we encounter in structured signal recovery either have close-form proximal operators or proximal operators that can be computed very efficiently via numerical means. We give a few examples below:

Proposition 8.2.3. *Proximal operators for the indicator function, ℓ^1 norm, and nuclear norm are given by:*

1. Let $g(\mathbf{x}) = I_{\mathcal{D}}$ be the indicator function for a closed convex set \mathcal{D} , namely, $I_{\mathcal{D}}(\mathbf{x}) = 0, \mathbf{x} \in \mathcal{D}$ otherwise $I_{\mathcal{D}}(\mathbf{x}) = \infty$. Then $\text{prox}_g(\mathbf{w})$ is the projection operator:

$$\text{prox}_g(\mathbf{w}) = \arg \min_{\mathbf{x} \in \mathcal{D}} \|\mathbf{x} - \mathbf{w}\|_2^2 = \mathcal{P}_{\mathcal{D}}[\mathbf{w}].$$

2. Let $g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$ be the ℓ^1 norm. Then $\text{prox}_g(\mathbf{w})$ is the soft-thresholding function applied element-wise:

$$\text{prox}_g(u_i) = \text{soft}(u_i, \lambda) \doteq \text{sign}(u_i) \max(|u_i| - \lambda, 0).$$

3. Let $g(\mathbf{x}) = \lambda \|\mathbf{X}\|_*$ be the matrix nuclear norm. Then $\text{prox}_g(\mathbf{W})$ is the singular-value thresholding function:

$$\text{prox}_g(\mathbf{W}) = \mathcal{S}(\mathbf{W}, \lambda) \doteq \mathbf{U} \text{soft}(\mathbf{\Sigma}, \lambda) \mathbf{V}^*,$$

where $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V})$ are the singular value decomposition (SVD) of \mathbf{W} . In other words, $\mathcal{S}(\mathbf{W}, \lambda)$ applies soft thresholding operator on the singular values of \mathbf{X} .

Proof. We prove the second assertion. The objective function reaches minimum when the subdifferential of $\lambda \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2$ contains zero,

$$0 \in (\mathbf{x} - \mathbf{w}) + \lambda \partial \|\mathbf{x}\|_1 = \begin{cases} x_i - w_i + \lambda, & x_i > 0 \\ -w_i + \lambda[-1, 1], & x_i = 0 \\ x_i - w_i - \lambda, & x_i < 0 \end{cases}, \quad i = 1, \dots, n.$$

Therefore, the solution to this equality constraint is the soft-thresholding function applied element-wise:

$$x_{i*} = \text{soft}(w_i, \lambda) \doteq \text{sign}(w_i) \max(|w_i| - \lambda, 0), \quad i = 1, \dots, n.$$

We leave the first and third assertions as exercises to the reader. [Hint: for the first, use the definition; for the third, use the subdifferential of $\|\cdot\|_*$.] \square

Thus, for problems of our interest, we can compute the proximal operator efficiently. In this setting, it provides an alternative replacement for the gradient step. Unlike the subgradient method, this *proximal gradient* algorithm enjoys a convergence rate of $O(1/k)$ – exactly the same as if the nonsmooth term was not present! More formally:

Theorem 8.2.4. *Let $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$, with f be a convex, differentiable function with L -Lipschitz continuous gradient, and g a convex function. Consider the following iterative update scheme:*

$$\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k), \quad \mathbf{x}_{k+1} \leftarrow \text{prox}_{g/L}(\mathbf{w}_k).$$

Assume $F(\mathbf{x})$ has a minimum at \mathbf{x}_* . Then for any $k \geq 1$,

$$F(\mathbf{x}_k) - F(\mathbf{x}_*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{2k}.$$

We will give a detailed proof to this theorem in Section 8.2.3. Thus, for a composite non-smooth convex function, under certain conditions, we can still obtain an efficient “gradient descent” like algorithm that has the same convergence rate $O(1/k)$ as that for a smooth function. As long as the non-smooth part has an easy-to-solve proximal operator, the proximal gradient algorithm has very cheap iterations hence it is typically much more scalable than second-order methods. We summarize properties of the iterative process that we have derived so far for minimizing the convex composite problem in Figure 8.1, which is also known as the *proximal gradient* method.

8.2.1 Proximal Gradient for the Lasso

For the rest of the section, we will see how to apply the proximal gradient algorithm to several important cases of structured signal recovery problems that we have encountered.

As the first instance, the Lasso problem (8.1.2) obviously falls into the class of problems that can be addressed by the proximal gradient method. We can view g to be the ℓ^1 norm function $\lambda \|\mathbf{x}\|_1$ whose proximal operator is given in Proposition 8.2.3; f is simply be the quadratic data term $\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ whose gradient is clearly Lipschitz: the Lipschitz constant L can be the largest eigenvalue of the matrix $\mathbf{A}^* \mathbf{A}$, which can be computed in advance.

Proximal Gradient (PG)**Problem Class:**

$$\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$$

$f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex, ∇f L -Lipschitz and g non-smooth.

Basic Iteration: set $\mathbf{x}_0 \in \mathbb{R}^n$.

Repeat:

$$\begin{aligned}\mathbf{w}_k &\leftarrow \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k), \\ \mathbf{x}_{k+1} &\leftarrow \text{prox}_{g/L}(\mathbf{w}_k).\end{aligned}$$

Convergence Guarantee:

$F(\mathbf{x}_k) - F(\mathbf{x}_*)$ converges at a rate of $O(1/k)$.

Figure 8.1. Properties of the Proximal Gradient Method.

The resulting proximal gradient descent algorithm for Lasso is sometimes referred to as the *iterative soft-thresholding algorithm* (ISTA), which is summarized in Algorithm 8.1. In terms of computational complexity, the

Algorithm 8.1 Proximal Gradient (PG) for Lasso

- 1: **Problem:** $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$, given $\mathbf{y} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times n}$.
- 2: **Input:** $\mathbf{x}_0 \in \mathbb{R}^n$ and $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$.
- 3: **while** \mathbf{x}_k not converged ($k = 1, 2, \dots$) **do**
- 4: $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y})$.
- 5: $\mathbf{x}_{k+1} \leftarrow \text{soft}(\mathbf{w}_k, \lambda/L)$.
- 6: **end while**
- 7: **Output:** $\mathbf{x}_* \leftarrow \mathbf{x}_k$.

main cost is calculating the gradient $\nabla f(\mathbf{x}) = \mathbf{A}^* \mathbf{A}\mathbf{x} - \mathbf{A}^* \mathbf{y}$ in the inner loop, which is in the order of $O(mn)$ in general.

Example 8.2.5. We randomly generate a sparse signal \mathbf{x} in an 1000- D space and then add a small Gaussian noise \mathbf{n} to all of its coefficients, as shown in Figure 8.2 Top. With the added Gaussian noise, the signal $\mathbf{w} = \mathbf{x} + \mathbf{n}$ is not sparse anymore. Then we may try to recover \mathbf{x} from \mathbf{w} by solving the following problem $\min_{\mathbf{x}} \lambda \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{w} - \mathbf{x}\|_2^2$, where λ is proportional to the noise level. We know the solution to this problem is simply the soft thresholding $\hat{\mathbf{x}} = \text{soft}(\mathbf{w}, \lambda)$. The result is shown in Figure

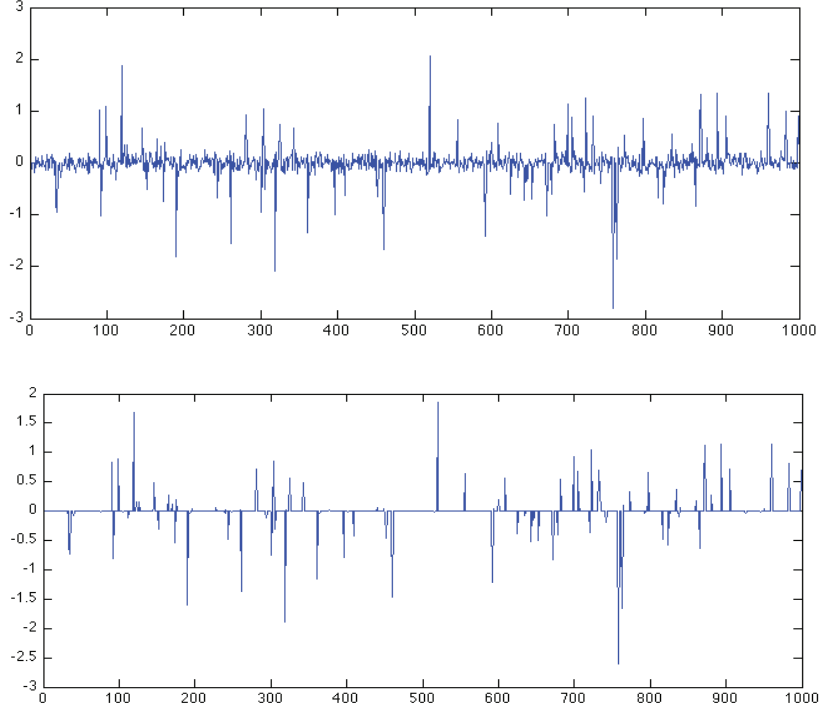


Figure 8.2. **Top:** A sparse signal \mathbf{x} perturbed by small Gaussian noise \mathbf{n} . **Bottom:** The output from a properly chosen soft-thresholding function $\text{soft}(\mathbf{w}, \lambda)$.

8.2 Bottom. We see that the operator successfully remove most of the noise in \mathbf{w} and returns a sparse estimate for \mathbf{x} .

8.2.2 Proximal Gradient for Stable PCP

According to Proposition 8.2.3, the nuclear norm $\|\mathbf{X}\|_*$ also has a simple proximal operator. Hence we could apply proximal gradient algorithm to solve low-rank matrix recovery problems. For instance, the stable principal component pursuit (PCP) program is also for the form that is amenable to proximal gradient method:

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2. \quad (8.2.26)$$

Notice that however, for this problem the non-smooth term $g(\mathbf{L}, \mathbf{S}) = \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1$ now contains two non-smooth functions $\|\mathbf{L}\|_*$ and $\lambda \|\mathbf{S}\|_1$, each having a simple proximal operator.

We leave as an exercise for the reader to prove the following simple fact about the proximal operator of a separable convex function, which comes handy for this problem: Let $\mathbf{x} = [\mathbf{x}_1; \mathbf{x}_2]$ and $g(\mathbf{x}) = g_1(\mathbf{x}_1) + g_2(\mathbf{x}_2)$ be a separable function. Then

$$\text{prox}_g(\mathbf{w}) = (\text{prox}_{g_1}(\mathbf{w}_1), \text{prox}_{g_2}(\mathbf{w}_2)),$$

where \mathbf{w}_1 and \mathbf{w}_2 in $\mathbf{w} = [\mathbf{w}_1; \mathbf{w}_2]$ correspond to the variables \mathbf{x}_1 and \mathbf{x}_2 in \mathbf{x} , respectively.

Hence, the proximal operator for $g(\mathbf{L}, \mathbf{S})$ can be computed separately from the proximal operators for the ℓ^1 norm for \mathbf{S} and nuclear norm for \mathbf{L} , respectively. The rest of the proximal gradient algorithm for the stable principal component pursuit program then is easy to derive (and we leave this as an exercise to the reader. See Exercise 8.2.). We summarize the overall algorithm below for clarity.

Algorithm 8.2 Proximal Gradient (PG) for Stable Principal Component Pursuit

- 1: **Problem:** $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2$, given \mathbf{Y} , $\lambda, \mu > 0$.
 - 2: **Input:** $\mathbf{L}_0 \in \mathbb{R}^{m \times n}$, $\mathbf{S}_0 \in \mathbb{R}^{m \times n}$.
 - 3: **while** $((\mathbf{L}_k, \mathbf{S}_k)$ not converged ($k = 1, 2, \dots$) **do**
 - 4: $\mathbf{W}_k \leftarrow \mathbf{Y} - \mathbf{S}_k$ and compute $\mathbf{W}_k = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$.
 - 5: $\mathbf{L}_{k+1} \leftarrow \mathbf{U}\text{soft}(\mathbf{\Sigma}, 1/\mu)\mathbf{V}^*$.
 - 6: $\mathbf{S}_{k+1} \leftarrow \text{soft}((\mathbf{Y} - \mathbf{L}_k), \lambda/\mu)$.
 - 7: **end while**
 - 8: **Output:** $\mathbf{L}_\star \leftarrow \mathbf{L}_k; \mathbf{S}_\star \leftarrow \mathbf{S}_k$.
-

8.2.3 Convergence of Proximal Gradient

In this subsection, we prove Theorem 8.2.4. We find it convenient to do this in two steps. In the first step, we proved an analysis of a simpler algorithm, known as the *proximal point algorithm*, which only consists of repeated application of the proximal operator. This algorithm is of independent interest, and we will reuse its analysis when we encounter Augmented Lagrangian techniques.

Proposition 8.2.6 (Proximal point algorithm). *Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, and \mathbf{x}_\star a minimizer of g . Let $\mathbf{x}_0 \in \mathbb{R}^n$ be arbitrary, and consider the iteration*

$$\mathbf{x}_{k+1} = \text{prox}_{\gamma_k g}(\mathbf{x}_k), \quad (8.2.27)$$

with $\gamma_k \in \mathbb{R}_+$. Then

$$g(\mathbf{x}_k) - g(\mathbf{x}_\star) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{2 \sum_{i=1}^k \gamma_i}. \quad (8.2.28)$$

Moreover, if $\sum_{i=1}^\infty \gamma_i = +\infty$, then $\mathbf{x}_k \rightarrow \bar{\mathbf{x}}$, a minimizer of g .

Proof. By construction,

$$\mathbf{0} \in \gamma_k \partial g(\mathbf{x}_{k+1}) + \mathbf{x}_{k+1} - \mathbf{x}_k. \quad (8.2.29)$$

Equivalently, $\mathbf{x}_k - \mathbf{x}_{k+1} \in \gamma_k \partial g(\mathbf{x}_{k+1})$. Using the subgradient inequality, we have that

$$\langle \mathbf{x}_k - \mathbf{x}_{k+1}, \mathbf{x}_k - \mathbf{x}_{k+1} \rangle \leq \gamma_k (g(\mathbf{x}_k) - g(\mathbf{x}_{k+1})). \quad (8.2.30)$$

Since left hand side is nonnegative, $g(\mathbf{x}_k) \geq g(\mathbf{x}_{k+1})$. So, the objective function is nondecreasing.

Using the subgradient inequality again, we have that

$$\langle \mathbf{x}_\star - \mathbf{x}_{k+1}, \mathbf{x}_k - \mathbf{x}_{k+1} \rangle \leq \gamma_k (g(\mathbf{x}_\star) - g(\mathbf{x}_{k+1})). \quad (8.2.31)$$

Let us use this fact to bound the distance of \mathbf{x}_{k+1} to the optimum. Notice that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_\star\|_2^2 &= \|\mathbf{x}_k - \mathbf{x}_\star + \mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_\star\|_2^2 + 2 \langle \mathbf{x}_k - \mathbf{x}_\star, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_\star\|_2^2 - \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + 2 \langle \mathbf{x}_{k+1} - \mathbf{x}_\star, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle \\ &\leq \|\mathbf{x}_k - \mathbf{x}_\star\|_2^2 + 2 \langle \mathbf{x}_{k+1} - \mathbf{x}_\star, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle \\ &\leq \|\mathbf{x}_k - \mathbf{x}_\star\|_2^2 + 2\gamma_k (g(\mathbf{x}_\star) - g(\mathbf{x}_{k+1})). \end{aligned} \quad (8.2.32)$$

Since $g(\mathbf{x}_{k+1}) \geq g(\mathbf{x}_\star)$, the distance of \mathbf{x}_k to \mathbf{x}_\star also does not increase. In fact, we can say slightly more. Summing the relationship (8.2.32), we obtain

$$\sum_{i=1}^k 2\gamma_i (g(\mathbf{x}_{i+1}) - g(\mathbf{x}_\star)) \leq \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2. \quad (8.2.33)$$

Since $g(\mathbf{x}_i)$ is nonincreasing, this implies that

$$2 \left(\sum_{i=1}^k \gamma_i \right) (g(\mathbf{x}_{k+1}) - g(\mathbf{x}_\star)) \leq \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2. \quad (8.2.34)$$

This gives convergence in function values, as in (8.2.28).

Since $\|\mathbf{x}_k - \mathbf{x}_\star\|_2$ is nonincreasing, the sequence \mathbf{x}_k is bounded, and hence has a cluster point $\bar{\mathbf{x}}$. Since $g(\mathbf{x}_k) \searrow g(\mathbf{x}_\star)$, $g(\bar{\mathbf{x}}) = g(\mathbf{x}_\star)$, and hence $\bar{\mathbf{x}}$ is optimal. Applying inequality (8.2.32) with \mathbf{x}_\star replaced by $\bar{\mathbf{x}}$, we obtain that $\|\mathbf{x}_k - \bar{\mathbf{x}}\|_2$ is also nonincreasing, whence the cluster point $\bar{\mathbf{x}}$ is the limit of the sequence $\{\mathbf{x}_k\}$. \square

The key idea in the above proof is to use the optimality condition (8.2.29) for the proximal operator, together with the subgradient inequality to relate the suboptimality $g(\mathbf{x}_{k+1}) - g(\mathbf{x}_*)$ in objective function to the distance to the feasible set. To prove Theorem 8.2.4, we follow a very similar program.

Proof. Notice that by construction, there exists $\gamma \in \partial g(\mathbf{x}_{k+1})$ such that

$$\nabla f(\mathbf{x}_k) + L(\mathbf{x}_{k+1} - \mathbf{x}_k) + \gamma = \mathbf{0}. \quad (8.2.35)$$

The subgradient and gradient inequalities for the convex functions f and g give that for any \mathbf{x} ,

$$f(\mathbf{x}) \geq f(\mathbf{x}_k) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle, \quad (8.2.36)$$

$$g(\mathbf{x}) \geq g(\mathbf{x}_{k+1}) + \langle \mathbf{x} - \mathbf{x}_{k+1}, \gamma \rangle, \quad (8.2.37)$$

whence

$$F(\mathbf{x}) \geq f(\mathbf{x}_k) + g(\mathbf{x}_{k+1}) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle + \langle \mathbf{x} - \mathbf{x}_{k+1}, \gamma \rangle. \quad (8.2.38)$$

So,

$$\begin{aligned} F(\mathbf{x}) - F(\mathbf{x}_{k+1}) &\geq F(\mathbf{x}) - \hat{F}(\mathbf{x}_{k+1}, \mathbf{x}_k) \\ &\geq f(\mathbf{x}_k) + g(\mathbf{x}_{k+1}) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle \\ &\quad + \langle \mathbf{x} - \mathbf{x}_{k+1}, \gamma \rangle - \hat{F}(\mathbf{x}_{k+1}, \mathbf{x}_k) \\ &= -\frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + \langle \mathbf{x} - \mathbf{x}_{k+1}, \nabla f(\mathbf{x}_k) + \gamma \rangle \\ &= -\frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + L \langle \mathbf{x} - \mathbf{x}_{k+1}, \mathbf{x}_k - \mathbf{x}_{k+1} \rangle \\ &= \frac{L}{2} \|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2 - \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2. \end{aligned} \quad (8.2.39)$$

Evaluating this expression at $\mathbf{x} = \mathbf{x}_*$, we see that $\|\mathbf{x}_k - \mathbf{x}_*\|_2$ is nonincreasing. Moreover, rearranging the relationship and summing from 0 to k , we obtain that

$$\frac{2}{L} \sum_{i=0}^{k-1} \{F(\mathbf{x}_{i+1}) - F(\mathbf{x}_*)\} \leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2. \quad (8.2.40)$$

Evaluating (8.2.39) at $\mathbf{x} = \mathbf{x}_k$, we obtain

$$F(\mathbf{x}_k) - F(\mathbf{x}_{k+1}) \geq \frac{L}{2} \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_2^2. \quad (8.2.41)$$

Hence, (8.2.40) implies that

$$\frac{2}{L} k \{F(\mathbf{x}_k) - F(\mathbf{x}_*)\} \leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2. \quad (8.2.42)$$

Rearranging, we get the desired conclusion. \square

8.3 Accelerated Proximal Gradient Methods

In the previous section, we have seen that by exploiting the fact that if a (non-smooth) convex function has an easily computable proximal operator, we are able to extend the gradient descent algorithm for smooth functions to the special class of composite objective functions that we encounter in structured signal recovery. The resulting algorithm enjoys the same $O(1/k)$ convergence rate as in the smooth case. Recognizing special structure in our problem of interest yields a significantly more accurate, efficient algorithm.

Interlude: Accelerating the gradient method.

With the taste of victory still on our lips, we might naturally ask whether there are further improvements to be had – is our proximal gradient algorithm optimal for this class of functions? For the question to be meaningful, we need to restrict our attention to methods with efficient iterations, such as gradient-like methods. For example, we could restrict our attention to *first order methods*, which base their future actions on the past iterates $\mathbf{x}_0, \dots, \mathbf{x}_k$, objective values at these iterates, and the gradients $\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_k)$. The corresponding question for *smooth functions* was studied in great depth in the late 1970's and early 1980's by Russian optimization theorists, including Nemirovski, Yudin, and Nesterov.⁴

They asked the very natural question: *for minimizing a smooth function f , is the gradient method optimal amongst first-order methods?* Again, to study this problem one needs a model for computation. They considered a *black box* model, in which the algorithm produces a sequence of iterates $\mathbf{x}_0, \dots, \mathbf{x}_k$. At each iteration, the algorithm is provided with the value $f(\mathbf{x}_i)$ and the gradient $\nabla f(\mathbf{x}_i)$. It produces the next iterate as a function of the history of iterates, gradients, and function values up to this time:

$$\mathbf{x}_{k+1} = \varphi_k(\mathbf{x}_0, \dots, \mathbf{x}_k, f(\mathbf{x}_0), \dots, f(\mathbf{x}_k), \nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_k)). \quad (8.3.1)$$

With this model in mind, one can begin to study algorithms from a worst-case perspective. To do so, we fix a class of functions \mathcal{F} , and ask how well the algorithm does on the “worst function” from this class:

$$\sup_{f \in \mathcal{F}, \mathbf{x}_0} \left\{ f(\mathbf{x}_k) - \inf_{\mathbf{x}} f(\mathbf{x}) \right\}. \quad (8.3.2)$$

One can study various classes of functions \mathcal{F} . However, for our purposes, one interesting class is the convex differentiable functions $f : \mathfrak{B}(0, r) \rightarrow \mathbb{R}$, with L -Lipschitz gradient:

$$\mathcal{F}_{L,r} \doteq \{f : \mathfrak{B}(0, r) \rightarrow \mathbb{R} \mid \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\| \leq L\|\mathbf{x} - \mathbf{x}'\| \forall \mathbf{x}, \mathbf{x}' \in \mathfrak{B}(0, r)\}. \quad (8.3.3)$$

⁴For a more comprehensive introduction to this circle of ideas, see [Nesterov, 2003a] or [Nemirovski, 1995].

The gradient method achieves a rate of $O(1/k)$ over this class:

$$\sup_{f \in \mathcal{F}_{L,r}, \mathbf{x}_0} \left\{ f(\mathbf{x}_k) - \inf_{\mathbf{x}} f(\mathbf{x}) \right\} \leq \frac{CLr^2}{k}. \quad (8.3.4)$$

However, the best lower bound that anyone could prove was of much lower order:

$$\sup_{f \in \mathcal{F}_{L,r}, \mathbf{x}_0} \left\{ f(\mathbf{x}_k) - \inf_{\mathbf{x}} f(\mathbf{x}) \right\} \geq \frac{cLr^2}{k^2}. \quad (8.3.5)$$

Was this merely a gap in the theory? Or might there actually exist a “faster” gradient method than gradient descent itself?

In 1983, Yuri Nesterov closed this gap, to remarkable effect [Nesterov, 1983]. He demonstrated a relatively simple first-order method, which achieved the optimal rate of convergence, $O(1/k^2)$. The analysis of this algorithm is straightforward to read – the 1983 paper is only five pages! However, it is not straightforward to build intuition into *how* the method achieves this rate. To gain some loose appreciation for what is going on, we start from a simpler idea.

Let us first consider the gradient descent method for a smooth function with Lipschitz gradient. At each iteration, the update simply follows the direction of the gradient:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k),$$

where a good choice of α for an L -Lipschitz function is $1/L$. The gradient $-\nabla f(\mathbf{x}_k)$ indicates the direction in which the function drops its value the fastest. Instead of updating along this most greedy direction at the current estimate \mathbf{x}_k , a slightly more conservative strategy is to update by keeping some momentum from the previous update direction: $\mathbf{x}_k - \mathbf{x}_{k+1}$. That leads to an update rule that is known as the *heavy ball method*:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) + \beta(\mathbf{x}_k - \mathbf{x}_{k-1}). \quad (8.3.6)$$

The heavy ball method often leads to smoother update trajectories and faster convergence.

Like the heavy ball method, Nesterov’s method uses a momentum step. It introduces an auxiliary point \mathbf{p}_{k+1} of the form similar to that in the heavy ball method:

$$\mathbf{p}_{k+1} \doteq \mathbf{x}_k + \beta_k(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

At each iteration, we move to this new point, and then descend from it:

$$\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \alpha \nabla f(\mathbf{p}_{k+1}). \quad (8.3.7)$$

The weights $\alpha = 1/L$ and β_k are carefully chosen to achieve the optimal convergence rate:

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}, \quad \beta_k = \frac{t_{k-1} - 1}{t_k}. \quad (8.3.8)$$

These particular values come from the convergence analysis. One can rigorously show that with this update scheme, the resulting algorithm achieves the theoretically optimal convergence rate $O(1/k^2)$ for the class of smooth functions with Lipschitz gradient.

Accelerating the proximal gradient method.

As we have seen in the previous section, convex programs that arise in the context of structured signal recovery are often of the composite form $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$, where f is a smooth term whose gradient is L -Lipschitz and $g(\mathbf{x})$ is convex but not necessarily smooth. In the proximal gradient method introduced in the previous section, we have shown that at the k -th iteration, the value of the objective function $F(\mathbf{x})$ can be upper-bounded by

$$\begin{aligned} \hat{F}(\mathbf{x}, \mathbf{x}_k) &\doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 + g(\mathbf{x}) \\ &\doteq \frac{L}{2} \|\mathbf{x} - \mathbf{w}_k\|_2^2 + g(\mathbf{x}) + \text{terms that do not depend on } \mathbf{x}, \end{aligned}$$

where $\mathbf{w}_k = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)$. As we have seen in the previous section, the gradient descent algorithm for the smooth part $f(\mathbf{x})$ corresponds to directly minimizing its quadratic approximation of $f(\mathbf{x})$ at \mathbf{x}_k .

In this language, if $g \equiv 0$, Nesterov's method corresponds to: *extrapolating* to find the point \mathbf{p}_{k+1} based on the past two iterates, and then *minimizing* a quadratic upper bound \hat{f} to f , taken at the new point \mathbf{p}_{k+1} , by taking a gradient step. Let us attempt to do the same thing for our composite function \hat{F} . Set

$$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (8.3.9)$$

instead of the current estimate \mathbf{x}_k . Then minimize $\hat{F}(\mathbf{x}, \mathbf{p}_{k+1})$ to obtain the next iterate \mathbf{x}_{k+1} :

$$\mathbf{x}_{k+1} = \text{prox}_g \left(\mathbf{p}_{k+1} - \frac{1}{L} \nabla f(\mathbf{p}_{k+1}) \right). \quad (8.3.10)$$

We summarize this scheme in Figure 8.3. Theorem 8.3.1 establishes that with this simple modification to the proximal gradient method, the resulting new algorithm achieves the theoretically optimal convergence rate $O(1/k^2)$ for this class of methods, despite the presence of a nonsmooth term in the objective function.

Theorem 8.3.1. *Let the sequence $\{\mathbf{x}_k\}$ be generated by the above accelerated proximal gradient scheme for the convex composite function $F(\mathbf{x}) =$*

Accelerated Proximal Gradient (APG)**Problem Class:**

$$\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$$

$f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex, ∇f L -Lipschitz and g non-smooth.

Basic Iteration: set $\mathbf{x}_0 \in \mathbb{R}^n$, $\mathbf{p}_1 \leftarrow \mathbf{x}_0$, and $t_0 \leftarrow 1$.

Repeat:

$$t_k \leftarrow \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}, \quad \beta_k \leftarrow \frac{t_{k-1} - 1}{t_k}.$$

$$\mathbf{p}_{k+1} \leftarrow \mathbf{x}_k + \beta_k(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

$$\mathbf{x}_{k+1} \leftarrow \text{prox}_g\left(\mathbf{p}_{k+1} - \frac{1}{L}\nabla f(\mathbf{p}_{k+1})\right).$$

Convergence Guarantee:

$F(\mathbf{x}_k) - F(\mathbf{x}_\star)$ converges at a rate of $O(1/k^2)$.

Figure 8.3. Properties of the Accelerated Proximal Gradient Method.

$f(\mathbf{x}) + g(\mathbf{x})$, where the gradient of f is L -Lipschitz. Let \mathbf{x}_\star be a minimizer of $F(\mathbf{x})$. Then for any $k \geq 1$,

$$F(\mathbf{x}_k) - F(\mathbf{x}_\star) \leq \frac{2L\|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{(k+1)^2}.$$

8.3.1 APG for Basis Pursuit Denoising

Applying the APG algorithm in Figure 8.3 to the basis pursuit denoising problem 8.1.2, we obtain Algorithm 8.3.

Algorithm 8.3 Accelerated Proximal Gradient (APG) for BPDN

- 1: **Problem:** $\min_{\mathbf{x}} \frac{1}{2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda\|\mathbf{x}\|_1$, given $\mathbf{y} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times n}$.
- 2: **Input:** $\mathbf{x}_0 \in \mathbb{R}^n$, $\mathbf{p}_1 \leftarrow \mathbf{x}_0$, and $t_1 \leftarrow 1$, and $L \geq \lambda_{\max}(\mathbf{A}^*\mathbf{A})$.
- 3: **while** \mathbf{x}_k not converged ($k = 1, 2, \dots$) **do**
- 4: $\mathbf{w}_{k+1} \leftarrow \mathbf{p}_k - \frac{1}{L}\mathbf{A}^*(\mathbf{A}\mathbf{p}_k - \mathbf{y})$.
- 5: $\mathbf{x}_{k+1} \leftarrow \text{soft}(\mathbf{w}_{k+1}, \lambda/L)$.
- 6: $t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}$; $\beta_{k+1} \leftarrow \frac{t_k - 1}{t_{k+1}}$.
- 7: $\mathbf{p}_{k+1} \leftarrow \mathbf{x}_{k+1} + \beta_{k+1}(\mathbf{x}_{k+1} - \mathbf{x}_k)$.
- 8: **end while**
- 9: **Output:** $\mathbf{x}_\star \leftarrow \mathbf{x}_k$.

8.3.2 APG for Stable Principal Component Pursuit

Similarly we could apply the APG algorithm in Figure 8.3 to solving the stable principal component pursuit (PCP) problem 8.2.26. Again, notice that the APG scheme respects the natural separable structure of the objective function.

Algorithm 8.4 Accelerated Proximal Gradient (APG) for Stable PCP

1: **Problem:** $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2$, given \mathbf{Y} , $\lambda, \mu > 0$.

2: **Input:** $\mathbf{L}_0 \in \mathbb{R}^{m \times n}$, $\mathbf{S}_0 \in \mathbb{R}^{m \times n}$, $t_1 \leftarrow 1$.

3: **while** $(\mathbf{L}_k, \mathbf{S}_k)$ not converged ($k = 1, 2, \dots$) **do**

4: $t_k \leftarrow \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$, $\beta_k \leftarrow \frac{t_{k-1} - 1}{t_k}$.

5: $\mathbf{P}_k^L \leftarrow \mathbf{L}_k + \beta_k(\mathbf{L}_k - \mathbf{L}_{k-1})$.

6: $\mathbf{P}_k^S \leftarrow \mathbf{S}_k + \beta_k(\mathbf{S}_k - \mathbf{S}_{k-1})$.

7: $\mathbf{W}_k \leftarrow \mathbf{Y} - \mathbf{P}_k^S$ and compute the SVD: $\mathbf{W}_k = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$.

8: $\mathbf{L}_{k+1} \leftarrow \mathbf{U}\text{soft}(\mathbf{\Sigma}, 1/\mu)\mathbf{V}^*$.

9: $\mathbf{S}_{k+1} \leftarrow \text{soft}((\mathbf{Y} - \mathbf{P}_k^L), \lambda/\mu)$.

10: **end while**

11: **Output:** $\mathbf{L}_\star \leftarrow \mathbf{L}_k$; $\mathbf{S}_\star \leftarrow \mathbf{S}_k$.

8.3.3 Convergence of APG

Our convergence analysis follows, almost verbatim, Beck and Teboulle [Beck and Teboulle, 2009]. Let $\varphi(\mathbf{y})$ denote the operator that takes a step in the direction of the gradient of f at \mathbf{y} , and then applies the proximal operator of g/L :

$$\varphi(\mathbf{y}) = \text{prox}_{g/L} \left(\mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y}) \right). \quad (8.3.11)$$

In this language, the accelerated proximal gradient iteration is

$$\mathbf{x}_{k+1} = \varphi(\mathbf{p}_{k+1}). \quad (8.3.12)$$

The following lemma allows us to compare the value of $F(\mathbf{x})$ at any point \mathbf{x} to the value at $\varphi(\mathbf{p})$ for an arbitrary point \mathbf{p} :

Lemma 8.3.2. For every $\mathbf{x}, \mathbf{p} \in \mathbb{R}^n$,

$$F(\mathbf{x}) - F(\varphi(\mathbf{p})) \geq \frac{L}{2} \|\varphi(\mathbf{p}) - \mathbf{p}\|_2^2 + L \langle \mathbf{p} - \mathbf{x}, \varphi(\mathbf{p}) - \mathbf{p} \rangle. \quad (8.3.13)$$

Proof. For it, we note that from the optimality condition for the proximal problem, $\mathbf{z} = \varphi(\mathbf{p})$ if and only if there exists $\boldsymbol{\gamma} \in \partial g(\mathbf{z})$ such that

$$\boldsymbol{\gamma} + L(\mathbf{z} - \mathbf{p}) + \nabla f(\mathbf{p}) = \mathbf{0}. \quad (8.3.14)$$

Using the subgradient inequalities for f and g , we obtain that

$$f(\mathbf{x}) \geq f(\mathbf{p}) + \langle \mathbf{x} - \mathbf{p}, \nabla f(\mathbf{p}) \rangle \quad (8.3.15)$$

$$g(\mathbf{x}) \geq g(\varphi(\mathbf{p})) + \langle \mathbf{x} - \varphi(\mathbf{p}), \gamma \rangle. \quad (8.3.16)$$

Hence,

$$\begin{aligned} F(\mathbf{x}) - F(\varphi(\mathbf{p})) &\geq f(\mathbf{x}) + g(\varphi(\mathbf{p})) + \langle \mathbf{x} - \mathbf{p}, \nabla f(\mathbf{p}) \rangle \\ &\quad + \langle \mathbf{x} - \varphi(\mathbf{p}), \gamma \rangle - F(\varphi(\mathbf{p})) \\ &\geq f(\mathbf{x}) + g(\varphi(\mathbf{p})) + \langle \mathbf{x} - \mathbf{p}, \nabla f(\mathbf{p}) \rangle \\ &\quad + \langle \mathbf{x} - \varphi(\mathbf{p}), \gamma \rangle - \hat{F}(\varphi(\mathbf{p}), \mathbf{p}) \\ &= -\frac{L}{2} \|\varphi(\mathbf{p}) - \mathbf{p}\|_2^2 + \langle \mathbf{x} - \varphi(\mathbf{p}), \nabla f(\mathbf{p}) + \gamma \rangle \\ &= -\frac{L}{2} \|\varphi(\mathbf{p}) - \mathbf{p}\|_2^2 + L \langle \mathbf{x} - \varphi(\mathbf{p}), \mathbf{p} - \varphi(\mathbf{p}) \rangle \\ &= \frac{L}{2} \|\varphi(\mathbf{p}) - \mathbf{p}\|_2^2 + L \langle \mathbf{p} - \mathbf{x}, \varphi(\mathbf{p}) - \mathbf{p} \rangle, \quad (8.3.17) \end{aligned}$$

as desired. \square

Using Lemma 8.3.2, we obtain a relationship between the suboptimality in function values and the distance of an interpolated point to the optimum:

Lemma 8.3.3. *Let $(\mathbf{x}_k, \mathbf{p}_k)$ be the sequence generated by the proximal gradient method. Set*

$$v_k = F(\mathbf{x}_k) - F(\mathbf{x}_\star) \quad (8.3.18)$$

and

$$\mathbf{u}_k = t_k \mathbf{x}_k - (t_k - 1) \mathbf{x}_{k-1} - \mathbf{x}_\star. \quad (8.3.19)$$

Then

$$\frac{2}{L} t_k^2 v_k - \frac{2}{L} t_{k+1}^2 v_{k+1} \geq \|\mathbf{u}_{k+1}\|_2^2 - \|\mathbf{u}_k\|_2^2. \quad (8.3.20)$$

Proof. Let us apply the previous lemma with $\mathbf{x} = \mathbf{x}_k$ and $\mathbf{p} = \mathbf{p}_{k+1}$. This gives

$$\frac{2}{L} (v_k - v_{k+1}) \geq \|\mathbf{x}_{k+1} - \mathbf{p}_{k+1}\|_2^2 + 2 \langle \mathbf{p}_{k+1} - \mathbf{x}_k, \mathbf{x}_{k+1} - \mathbf{p}_{k+1} \rangle. \quad (8.3.21)$$

Applying the lemma with $\mathbf{x} = \mathbf{x}_\star$ and $\mathbf{p} = \mathbf{p}_{k+1}$, we get

$$-\frac{2}{L} v_{k+1} \geq \|\mathbf{x}_{k+1} - \mathbf{p}_{k+1}\|_2^2 + 2 \langle \mathbf{p}_{k+1} - \mathbf{x}_\star, \mathbf{x}_{k+1} - \mathbf{p}_{k+1} \rangle. \quad (8.3.22)$$

Combining these two inequalities, we get

$$\begin{aligned} &\frac{2}{L} ((t_{k+1} - 1) v_k - t_{k+1} v_{k+1}) \\ &\geq t_{k+1} \|\mathbf{x}_{k+1} - \mathbf{p}_{k+1}\|_2^2 + 2 \langle \mathbf{x}_{k+1} - \mathbf{p}_{k+1}, t_{k+1} \mathbf{p}_{k+1} - (t_{k+1} - 1) \mathbf{x}_k - \mathbf{x}_\star \rangle. \end{aligned}$$

Multiplying both sides by t_{k+1} , and using that $t_k^2 = t_{k+1}(t_{k+1} - 1)$, we get

$$\begin{aligned}
& \frac{2}{L} (t_k^2 v_k - t_{k+1}^2 v_{k+1}) \\
& \geq \|t_{k+1}(\mathbf{x}_{k+1} - \mathbf{p}_{k+1})\|_2^2 + 2t_{k+1} \langle \mathbf{x}_{k+1} - \mathbf{p}_{k+1}, t_{k+1}\mathbf{p}_{k+1} - (t_{k+1} - 1)\mathbf{x}_k - \mathbf{x}_\star \rangle \\
& = \|t_{k+1}\mathbf{x}_{k+1} - (t_{k+1} - 1)\mathbf{x}_k - \mathbf{x}_\star\|_2^2 - \|t_{k+1}\mathbf{p}_{k+1} - (t_{k+1} - 1)\mathbf{x}_k - \mathbf{x}_\star\|_2^2 \\
& = \|\mathbf{u}_{k+1}\|_2^2 - \|\mathbf{u}_k\|_2^2,
\end{aligned}$$

as desired. \square

To prove the desired result, we note two simple facts, and then perform a calculation. First,

Lemma 8.3.4. *Let (a_k, b_k) be sequences of positive real numbers satisfying*

$$a_k - a_{k+1} \geq b_{k+1} - b_k \quad \forall k, \quad a_1 + b_1 \leq c. \quad (8.3.23)$$

Then $a_k \geq c$ for every k .

Second,

Lemma 8.3.5. *The sequence t_k generated by the accelerated proximal gradient method satisfies*

$$t_k \geq \frac{k+1}{2} \quad \forall k \geq 1. \quad (8.3.24)$$

With these facts in mind, we calculate:

Proof. Define

$$a_k = \frac{2}{L} t_k^2 v_k \quad (8.3.25)$$

$$b_k = \|\mathbf{u}_k\|_2^2 \quad (8.3.26)$$

$$c = \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2. \quad (8.3.27)$$

By Lemma 8.3.3, for every k ,

$$a_k - a_{k+1} \geq b_{k+1} - b_k \quad (8.3.28)$$

so, provided $a_1 + b_1 \geq c$, we obtain that $a_k \leq c$ for every k , whence

$$\frac{2}{L} t_k^2 v_k \leq \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2. \quad (8.3.29)$$

Since $t_k \geq (k+1)/2$, this gives

$$F(\mathbf{x}_k) - F(\mathbf{x}_\star) \leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{(k+1)^2}. \quad (8.3.30)$$

So, it just remains to check that $a_1 + b_1 \leq c$. Since $t_1 = 1$, $a_1 = \frac{2}{L}v_1$, while $b_1 = \|\mathbf{x}_1 - \mathbf{x}_\star\|_2^2$. In Lemma 8.3.2, set $\mathbf{x} = \mathbf{x}_\star$, $\mathbf{p} = \mathbf{p}_1$, to obtain

$$F(\mathbf{x}_\star) - F(\mathbf{x}_1) \geq \frac{L}{2} \|\mathbf{x}_1 - \mathbf{p}_1\|_2^2 + L \langle \mathbf{p}_1 - \mathbf{x}_\star, \mathbf{x}_1 - \mathbf{p}_1 \rangle \quad (8.3.31)$$

$$= \frac{L}{2} \left(\|\mathbf{x}_1 - \mathbf{x}_\star\|_2^2 - \|\mathbf{p}_1 - \mathbf{x}_\star\|_2^2 \right). \quad (8.3.32)$$

Since $\|\mathbf{p}_1 - \mathbf{x}_\star\|_2^2 = \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2$, this gives the result. \square

8.4 Augmented Lagrange Multipliers

So far, we have described how to solve certain classes of *unconstrained* convex optimization problems arising in structured signal recovery. However, in some scenarios – e.g., if the noise level is low, or if the target solution \mathbf{x} is known ahead of time to possess additional application-specific structure – it may be desirable to exactly enforce equality constraints such as in the exact BP program (8.1.1) or the PCP (8.1.5).

In this section, we describe a framework for solving *equality constrained* problems of the form

$$\begin{aligned} & \text{minimize} && g(\mathbf{x}), \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y}, \end{aligned} \quad (8.4.1)$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a matrix and $\mathbf{y} \in \text{range}(\mathbf{A})$ (so that the problem is feasible). One very intuitive approach to producing an approximate solution to (8.4.1) is to simply replace the inequality constraint $\mathbf{A}\mathbf{x} = \mathbf{y}$ with a penalty function $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$, and solve the unconstrained problem

$$\text{minimize} \quad g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \quad (8.4.2)$$

for a very large value of μ . As μ increases to $+\infty$, the solution set of this problem approaches the solution set of equality constrained problem (8.4.1). This is known as the *penalty method* in the optimization literature, and has a long history, with many variants. Its main advantage is that it leaves us with a simpler unconstrained problem, to which we can directly apply scalable first order methods such as the proximal gradient methods of the previous two sections.

However, there is a serious drawback to this approach. For first order methods such as PG and APG, the rate of convergence is dictated by how quickly the gradient $\nabla(\mu f) = \mu \mathbf{A}^*(\mathbf{A}\mathbf{x} - \mathbf{y})$ can change from point-to-point, which is measured through the Lipschitz constant

$$L_{\nabla \mu f} = \mu \|\mathbf{A}\|_2^2.$$

This increases linearly with μ : *The larger μ is, the harder the unconstrained problem (8.4.2) is to solve!* One practical approach to mitigating this effect

is to solve a sequence of unconstrained problems, with increasing μ , and use the solution to each as an initial guess for the next. This *continuation* technique is often valuable in practice. Nevertheless, it suffers from the same drawback: as μ increases, accurate solutions become increasingly difficult to obtain.

Lagrange duality gives a more principled mechanism for studying and solving the constrained problem (8.4.1) via unconstrained problems. In particular, it will give us a mechanism for exactly solving the constrained problem (8.4.1) by solving a sequence of unconstrained problems *whose difficulty does not increase*. The central object in Lagrange duality is the Lagrangian

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \doteq g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle, \quad (8.4.3)$$

where $\boldsymbol{\lambda} \in \mathbb{R}^m$ is a vector of *Lagrange multipliers* corresponding to the equality constraint $\mathbf{Ax} = \mathbf{y}$. In particular, we can characterize optimal solutions $(\mathbf{x}, \boldsymbol{\lambda})$ as saddle points of the Lagrangian

$$\sup_{\boldsymbol{\lambda}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \sup_{\boldsymbol{\lambda}} \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle. \quad (8.4.4)$$

If we define the dual function

$$d(\boldsymbol{\lambda}) \doteq \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle. \quad (8.4.5)$$

The saddle point characterization of optimal solutions suggests a natural computational approach to finding $(\mathbf{x}, \boldsymbol{\lambda})$:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (8.4.6)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + t_{k+1}(\mathbf{Ax}_{k+1} - \mathbf{y}). \quad (8.4.7)$$

It is not difficult to show that $\mathbf{Ax}_{k+1} - \mathbf{y}$ is a subgradient⁵ of the dual function $d(\boldsymbol{\lambda})$. This iteration corresponds to a subgradient ascent algorithm for maximizing the dual function, and hence is called *dual ascent*. In (8.4.7), t_{k+1} is the step size. For certain problem classes, dual ascent yields efficient, convergent algorithms, which produce an optimal primal-dual pair $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$. However, for problems arising in structured signal recovery, the straightforward iteration (8.4.6)-(8.4.7) may fail:

Example 8.4.1. *Show that*

$$\inf_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle = \begin{cases} -\infty & \|\mathbf{A}^* \boldsymbol{\lambda}\|_\infty > 1, \\ -\langle \boldsymbol{\lambda}, \mathbf{y} \rangle & \|\mathbf{A}^* \boldsymbol{\lambda}\|_\infty \leq 1. \end{cases} \quad (8.4.8)$$

So, for basis pursuit, if Dual Ascent happens to produce a $\boldsymbol{\lambda}$ such that $\|\mathbf{A}^ \boldsymbol{\lambda}\|_\infty > 1$, the algorithm will break down. Notice, in particular, that when $\|\mathbf{A}^* \boldsymbol{\lambda}\|_\infty > 1$, we can produce arbitrarily large (in magnitude)*

⁵Strictly, a *supergradient* since the dual $d(\boldsymbol{\lambda})$ is concave.

negative values of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ by choosing \mathbf{x} far away from the feasible set $\{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{y}\}$.

This bad behavior occurs more generally. Thus, for structured signal recovery, the classical Lagrangian is sufficient for characterizing optimality conditions, but it does not penalize the equality $\mathbf{A}\mathbf{x} = \mathbf{y}$ strongly enough for (8.4.6)-(8.4.7) to lead to a useful algorithm. A natural remedy is to *augment* the Lagrangian with an extra penalty term, by introducing the function

$$\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}) \doteq g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2. \quad (8.4.9)$$

This function is known as the *Augmented Lagrangian* [Hestenes, 1969, Rockafellar, 1973, Powell, 1969]. As before, $\mu > 0$ is a penalty parameter. The augmented Lagrangian can be regarded as the Lagrangian for the constrained problem

$$\begin{aligned} & \text{minimize} && g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y}. \end{aligned} \quad (8.4.10)$$

Since the penalty term $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ is zero for all feasible \mathbf{x} , the optimal solutions of this problem coincide with the optimal solutions of the original problem (8.4.1).

Despite this formal equivalence, augmentation has dramatic consequences for numerical optimization. In particular, it can render the dual ascent iteration provably convergent, under very weak assumptions on the objective function g . To achieve this, we apply dual ascent to the penalized problem (8.4.10), and make a very particular choice of step size, $t_k = \mu$, yielding the iteration

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (8.4.11)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}). \quad (8.4.12)$$

This iteration, with the particular choice $t_k = \mu$ is known as the *Method of Multipliers*. The update step (8.4.11) for \mathbf{x} is a convex optimization problem itself and can typically be solved via the proximal gradient methods introduced in the previous sections.

Remark 8.4.2. The choice $t_k = \mu$ is important, because it allows us to avoid the breakdown described in Example 8.4.1: To see this, since \mathbf{x}_{k+1} minimizes the convex function \mathcal{L}_μ ,

$$\begin{aligned} \mathbf{0} & \in \partial \mathcal{L}_\mu(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_k), \\ & = \partial g(\mathbf{x}_{k+1}) + \mathbf{A}^* \boldsymbol{\lambda}_k + \mu \mathbf{A}^* (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}), \\ & = \partial g(\mathbf{x}_{k+1}) + \mathbf{A}^* \boldsymbol{\lambda}_{k+1}, \\ & = \partial \mathcal{L}(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}). \end{aligned}$$

Thus, \mathbf{x}_{k+1} minimizes the unaugmented Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_{k+1})$ with $\boldsymbol{\lambda} = \boldsymbol{\lambda}_{k+1}$ fixed. This means that $d(\boldsymbol{\lambda}_{k+1}) > -\infty$, and $\boldsymbol{\lambda}_{k+1}$ is dual feasible for

the original problem. In particular, $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_{k+1})$ is bounded below. Because $\boldsymbol{\lambda}_{k+1}$ is always dual feasible, the bad behavior in Example 8.4.1 cannot occur.

Under appropriate assumptions on g , the iterates \mathbf{x}_k produced by this modified algorithm converge to an optimal solution \mathbf{x}_* to the constrained problem (8.4.1). We state a slightly more general result that allows the penalty parameters μ to vary from iteration to iteration, as long as they remain bounded away from zero:

Theorem 8.4.3 (Convergence of Augmented Lagrangian). *Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex, coercive function, $\mathbf{A} \in \mathbb{R}^{m \times n}$ an arbitrary matrix, and $\mathbf{y} \in \text{range}(\mathbf{A})$. Then the problem*

$$\begin{aligned} & \text{minimize} && g(\mathbf{x}) \\ & \text{subject to} && \mathbf{Ax} = \mathbf{y}, \end{aligned} \tag{8.4.13}$$

has at least one optimal solution. Moreover, the ALM iteration

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_{\mu_k}(\mathbf{x}, \boldsymbol{\lambda}_k), \tag{8.4.14}$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu_k (\mathbf{Ax}_{k+1} - \mathbf{y}). \tag{8.4.15}$$

with sequence (μ_k) bounded away from zero produces a sequence $\boldsymbol{\lambda}_k$ that converges to a dual optimal solution. Moreover, every limit point of the sequence \mathbf{x}_k is optimal for (8.4.13).

Figure 8.4 summarizes our general observations on the ALM method up to this point.

Remark 8.4.4 (More general convergence theorems). *The statement of Theorem 8.4.3 represents a deliberate tradeoff between simplicity and generality. With somewhat more technical analysis, it is possible to show convergence ALM for more much general classes of g . The most practically important extension allows g to be an extended real valued function (a function from \mathbb{R}^n to $\mathbb{R} \cup \{+\infty\}$). For example, if we wish to optimize a real-valued convex function g_0 over the set of \mathbf{x} that satisfy the equality constraint $\mathbf{Ax} = \mathbf{y}$, and reside in some additional (nonempty closed, convex) constraint set C :*

$$\begin{aligned} & \text{minimize} && g_0(\mathbf{x}) \\ & \text{subject to} && \mathbf{Ax} = \mathbf{y}, \mathbf{x} \in C, \end{aligned} \tag{8.4.16}$$

we can apply ALM to the problem

$$\begin{aligned} & \text{minimize} && g(\mathbf{x}) \doteq g_0(\mathbf{x}) + \mathcal{I}_{\mathbf{x} \in C} \\ & \text{subject to} && \mathbf{Ax} = \mathbf{y} \end{aligned} \tag{8.4.17}$$

where $\mathcal{I}_{\mathbf{x} \in C}$ is the indicator function for C :

$$\mathcal{I}_{\mathbf{x} \in C} = \begin{cases} 0 & \mathbf{x} \in C \\ +\infty & \mathbf{x} \notin C. \end{cases} \tag{8.4.18}$$

Augmented Lagrange Multiplier (ALM)

Problem Class:

$$\begin{aligned} & \text{minimize} && g(\mathbf{x}) \\ & \text{subject to} && \mathbf{Ax} = \mathbf{y}. \end{aligned}$$

$$g : \mathbb{R}^n \rightarrow \mathbb{R} \text{ convex, } \mathbf{y} \in \text{range}(\mathbf{A}).$$

Basic Iteration: set

$$\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}) = g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2.$$

Repeat:

$$\begin{aligned} \mathbf{x}_{k+1} &\in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}_k), \\ \boldsymbol{\lambda}_{k+1} &= \boldsymbol{\lambda}_k + \mu (\mathbf{Ax}_{k+1} - \mathbf{y}). \end{aligned}$$

Convergence Guarantee:

If g is coercive, every limit point of (\mathbf{x}_k) is optimal.

Figure 8.4. Properties of the Augmented Lagrangian Method of Multipliers.

The survey of Eckstein [Eckstein, 2012] and the monograph of Bertsekas [Bertsekas, 1982b] are good introductory points for the more general theory, which enables such modifications.

Implementation Considerations.

The most important practical consideration is how to choose the sequence of penalty parameters $\mu^{(k)}$. As discussed above, this choice induces a trade-off between the cost of solving subproblems and the overall number of outer iterations – larger μ leaves us with fewer outer iterations, but harder subproblems. A typical strategy is to increase μ geometrically, up to some pre-fixed ceiling:

$$\mu_k = \min \{ \beta \mu_k, \mu_{\max} \},$$

where $\beta \approx 1.25$ is typical. The ceiling μ_{\max} is strongly problem dependent; choosing it “optimally” is something of a black art.

Our description and analysis of ALM assume that each of the subproblems is solved exactly. However, practically speaking, it may not be necessary to obtain high-accuracy solutions to the subproblems, especially in the early iterations. This can be justified theoretically. The choice of iterative method for solving the unconstrained subproblems is largely problem dependent. However, because the penalty term is quadratic, for many prob-

lems of interest in this text, the subproblems have composite form, and the APG algorithm applies.

In using APG (or any other iterative solver) to solve the unconstrained subproblems, it is highly advisable to use the previous iterate \mathbf{x}_k as an initialization to solve for the subsequent iterate \mathbf{x}_{k+1} . While the subproblems are convex, and the global optimality of iterative algorithms does not depend on initialization, choosing an appropriate initializer can drastically reduce the overall number of iterations.

8.4.1 ALM for Basis Pursuit

We may apply ALM to the exact BP problem (8.1.1), which we summarize as Algorithm 8.5. This algorithm was introduced by [Yin et al., 2008], where it was interpreted as a *Bregman iteration*.

Algorithm 8.5 Augmented Lagrange Multiplier (ALM) for BP

- 1: **Problem:** $\min_{\mathbf{x}} \|\mathbf{x}\|_1$ subj. to $\mathbf{y} = \mathbf{A}\mathbf{x}$, given $\mathbf{y} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$.
 - 2: **Input:** $\mathbf{x}_0 \in \mathbb{R}^n$ and $\beta > 1$.
 - 3: **while** $\{\mathbf{x}_k, \boldsymbol{\lambda}_k\}$ not converged ($k = 1, 2, \dots$) **do**
 - 4: $\mathbf{x}_{k+1} \leftarrow \arg \min \mathcal{L}_{\mu_k}(\mathbf{x}, \boldsymbol{\lambda}_k)$ using APG.
 - 5: $\boldsymbol{\lambda}_{k+1} \leftarrow \boldsymbol{\lambda}_k + \mu_k(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y})$.
 - 6: $\mu_{k+1} \leftarrow \min \{\beta\mu_k, \mu_{\max}\}$.
 - 7: **end while**
 - 8: **Output:** $\mathbf{x}_\star \leftarrow \mathbf{x}_k$.
-

8.4.2 ALM for Principal Component Pursuit

We recall principal component pursuit (PCP) solves the following program:

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \quad \text{subj. to} \quad \mathbf{Y} = \mathbf{L} + \mathbf{S}. \quad (8.4.19)$$

First, we rewrite the above program as a standard ALM objective function:

$$\mathcal{L}_\mu(\mathbf{L}, \mathbf{S}, \boldsymbol{\Lambda}) \doteq \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \langle \boldsymbol{\Lambda}, \mathbf{Y} - \mathbf{L} - \mathbf{S} \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2,$$

where $\mathcal{L}_\mu(\cdot)$ consists of a Lagrangian term with a Lagrange multiplier matrix $\boldsymbol{\Lambda}$ of the same size as \mathbf{Y} and an augmented quadratic term that encourages the equality condition $\mathbf{Y} = \mathbf{L} + \mathbf{S}$.

Algorithm 8.6 Augmented Lagrange Multiplier (ALM) for PCP

```

1: Problem:  $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1$  subj to  $\mathbf{Y} = \mathbf{L} + \mathbf{S}$ .
2: Input:  $\mathbf{L}_0, \mathbf{S}_0 \in \mathbb{R}^{m \times n}$  and  $\beta > 1$ .
3: while  $\{\mathbf{L}_k, \mathbf{S}_k\}$  not converged ( $k = 1, 2, \dots$ ). do
4:    $\{\mathbf{L}_{k+1}, \mathbf{S}_{k+1}\} \leftarrow \arg \min \mathcal{L}_{\mu_k}(\mathbf{L}, \mathbf{S}, \boldsymbol{\Lambda}_k)$  using APG.
5:    $\boldsymbol{\Lambda}_{k+1} \leftarrow \boldsymbol{\Lambda}_k + \mu_k(\mathbf{Y} - \mathbf{L}_{k+1} - \mathbf{S}_{k+1})$ .
6:    $\mu_{k+1} \leftarrow \min \{\beta \mu_k, \mu_{\max}\}$ .
7: end while
8: Output:  $\mathbf{L}_* \leftarrow \mathbf{L}_k, \mathbf{S}_* \leftarrow \mathbf{S}_k$ .

```

8.4.3 Convergence of ALM

In this subsection, we prove Theorem 8.4.3. The proof will actually reveal another interpretation of the Augmented Lagrangian method, as an application of the proximal point algorithm to the dual problem.

Proof. Let $d(\boldsymbol{\lambda})$ denote the dual function

$$d(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle. \quad (8.4.20)$$

The dual function is concave, and so its negative

$$q(\boldsymbol{\lambda}) = -d(\boldsymbol{\lambda}) \quad (8.4.21)$$

is convex.

Note that for any $\boldsymbol{\lambda}$,

$$\begin{aligned} d(\boldsymbol{\lambda}) &\leq g(\mathbf{x}_{k+1}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x}_{k+1} - \mathbf{y} \rangle \\ &= g(\mathbf{x}_{k+1}) + \langle \boldsymbol{\lambda}_{k+1}, \mathbf{A}\mathbf{x}_{k+1} - \mathbf{y} \rangle + \langle \boldsymbol{\lambda} - \boldsymbol{\lambda}_{k+1}, \mathbf{A}\mathbf{x}_{k+1} - \mathbf{y} \rangle. \end{aligned} \quad (8.4.22)$$

Now recall from Remark 8.4.2 that the augmented Lagrangian method ensures that \mathbf{x}_{k+1} minimizes the unaugmented Lagrangian $g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle$ with $\boldsymbol{\lambda} = \boldsymbol{\lambda}_{k+1}$ fixed. Hence, by definition of the function $d(\boldsymbol{\lambda})$, we have $d(\boldsymbol{\lambda}_{k+1}) = g(\mathbf{x}_{k+1}) + \langle \boldsymbol{\lambda}_{k+1}, \mathbf{A}\mathbf{x}_{k+1} - \mathbf{y} \rangle$. Applying this to the above inequality, we obtain

$$d(\boldsymbol{\lambda}) \leq d(\boldsymbol{\lambda}_{k+1}) + \langle \boldsymbol{\lambda} - \boldsymbol{\lambda}_{k+1}, \mathbf{A}\mathbf{x}_{k+1} - \mathbf{y} \rangle. \quad (8.4.23)$$

As $q(\boldsymbol{\lambda}) = -d(\boldsymbol{\lambda})$, we have

$$q(\boldsymbol{\lambda}) \geq q(\boldsymbol{\lambda}_{k+1}) + \langle \boldsymbol{\lambda} - \boldsymbol{\lambda}_{k+1}, \mathbf{y} - \mathbf{A}\mathbf{x}_{k+1} \rangle. \quad (8.4.24)$$

Hence $\mathbf{y} - \mathbf{A}\mathbf{x}_{k+1}$ is in the subgradient of $q(\cdot)$ at $\boldsymbol{\lambda}_{k+1}$, and

$$\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1} = \mu_k(\mathbf{y} - \mathbf{A}\mathbf{x}_{k+1}) \in \mu_k \partial q(\boldsymbol{\lambda}_{k+1}), \quad (8.4.25)$$

and so

$$\boldsymbol{\lambda}_{k+1} = \text{prox}_{\mu_k q}(\boldsymbol{\lambda}_k). \quad (8.4.26)$$

Thus, dual ascent corresponds to the proximal point iteration applied to $q(\cdot)$. Under our assumptions, strong duality obtains, the dual optimal value

$\sup_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}) > -\infty$ is finite, and a dual optimal solution $\bar{\boldsymbol{\lambda}}$ exists. Proposition 8.2.6 then implies that $\boldsymbol{\lambda}_k \rightarrow \boldsymbol{\lambda}_*$, where $\boldsymbol{\lambda}_*$ is some dual optimal point. This and the fact that μ_k is bounded away from zero give that

$$\|\mathbf{A}\mathbf{x}_k - \mathbf{y}\|_2 = \frac{\|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k-1}\|_2}{\mu_k} \rightarrow 0, \quad (8.4.27)$$

and so the sequence (\mathbf{x}_k) approaches the feasible set. From coercivity of g , there exists at least one primal optimal solution \mathbf{x}_* . By optimality of \mathbf{x}_{k+1} , we have

$$g(\mathbf{x}_{k+1}) + \langle \boldsymbol{\lambda}_k, \mathbf{A}\mathbf{x}_{k+1} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}\|_2^2 \leq g(\mathbf{x}_*). \quad (8.4.28)$$

For any cluster point $\bar{\mathbf{x}}$, continuity of g and $\mathbf{A}\mathbf{x}_k - \mathbf{y} \rightarrow \mathbf{0}$ imply that $g(\bar{\mathbf{x}}) \leq g(\mathbf{x}_*)$, whence $g(\bar{\mathbf{x}}) = g(\mathbf{x}_*)$. Hence, every cluster point is optimal. \square

8.5 Alternating Direction Method of Multipliers

The previous section showed how augmented Lagrangian techniques could be used to solve equality constrained convex optimization problems, by reducing them to a sequence of unconstrained subproblems. These subproblems may still be challenging optimization problems if we need to minimize against all the variables simultaneously. In many situations, though, it is possible to exploit special separable structures of the objective function and to alleviate the difficulty by reducing the overall optimization to multiple parallel subproblems of smaller sizes. We give two prominent examples:

Example 8.5.1 (Principal Component Pursuit). *We solve*

$$\begin{aligned} &\text{minimize} && \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \\ &\text{subject to} && \mathbf{L} + \mathbf{S} = \mathbf{Y}. \end{aligned} \quad (8.5.1)$$

The objective function is separable into two terms, $\|\cdot\|_$ and $\|\cdot\|_1$, each of which has an efficient proximal operator.*

Example 8.5.2 (Distributed Learning). *Special structure also arises in many large scale learning problems, where the goal is to fit a parametric model to a collection of observation vectors $\mathbf{y}_1, \dots, \mathbf{y}_p$. Typically, we are provided with a loss $L(\mathbf{y}, \mathbf{x})$, which could be, e.g., the negative logarithm of the likelihood of observation \mathbf{y} given parameters \mathbf{x} , and our goal is to minimize $\sum_i L(\mathbf{y}_i, \mathbf{x})$ over \mathbf{x} .*

In very large scale applications, it may be prohibitively expensive to store the \mathbf{y}_i centrally, or to transmit them during the operation of an iterative algorithm. Rather, we can assume that they are stored in a distributed fashion, in N locations: the j -th location stores \mathbf{y}_i , $i \in I_j$. The loss on this

subset is $f_j(\mathbf{x}) = \sum_{i \in I_j} L(\mathbf{y}_i, \mathbf{x})$. Our overall goal is then to solve

$$\text{minimize} \quad \sum_{i=1}^N f_i(\mathbf{x}). \quad (8.5.2)$$

Again, this objective function appears to separate into independent terms. To exploit this structure, we can introduce N additional parameter vectors \mathbf{x}_i , which are constrained to coincide with \mathbf{x} :

$$\begin{aligned} &\text{minimize} \quad \sum_{i=1}^N f_i(\mathbf{x}_i) \\ &\text{subject to} \quad \mathbf{x}_i = \mathbf{x} \quad i = 1, \dots, N. \end{aligned} \quad (8.5.3)$$

In this section, we study a family of augmented Lagrangian algorithms that can exploit this special, separable structure. We begin by treating a generic problem of the form

$$\begin{aligned} &\text{minimize} \quad g(\mathbf{x}) + h(\mathbf{z}) \\ &\text{subject to} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}, \end{aligned} \quad (8.5.4)$$

where g and h are convex functions, \mathbf{A} and \mathbf{B} are matrices, and $\mathbf{y} \in \text{range}([\mathbf{A} \mid \mathbf{B}])$. The Lagrangian $\mathcal{L}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$ associated with this problem simply is:

$$\mathcal{L}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = g(\mathbf{x}) + h(\mathbf{z}) + \langle \boldsymbol{\lambda}, \mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{z} \rangle. \quad (8.5.5)$$

As in the previous section, we form the augmented Lagrangian $\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$ associated with this problem:

$$\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = g(\mathbf{x}) + h(\mathbf{z}) + \langle \boldsymbol{\lambda}, \mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{z} \rangle + \frac{\mu}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{z}\|_2^2.$$

In many applications, including the examples listed above, it is easy to minimize \mathcal{L}_μ with respect to \mathbf{x} , when $\boldsymbol{\lambda}$ and \mathbf{z} are fixed, and also easy to minimize it with respect to \mathbf{z} when $\boldsymbol{\lambda}$ and \mathbf{x} are fixed. This suggests a simple, alternating iteration

$$\mathbf{z}_{k+1} \in \arg \min_{\mathbf{z}} \mathcal{L}_\mu(\mathbf{x}_k, \mathbf{z}, \boldsymbol{\lambda}_k), \quad (8.5.6)$$

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}_{k+1}, \boldsymbol{\lambda}_k), \quad (8.5.7)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{y} - \mathbf{A}\mathbf{x}_{k+1} - \mathbf{B}\mathbf{z}_{k+1}). \quad (8.5.8)$$

This is known as the *alternating directions method of multipliers*. In the numerical analysis literature, this style of updating is referred to as a *Gauss-Seidel iteration*. We recommend [Boyd et al., 2011a] for a friendly introduction to these methods, as well as useful recommendations on stopping criteria, parameter setting, etc.

8.5.1 ADMM for Principal Component Pursuit

When applied to the principal component pursuit program, the ADMM iteration takes on a particularly simple form. Here, the two groups of variables are the unknown low-rank matrix \mathbf{L} and the unknown sparse error \mathbf{S} . The augmented Lagrangian is

$$\mathcal{L}_\mu(\mathbf{L}, \mathbf{S}, \mathbf{\Lambda}) = \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \langle \mathbf{\Lambda}, \mathbf{Y} - \mathbf{L} - \mathbf{S} \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2. \quad (8.5.9)$$

The ADMM iteration sequentially updates \mathbf{L} , then \mathbf{S} , then $\mathbf{\Lambda}$. Each of these updates has a very simple familiar form. For example,

$$\begin{aligned} \mathbf{L}_{k+1} &= \arg \min_{\mathbf{L}} \mathcal{L}_\mu(\mathbf{L}, \mathbf{S}_k, \mathbf{\Lambda}_k) \\ &= \arg \min_{\mathbf{L}} \|\mathbf{L}\|_* + \langle \mathbf{\Lambda}_k, \mathbf{Y} - \mathbf{L} - \mathbf{S}_k \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}_k\|_F^2 \\ &= \arg \min_{\mathbf{L}} \|\mathbf{L}\|_* + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{S}_k + \mu^{-1} \mathbf{\Lambda}_k - \mathbf{L}\|_F^2 + \varphi(\mathbf{S}_k, \mathbf{\Lambda}_k) \\ &= \text{prox}_{\mu^{-1} \|\cdot\|_*}(\mathbf{Y} - \mathbf{S}_k + \mu^{-1} \mathbf{\Lambda}_k). \end{aligned} \quad (8.5.10)$$

Thus, the update step for the low-rank term can be evaluated simply by computing the proximal operator for the nuclear norm.

A similar simple rule can be derived for the sparse term:

$$\begin{aligned} \mathbf{S}_{k+1} &= \arg \min_{\mathbf{S}} \mathcal{L}_\mu(\mathbf{L}_{k+1}, \mathbf{S}, \mathbf{\Lambda}_k) \\ &= \arg \min_{\mathbf{S}} \lambda \|\mathbf{S}\|_1 + \langle \mathbf{\Lambda}_k, \mathbf{Y} - \mathbf{L}_{k+1} - \mathbf{S} \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L}_{k+1} - \mathbf{S}\|_F^2 \\ &= \arg \min_{\mathbf{S}} \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L}_{k+1} + \mu^{-1} \mathbf{\Lambda}_k - \mathbf{S}\|_F^2 + \varphi(\mathbf{L}_{k+1}, \mathbf{\Lambda}_k) \\ &= \text{prox}_{\lambda \mu^{-1} \|\cdot\|_1}(\mathbf{Y} - \mathbf{L}_{k+1} + \mu^{-1} \mathbf{\Lambda}_k). \end{aligned} \quad (8.5.11)$$

Combining these two observations, we obtain a simple, lightweight algorithm for solving the Principal Component Pursuit program.

Algorithm 8.7 ADMM for Principal Component Pursuit

- 1: **Problem:** $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \langle \mathbf{\Lambda}, \mathbf{Y} - \mathbf{L} - \mathbf{S} \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2$,
given \mathbf{Y} , $\lambda, \mu > 0$.
 - 2: **Input:** $\mathbf{L}_0 \in \mathbb{R}^{m \times n}$, $\mathbf{S}_0 \in \mathbb{R}^{m \times n}$.
 - 3: **while** $((\mathbf{L}_k, \mathbf{S}_k)$ not converged $(k = 1, 2, \dots)$ **do**
 - 4: $\mathbf{L}_{k+1} \leftarrow \text{prox}_{\mu^{-1} \|\cdot\|_*}(\mathbf{Y} - \mathbf{S}_k + \mu^{-1} \mathbf{\Lambda}_k)$.
 - 5: $\mathbf{S}_{k+1} \leftarrow \text{prox}_{\lambda \mu^{-1} \|\cdot\|_1}(\mathbf{Y} - \mathbf{L}_{k+1} + \mu^{-1} \mathbf{\Lambda}_k)$.
 - 6: $\mathbf{\Lambda}_{k+1} \leftarrow \mathbf{\Lambda}_k + \mu(\mathbf{Y} - \mathbf{L}_{k+1} - \mathbf{S}_{k+1})$.
 - 7: **end while**
 - 8: **Output:** $\mathbf{L}_* \leftarrow \mathbf{L}_k$; $\mathbf{S}_* \leftarrow \mathbf{S}_k$.
-

8.5.2 Convergence of ADMM

There has been a rich history and literature on characterizing the convergence and convergence rates of the ADMM algorithm under various conditions [Deng and Yin, 2016]. The ADMM can be naturally viewed as an approximation to the classical ALM method studied in the previous section: In the case the objective function is separable, one uses a single pass of “Gauss-Seidel” block minimization to substitute for full minimization of the augmented Lagrangian in each iteration (8.4.11). However, as pointed out in [Eckstein, 2012], this interpretation does not seem to lead to any known convergence proof for the ADMM that is closely related the convergence proof for ALM given in Section 8.4.3. Proofs for strong convergence of ADMM under general assumptions and for many of its variants are considerably more involved than that for ALM.

In this section, we try to establish a convergence result of the ADMM algorithm for a very basic case that nevertheless is applicable to problems of interest to us, such as the principal component pursuit. Our proof follows the approach of [Deng and Yin, 2016]. In our analysis below, we will use the following monotonicity property of convex functions.

Lemma 8.5.3 (Monotonicity Property). *Given a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ and any $\mathbf{x}, \mathbf{x}', \mathbf{v}, \mathbf{v}' \in \mathbb{R}^n$ such that $\mathbf{v} \in \partial f(\mathbf{x})$ and $\mathbf{v}' \in \partial f(\mathbf{x}')$, we have*

$$\langle \mathbf{x} - \mathbf{x}', \mathbf{v} - \mathbf{v}' \rangle \geq 0. \quad (8.5.12)$$

Proof. From the definition of subgradient, we have

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{x}' - \mathbf{x} \rangle, \quad f(\mathbf{x}) \geq f(\mathbf{x}') + \langle \mathbf{v}', \mathbf{x} - \mathbf{x}' \rangle. \quad (8.5.13)$$

Adding these two inequalities together we obtain:

$$f(\mathbf{x}) + f(\mathbf{x}') \geq f(\mathbf{x}) + f(\mathbf{x}') + \langle \mathbf{v} - \mathbf{v}', \mathbf{x}' - \mathbf{x} \rangle. \quad (8.5.14)$$

Canceling $f(\mathbf{x}) + f(\mathbf{x}')$ from both sides obtains the desired result. \square

To simplify the proof of convergence, we will assume our problem satisfies all necessary nice properties. For instance, the functions $g(\mathbf{x})$ and $h(\mathbf{z})$ are both convex and have Lipschitz continuous gradient with constant L . We also assume that the problem (8.5.4) has a unique solution $(\mathbf{x}_*, \mathbf{z}_*, \boldsymbol{\lambda}_*)$ that satisfies the KKT conditions, or in other words, a saddle point of the Lagrangian (8.5.5):

$$\mathbf{A}^* \boldsymbol{\lambda}_* \in \partial g(\mathbf{x}_*), \quad \mathbf{B}^* \boldsymbol{\lambda}_* \in \partial h(\mathbf{z}_*), \quad \mathbf{A} \mathbf{x}_* + \mathbf{B} \mathbf{z}_* - \mathbf{y} = 0. \quad (8.5.15)$$

Theorem 8.5.4. *Under the above assumptions, the ADMM iterates (8.5.7) to (8.5.8) converge to the KKT solution $(\mathbf{x}_*, \mathbf{z}_*, \boldsymbol{\lambda}_*)$ in the following sense:*

$$\boldsymbol{\lambda}_k \rightarrow \boldsymbol{\lambda}_*, \quad \mathbf{A} \mathbf{x}_k \rightarrow \mathbf{A} \mathbf{x}_*, \quad \mathbf{B} \mathbf{z}_k \rightarrow \mathbf{B} \mathbf{z}_*. \quad (8.5.16)$$

In the case when \mathbf{A} and \mathbf{B} both have full column rank, we have $\mathbf{x}_k \rightarrow \mathbf{x}_\star$ and $\mathbf{z}_k \rightarrow \mathbf{z}_\star$.

Proof. Recall from (8.5.8) that we have

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu(\mathbf{y} - \mathbf{A}\mathbf{x}_{k+1} - \mathbf{B}\mathbf{z}_{k+1}). \quad (8.5.17)$$

Combining this with the conditions for the KKT solution (8.5.15), we have

$$\mathbf{A}(\mathbf{x}_{k+1} - \mathbf{x}_\star) + \mathbf{B}(\mathbf{z}_{k+1} - \mathbf{z}_\star) = \frac{1}{\mu}(\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1}). \quad (8.5.18)$$

From the optimality of the subproblems (8.5.7) and (8.5.6), we have

$$\mathbf{A}^* \boldsymbol{\lambda}_{k+1} \in \partial g(\mathbf{x}_{k+1}), \quad (8.5.19)$$

$$\mathbf{B}^*(\boldsymbol{\lambda}_{k+1} - \mu \mathbf{A}(\mathbf{x}_k - \mathbf{x}_{k+1})) \in \partial h(\mathbf{z}_{k+1}). \quad (8.5.20)$$

Applying the Lemma 8.5.3 the iterates $(\mathbf{x}_{k+1}, \mathbf{z}_{k+1})$ and the KKT solution $(\mathbf{x}_\star, \mathbf{z}_\star)$, we get:

$$\langle \mathbf{x}_{k+1} - \mathbf{x}_\star, \mathbf{A}^*(\boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}_\star) \rangle \geq 0, \quad (8.5.21)$$

$$\langle \mathbf{z}_{k+1} - \mathbf{z}_\star, \mathbf{B}^*(\boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}_\star - \mu \mathbf{A}(\mathbf{x}_k - \mathbf{x}_{k+1})) \rangle \geq 0. \quad (8.5.22)$$

Adding the above two inequalities and using (8.5.18), we get

$$\begin{aligned} & \frac{1}{\mu} \langle \boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}_\star, \boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1} \rangle + \langle \mathbf{x}_{k+1} - \mathbf{x}_\star, \mu \mathbf{A}^* \mathbf{A}(\mathbf{x}_k - \mathbf{x}_{k+1}) \rangle \\ & \geq \langle \mathbf{A}(\mathbf{x}_k - \mathbf{x}_{k+1}), \boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1} \rangle \geq 0. \end{aligned}$$

In the last inequality, we have used the fact that $\mathbf{A}^* \boldsymbol{\lambda}_k$ is the subgradient of $g(\mathbf{x})$ at \mathbf{x}_k from equation (8.5.19) and the monotonicity Lemma 8.5.3. Define $\|\mathbf{x}\|_Q^2 = \mathbf{x}^* \mathbf{Q} \mathbf{x}$ and the above inequality leads to:

$$\begin{aligned} & \frac{1}{\mu} \langle \boldsymbol{\lambda}_k - \boldsymbol{\lambda}_\star, \boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1} \rangle + \langle \mathbf{x}_k - \mathbf{x}_\star, \mu \mathbf{A}^* \mathbf{A}(\mathbf{x}_k - \mathbf{x}_{k+1}) \rangle \\ & \geq \frac{1}{\mu} \|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1}\|^2 + \mu \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_{\mathbf{A}^* \mathbf{A}}^2. \end{aligned}$$

Notice that we have the identity

$$\|\mathbf{a} - \mathbf{c}\|_Q^2 - \|\mathbf{b} - \mathbf{c}\|_Q^2 = 2(\mathbf{a} - \mathbf{c})^* \mathbf{Q}(\mathbf{a} - \mathbf{b}) - \|\mathbf{a} - \mathbf{b}\|_Q^2.$$

Let \mathbf{a} be the k -th iterate, \mathbf{b} be the $(k+1)$ -th iterate, and \mathbf{c} be the KKT solution. Substitute the left side of the inequality into the right side of this identity, we obtain:

$$\begin{aligned} & \frac{1}{\mu} \|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_\star\|^2 + \mu \|\mathbf{x}_k - \mathbf{x}_\star\|_{\mathbf{A}^* \mathbf{A}}^2 - \left(\frac{1}{\mu} \|\boldsymbol{\lambda}_{k+1} - \boldsymbol{\lambda}_\star\|^2 + \mu \|\mathbf{x}_{k+1} - \mathbf{x}_\star\|_{\mathbf{A}^* \mathbf{A}}^2 \right) \\ & \geq \frac{1}{\mu} \|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1}\|^2 + \mu \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_{\mathbf{A}^* \mathbf{A}}^2. \end{aligned}$$

From the above inequality, we see that the sequence

$$\frac{1}{\mu} \|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_\star\|^2 + \mu \|\mathbf{x}_k - \mathbf{x}_\star\|_{\mathbf{A}^* \mathbf{A}}^2$$

is monotonically non-increasing and thus must converge. This implies that the right hand side must converge to zero:

$$\frac{1}{\mu} \|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1}\|^2 + \mu \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_{\mathbf{A}^* \mathbf{A}}^2 \rightarrow 0. \quad (8.5.23)$$

Therefore, $\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_{k+1} \rightarrow 0$ and $\mathbf{A}\mathbf{x}_k - \mathbf{A}\mathbf{x}_{k+1} \rightarrow 0$.

It is easy to check that any converging subsequence of $(\mathbf{x}_k, \mathbf{z}_k, \boldsymbol{\lambda}_k)$ converges to a point $(\bar{\mathbf{x}}, \bar{\mathbf{z}}, \bar{\boldsymbol{\lambda}})$ that satisfies the KKT condition. As the solution that satisfies the KKT condition is unique, we have

$$\boldsymbol{\lambda}_k \rightarrow \boldsymbol{\lambda}_\star, \quad \mathbf{A}\mathbf{x}_k \rightarrow \mathbf{A}\mathbf{x}_\star, \quad \mathbf{B}\mathbf{z}_k \rightarrow \mathbf{B}\mathbf{z}_\star.$$

Finally, when the matrices \mathbf{A} and \mathbf{B} both have full column rank, such as the case with principal component pursuit, we have

$$\mathbf{x}_k \rightarrow \mathbf{x}_\star, \quad \mathbf{z}_k \rightarrow \mathbf{z}_\star.$$

□

The above theorem is not the most general and strongest result that one can obtain for the convergence of ADMM. We recommend [Eckstein, 2012] for a tutorial introduction to a more formal convergence analysis of ADMM. For more recent analysis of generalized version of ADMM including their convergence rates, we recommend the work of [Deng and Yin, 2016]. Like the ALM algorithm in the previous section, the most natural approach is to recognize ADMM as some known algorithm applied to the dual. In fact, ADMM turns out to be equivalent to Douglas-Rachford splitting applied to the dual. For more details on this, see [Eckstein and Bertsekas, 1992, Combettes and Wajs, 2005]. ADMM has also been widely applied to problems when the number of terms are more than three. The convergence analysis of ADMM with more than three terms is much more difficult and it has been shown to diverge in many cases.

8.6 Frank-Wolfe Methods for Scalable Optimization

In the previous sections, we showed how the special structure of optimization problems arising in sparse and low-dimensional data analysis can be leveraged to obtain efficient and scalable algorithms. One key piece of structure was the existence of an easy-to-compute proximal operator. For example, for nuclear norm minimization we showed that at a point

$$\mathbf{Z} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*,$$

$$\text{prox}_{\gamma\|\cdot\|_*}(\mathbf{Z}) = \mathbf{U}\mathcal{S}_\lambda[\mathbf{\Sigma}]\mathbf{V}^*, \quad (8.6.1)$$

where $\mathcal{S}_\lambda[\cdot]$ is the soft thresholding operator. Using the proximal operator, we obtain proximal gradient methods that enjoy the same convergence rate as if the objective were smooth, even though it is nonsmooth. Each iteration consists of simple linear operations, followed by the application of $\text{prox}_{\gamma\|\cdot\|_*}$. Each iteration can be computed in time polynomial in the size of the target matrix: the proximal operator can be computed in time $O(n_1 n_2 \max\{n_1, n_2\})$ in the worst case. This is sufficient for moderate-sized datasets where n_1 and n_2 are each in the thousands.

However, many problems in data science and imaging require even more scalable solutions. In this section, we introduce a classical method from optimization, known as the *Frank-Wolfe* or *conditional gradient* algorithm, which is scalable enough to solve extremely large sparse and low-rank recovery problems. The key property of this method is that in each iteration, it solves a subproblem which is simpler and easier to compute than the proximal operator.

In its classical form, the Frank-Wolfe algorithm applies to the problem of optimizing a smooth, convex function over a compact convex set:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}), \\ & \text{subject to} && \mathbf{x} \in C. \end{aligned} \quad (8.6.2)$$

Here, the objective function f is assumed to be a convex, differentiable function whose gradient $\nabla f(\mathbf{x})$ is L -Lipschitz. The constraint set C is assumed to be a closed, bounded convex set with diameter

$$\text{diam}(C) \doteq \max\{\|\mathbf{x} - \mathbf{x}'\|_2 \mid \mathbf{x}, \mathbf{x}' \in C\}. \quad (8.6.3)$$

Constrained formulations of sparse and low-rank recovery.

Many of the sparse and low-rank recovery problems that we have considered thus far can be reformulated in terms of (8.6.2). For example, for sparse recovery, we can choose $C = \{\mathbf{x} \mid \|\mathbf{x}\|_1 \leq \tau\}$ to be an ℓ^1 ball, and solve

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2, \\ & \text{subject to} && \|\mathbf{x}\|_1 \leq \tau. \end{aligned} \quad (8.6.4)$$

Similarly, for low-rank matrix completion, we can choose C to be a nuclear norm ball and solve

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathcal{P}_\Omega[\mathbf{X}] - \mathbf{Y}\|_F^2, \\ & \text{subject to} && \|\mathbf{X}\|_* \leq \tau. \end{aligned} \quad (8.6.5)$$

Exercises 8.6 – 8.7 explore further reformulations of unconstrained sparse and low-rank optimization in the form (8.6.2).

The Frank-Wolfe Method.

Similar to the other methods we have discussed thus far, Frank-Wolfe is an iterative method, which generates a sequence of iterates $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k, \dots$ as follows. At each iteration, we generate a new point \mathbf{v}_k by solving a constrained optimization problem

$$\mathbf{v}_k \in \arg \min_{\mathbf{v} \in C} \langle \mathbf{v}, \nabla f(\mathbf{x}_k) \rangle. \quad (8.6.6)$$

We then set

$$\mathbf{x}_{k+1} = (1 - \gamma_k)\mathbf{x}_k + \gamma_k\mathbf{v}_k \in C, \quad (8.6.7)$$

where $\gamma_k \in (0, 1)$ is a specially chosen step size. Figure 8.5 summarizes the properties of this method.

Interpretation as Minimizing a First-Order Approximation.

The Frank-Wolfe method can be interpreted as follows. At a given point \mathbf{x}_k , we form a first order approximation to the objective function f :

$$f(\mathbf{v}) \approx \hat{f}(\mathbf{v}, \mathbf{x}_k) \doteq f(\mathbf{x}_k) + \langle \mathbf{v} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle. \quad (8.6.8)$$

We minimize the approximation $\hat{f}(\mathbf{v}, \mathbf{x}_k)$ over $\mathbf{v} \in C$ to produce \mathbf{v}_k . We then take a step in the direction $\mathbf{w}_k = \mathbf{v}_k - \mathbf{x}_k$:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \gamma_k\mathbf{w}_k. \quad (8.6.9)$$

Computing the Step Direction.

The crucial subproblem in the Frank-Wolfe method involves minimizing a linear function over a compact convex set C :

$$\min_{\mathbf{v} \in C} \langle \mathbf{v}, \nabla f(\mathbf{x}) \rangle. \quad (8.6.10)$$

Depending on the constraint set C , this could itself be a challenging (or even intractable!) optimization problem. Fortunately, for the problems of interest in this book, this subproblem can be solved in an efficient and scalable manner. We give two examples below:

Example 8.6.1 (Frank-Wolfe Subproblem over an ℓ^1 Ball). *Consider the problem*

$$\min_{\|\mathbf{v}\|_1 \leq \tau} \langle \mathbf{v}, \mathbf{g} \rangle. \quad (8.6.11)$$

Let i be any index for which $\mathbf{g}_i = \|\mathbf{g}\|_\infty$, and $\sigma_i = \text{sign}(\mathbf{g}_i)$. Then (8.6.11) has a solution

$$\mathbf{v}_\star = -\tau\sigma_i\mathbf{e}_i, \quad (8.6.12)$$

where \mathbf{e}_i is the i -th standard basis vector. The solution \mathbf{v}_\star can be computed in linear time, simply by finding the largest magnitude entry of \mathbf{g} .

Example 8.6.2 (Frank-Wolfe Subproblem of a Nuclear Norm Ball). Consider the problem

$$\min_{\|\mathbf{V}\|_* \leq \tau} \langle \mathbf{V}, \mathbf{G} \rangle. \quad (8.6.13)$$

Let $\mathbf{G} = \mathbf{U}\Sigma\mathbf{V}^* = \sum_{i=1}^{n_1} \mathbf{u}_i \sigma_i \mathbf{v}_i$ denote the singular value decomposition of \mathbf{G} . Then (8.6.13) has an optimal solution

$$\mathbf{V}_* = -\tau \mathbf{u}_1 \mathbf{v}_1^*. \quad (8.6.14)$$

This optimal solution can be computed in time $O(n_1 n_2)$ by computing (only) the lead singular vector pair $(\mathbf{u}_1, \mathbf{v}_1)$ of \mathbf{G} .

The latter example illustrates the special virtue of the Frank-Wolfe method in nuclear norm minimization: the key subproblem only requires us to compute *one* singular value/vector triple. For a problem involving $n_1 \times n_2$ matrices, this can be done in time $O(n_1 n_2)$ – a dramatic improvement over proximal gradient methods, which require a full singular value decomposition in each iteration.⁶

This scalability comes at a price, though. Compared to accelerated proximal gradient methods, which converge at a rate of $O(1/k^2)$ in function values, the Frank-Wolfe method only achieves a rate of $O(1/k)$. The following theorem gives a precise bound on the worst case rate of convergence for Frank-Wolfe over the class of convex functions with Lipschitz gradient.

Theorem 8.6.3. Let $\mathbf{x}_0, \mathbf{x}_1, \dots$ denote the sequence of iterates generated by the Frank-Wolfe method, with step size $\gamma_k = \frac{2}{k+2}$. Then

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{2L \text{diam}^2(C)}{k+2}. \quad (8.6.15)$$

We prove this theorem in the next subsection. In the context of nuclear norm minimization, this result can be viewed as follows: Frank-Wolfe allows us to derive methods that produce moderate-quality solutions to extremely large problems, for which methods with better worst cases rates simply take too long to compute even a single iterate.

8.6.1 Convergence of Frank-Wolfe

In this section, we prove Theorem 8.6.3 on the convergence rate of the Frank-Wolfe method.

Proof. For ease of notation, write $d = \text{diam}(C)^2$, $\mathbf{x} = \mathbf{x}_k$, $\mathbf{x}^+ = \mathbf{x}_{k+1}$, $\gamma = \gamma_k$, and $\mathbf{v} = \mathbf{v}_k$. Note that

$$\mathbf{x}^+ - \mathbf{x} = \gamma(\mathbf{v} - \mathbf{x}). \quad (8.6.16)$$

⁶See Exercise ?? for a different attempt to address this issue in proximal gradient methods.

Frank-Wolfe Method (FW)**Problem Class:**

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}), \\ & \text{subject to} && \mathbf{x} \in C. \end{aligned}$$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ convex, differentiable, $\nabla f(\mathbf{x})$ L -Lipschitz.
 C a compact convex set.

Basic Iteration: Repeat

$$\begin{aligned} \mathbf{v}_k & \in \arg \min_{\mathbf{v} \in C} \langle \mathbf{v}, \nabla f(\mathbf{x}_k) \rangle, \\ \mathbf{x}_{k+1} & = \mathbf{x}_k + \gamma_k (\mathbf{v}_k - \mathbf{x}_k), \end{aligned}$$

with $\gamma_k = \frac{2}{k+2}$.

Convergence Guarantee:

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{2L \text{diam}^2(C)}{k+2}.$$

Figure 8.5. Properties of the Frank-Wolfe Method.

Because $\nabla f(\mathbf{x})$ is L -Lipschitz, we can use the upper bound (8.2.8) to obtain

$$\begin{aligned} f(\mathbf{x}^+) & \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}\|_2^2 \\ & \leq f(\mathbf{x}) + \gamma \langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle + \frac{\gamma^2 L}{2} \|\mathbf{v} - \mathbf{x}\|_2^2 \\ & \leq f(\mathbf{x}) + \gamma \langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle + \frac{\gamma^2 L}{2} d^2. \end{aligned} \quad (8.6.17)$$

Meanwhile, by convexity,

$$f(\mathbf{x}_*) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}_* - \mathbf{x} \rangle \quad (8.6.18)$$

$$\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle, \quad (8.6.19)$$

where the final line follows because \mathbf{v} is chosen to minimize $\langle \nabla f(\mathbf{x}), \mathbf{v} \rangle$. Combining these two inequalities, we find that

$$\langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle \leq -\left(f(\mathbf{x}) - f(\mathbf{x}_*)\right), \quad (8.6.20)$$

whence, plugging into (8.6.17) and subtracting $f(\mathbf{x}_*)$ from both sides, we obtain

$$f(\mathbf{x}^+) - f(\mathbf{x}_*) \leq (1 - \gamma) \left(f(\mathbf{x}) - f(\mathbf{x}_*)\right) + \frac{\gamma^2}{2} L d^2. \quad (8.6.21)$$

We use this basic relationship together with an inductive argument to bound the rate of convergence of the Frank-Wolfe method. Let ε_k denote

the suboptimality (in function values) at iteration k :

$$\varepsilon_k = f(\mathbf{x}_k) - f(\mathbf{x}_*). \quad (8.6.22)$$

Set $\gamma_k = \frac{2}{k+2}$, so $\gamma_0 = 1$. Applying (8.6.21), we find that

$$\varepsilon_1 \leq \frac{1}{2} L d^2. \quad (8.6.23)$$

Suppose now that for $\ell = 1, \dots, k$, $\varepsilon_\ell \leq \frac{2}{\ell+2} L d^2$. Applying (8.6.21) again, we find that

$$\begin{aligned} \varepsilon_{k+1} &\leq \frac{k}{k+2} \varepsilon_k + \frac{2}{(k+2)^2} L d^2 \\ &\leq \frac{k+1}{(k+2)^2} \times 2 L d^2 \\ &\leq \frac{2 L d^2}{(k+1)+2}. \end{aligned} \quad (8.6.24)$$

Hence, the relationship $\varepsilon_\ell \leq \frac{2}{\ell+2} L d^2$ holds for all iterations ℓ , as claimed. \square

8.6.2 Frank-Wolfe for Stable Matrix Completion

We illustrate the general Frank-Wolfe method for the particular problem of recovering a low-rank matrix from incomplete and noisy observations

$$\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{X}_{\text{true}} + \mathbf{Z}], \quad (8.6.25)$$

where $\mathbf{X}_{\text{true}} \in \mathbb{R}^{n_1 \times n_2}$ has low rank, $\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}$ is a matrix of small, dense noise, and $\Omega \subseteq [n_1] \times [n_2]$ is the set of observed entries. One approach to approximately recovering \mathbf{X}_{true} is to minimize the reconstruction error over the set of all matrices of small nuclear norm:

$$\begin{aligned} &\text{minimize} && f(\mathbf{X}) \equiv \frac{1}{2} \|\mathcal{P}_\Omega[\mathbf{X}] - \mathbf{Y}\|_F^2, \\ &\text{subject to} && \|\mathbf{X}\|_* \leq \tau. \end{aligned} \quad (8.6.26)$$

Here, the constraint $\|\mathbf{X}\|_* \leq \tau$ encourages \mathbf{X} to have low rank. The constraint set $C = \{\mathbf{X} \mid \|\mathbf{X}\|_* \leq \tau\}$ is closed and bounded. Moreover, the gradient

$$\nabla f(\mathbf{X}) = \mathcal{P}_\Omega[\mathbf{X} - \mathbf{Y}] \quad (8.6.27)$$

is 1-Lipschitz, and so the Frank-Wolfe method indeed applies to this problem.

The key step in the Frank-Wolfe method is to minimize a linear function $\langle \mathbf{V}, \nabla f(\mathbf{X}) \rangle$ over the constraint set C . As described above, this problem can be solved in closed form: if

$$\nabla f(\mathbf{X}) = \sum_{i=1}^{n_1} \mathbf{u}_i \sigma_i \mathbf{v}_i^* \quad (8.6.28)$$

is the singular value decomposition of f , then

$$-\tau \mathbf{u}_1 \mathbf{v}_1 \in \arg \min_{\mathbf{V} \in C} \langle \mathbf{V}, \nabla f(\mathbf{X}) \rangle. \quad (8.6.29)$$

The lead singular value / vectors can be extracted from the matrix $\nabla f(\mathbf{X})$ efficiently, without computing the entire SVD (8.6.28). Typically, this is done using the power method, which was described in some detail in Chapter 4 and Exercise 4.6. To cleanly describe the method, we simply let

$$(\mathbf{u}_1, \sigma_1, \mathbf{v}_1) = \text{LeadSV}(\mathbf{G}) \quad (8.6.30)$$

denote the operation which extracts a leading singular value / vector triple from a matrix \mathbf{G} . Using this notation, the complete Frank-Wolfe algorithm for stable matrix completion is describe in Algorithm 8.8.

Algorithm 8.8 Frank-Wolfe for Stable Matrix Completion

- 1: **Problem:** $\min_{\mathbf{X}} \frac{1}{2} \|\mathcal{P}_\Omega[\mathbf{X}] - \mathbf{Y}\|_F^2$ subj. to $\|\mathbf{X}\|_* \leq \tau$, given $\mathbf{Y} \in \mathbb{R}^{n_1 \times n_2}$ and $\Omega \subseteq [n_1] \times [n_2]$.
 - 2: **Input:** $\mathbf{X}_0 \in \mathbb{R}^{n_1 \times n_2}$ satisfying $\|\mathbf{X}_0\|_* \leq \tau$.
 - 3: **while** \mathbf{X}_k not converged ($k = 1, 2, \dots$) **do**
 - 4: $(\mathbf{u}_1, \sigma_1, \mathbf{v}_1) \leftarrow \text{LeadSV}(\mathcal{P}_\Omega[\mathbf{X}_k - \mathbf{Y}])$.
 - 5: $\mathbf{V}_k \leftarrow -\tau \mathbf{u}_1 \mathbf{v}_1^*$.
 - 6: $\mathbf{X}_{k+1} \leftarrow \frac{k}{k+2} \mathbf{X}_k + \frac{2}{k+2} \mathbf{V}_k$.
 - 7: **end while**
 - 8: **Output:** $\mathbf{X}_\star \leftarrow \mathbf{X}_k$.
-

The Frank-Wolfe method requires only a single singular value/vector triple at each iteration. Moreover, since $\mathbf{V}_k = -\tau \mathbf{u}_1 \mathbf{v}_1^*$ has rank one, the rank of \mathbf{X}_k increases by at most one at each iteration. In this sense, Frank-Wolfe can be viewed as a *greedy method*. It constructs a low-rank matrix \mathbf{X}_\star by adding on one (optimally chosen) rank-one factor at a time.

8.6.3 Connection to Greedy Methods for Sparsity

In sparse and low-rank approximation, *greedy methods* are sometimes favored for their simplicity and scalability. For sparse approximation, the Frank-Wolfe method gives one such greedy algorithm. Consider the problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2, \\ &\text{subject to} && \|\mathbf{x}\|_1 \leq \tau. \end{aligned} \quad (8.6.31)$$

Notice that

$$\nabla f(\mathbf{x}) = \mathbf{A}^*(\mathbf{A}\mathbf{x} - \mathbf{y}). \quad (8.6.32)$$

The Frank-Wolfe subproblem

$$\min_{\|\mathbf{v}\|_1 \leq \tau} \langle \mathbf{v}, \nabla f(\mathbf{x}) \rangle \quad (8.6.33)$$

has an especially simple solution: letting i be the index of the largest magnitude entry of ∇f , and σ its sign,

$$\mathbf{v}_* = -\tau \sigma \mathbf{e}_i. \quad (8.6.34)$$

Algorithm 8.9 Frank-Wolfe for Noisy Sparse Recovery

- 1: **Problem:** $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$ subj. to $\|\mathbf{x}\|_1 \leq \tau$, given $\mathbf{y} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$.
 - 2: **Input:** $\mathbf{x}_0 \in \mathbb{R}^n$ satisfying $\|\mathbf{x}_0\|_1 \leq \tau$.
 - 3: **while** \mathbf{x}_k not converged ($k = 1, 2, \dots$) **do**
 - 4: $i_k \leftarrow \arg \max_i |[\mathbf{A}^*(\mathbf{A}\mathbf{x}_k - \mathbf{y})]_i|$.
 - 5: $\sigma \leftarrow \text{sign}([\mathbf{A}^*(\mathbf{A}\mathbf{x}_k - \mathbf{y})]_{i_k})$.
 - 6: $\mathbf{v}_k \leftarrow -\tau \sigma \mathbf{e}_{i_k}$.
 - 7: $\mathbf{x}_{k+1} \leftarrow \frac{k}{k+2} \mathbf{x}_k + \frac{2}{k+2} \mathbf{v}_k$.
 - 8: **end while**
 - 9: **Output:** $\mathbf{x}_* \leftarrow \mathbf{x}_k$.
-

Algorithm 8.9 describes in detail the Frank-Wolfe method for problem (8.6.31). At each iteration, it increases the number of nonzero entries in the vector \mathbf{x} by at most one, by adding on a multiple of \mathbf{e}_{i_k} . Let

$$I_k = \{i_1, \dots, i_{k-1}\} = \text{supp}(\mathbf{x}_k) \quad (8.6.35)$$

denote the collection of indices that have been chosen up to time k . We generate I_{k+1} from I_k by introducing a (potentially) new index

$$I_{k+1} = I_k \cup \{i_k\}. \quad (8.6.36)$$

This new index is chosen according to the largest-magnitude entry in the gradient ∇f . Write

$$\mathbf{A} = [\mathbf{a}_1 \mid \dots \mid \mathbf{a}_n] \quad (8.6.37)$$

for the columns of \mathbf{A} , and let

$$\mathbf{r}_k = \mathbf{A}\mathbf{x}_k - \mathbf{y} \quad (8.6.38)$$

denote the measurement residual at point \mathbf{x}_k . Since $\nabla f(\mathbf{x}_k) = \mathbf{A}^* \mathbf{r}_k$, the Frank-Wolfe method chooses the index i_k corresponding to the column \mathbf{a}_{i_k} that is most correlated with the residual \mathbf{r}_k .

Matching Pursuit.

A number of classical *greedy methods* for sparse approximation have this basic structure. A canonical example is the *Matching Pursuit* algorithm.

Algorithm 8.10 Matching Pursuit for Sparse Approximation

```

1: Problem: find a sparse  $\mathbf{x}$  such that  $f(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2$  is small.
2:  $\mathbf{x}_0 \leftarrow \mathbf{0}$ .
3: while  $\mathbf{x}_k$  not converged ( $k = 1, 2, \dots$ ) do
4:    $i_k \leftarrow \arg \max_i |[\mathbf{A}^*(\mathbf{Ax}_k - \mathbf{y})]_i|$ .
5:    $t_k \leftarrow -\frac{\langle \mathbf{a}_{i_k}, \mathbf{r}_k \rangle}{\|\mathbf{a}_{i_k}\|_2^2}$ .
6:    $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + t_k \mathbf{e}_{i_k}$ .
7: end while
8: Output:  $\mathbf{x}_\star \leftarrow \mathbf{x}_k$ .

```

This algorithm generates a sequence of iterates $\mathbf{x}_0 = \mathbf{0}, \mathbf{x}_1, \mathbf{x}_2, \dots$, by repeatedly choosing a column of \mathbf{a}_{i_k} of \mathbf{A} that is most correlated with the residual \mathbf{r}_k . Similar to Frank-Wolfe, Matching Pursuit⁷ sets

$$i_k = \arg \max_i |[\nabla f(\mathbf{x}_k)]_i| = \arg \max_i |\mathbf{a}_i^* \mathbf{r}_k|. \quad (8.6.39)$$

However, rather than stepping a predetermined length in the \mathbf{e}_{i_k} direction, Matching Pursuit chooses the step size t_k by solving a one-dimensional minimization problem:

$$t_k = \arg \min_t f(\mathbf{x}_k + t \mathbf{e}_{i_k}) = -\frac{\langle \mathbf{a}_{i_k}, \mathbf{r}_k \rangle}{\|\mathbf{a}_{i_k}\|_2^2}. \quad (8.6.40)$$

This can be viewed as a form of *exact line search* and typically leads to more rapid convergence in practice. The overall Matching Pursuit algorithm is specified as Algorithm 8.10.

Orthogonal Matching Pursuit.

Matching pursuit achieves better convergence by choosing the step size t_k in an optimal manner. Since

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{e}_{i_k}, \quad (8.6.41)$$

this is equivalent to making an optimal choice of the i_k -th entry in \mathbf{x}_{k+1} , while leaving all of the other entries fixed. It is possible to further improve the rate of convergence of this approach by choosing *all* of the nonzero entries of \mathbf{x}_{k+1} optimally (rather than just the i_k -th entry). In notation, let $I_k = \{i_1, i_2, \dots, i_{k-1}\}$ denote the collection of indices that have been chosen up to step k . The *Orthogonal Matching Pursuit* method selects an

⁷Despite the strong parallels to the Frank-Wolfe method, Matching Pursuit was originally motivated from completely different principles. Its introduction to the literature on sparse approximation predates the use of the Frank-Wolfe method; see the Notes and References to this chapter for more historical details and connections to classical greedy algorithms in other fields.

Algorithm 8.11 Orthogonal Matching Pursuit for Sparse Approximation

```

1: Problem: find a sparse  $\mathbf{x}$  such that  $f(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2$  is small.
2:  $\mathbf{x}_0 \leftarrow \mathbf{0}$ ,  $I_0 \leftarrow \emptyset$ .
3: while  $\mathbf{x}_k$  not converged ( $k = 1, 2, \dots$ ) do
4:    $\mathbf{r}_k \leftarrow \mathbf{Ax}_k - \mathbf{y}$ .
5:    $i_k \leftarrow \arg \max_i |\mathbf{a}_i^* \mathbf{r}_k|$ .
6:    $I_{k+1} \leftarrow I_k \cup \{i_k\}$ .
7:    $[\mathbf{x}_{k+1}]_{I_{k+1}} \leftarrow (\mathbf{A}_{I_{k+1}}^* \mathbf{A}_{I_{k+1}})^{-1} \mathbf{A}_{I_{k+1}}^* \mathbf{y}$ .
8:    $[\mathbf{x}_{k+1}]_{I_{k+1}^c} \leftarrow \mathbf{0}$ .
9: end while
10: Output:  $\mathbf{x}_\star \leftarrow \mathbf{x}_k$ .

```

index i_k that maximizes the correlation $|\mathbf{a}_i^* \mathbf{r}_k|$ of a column of \mathbf{A} with the residual $\mathbf{r}_k = \mathbf{Ax}_k - \mathbf{y}$. It sets $I_{k+1} = I_k \cup \{i_k\}$, and then updates all of the nonzero entries in \mathbf{x} by setting

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 \quad \text{subject to} \quad \text{supp}(\mathbf{x}) \subseteq I_{k+1}. \quad (8.6.42)$$

This problem can be solved in closed form:

$$[\mathbf{x}_{k+1}]_{I_{k+1}} = (\mathbf{A}_{I_{k+1}}^* \mathbf{A}_{I_{k+1}})^{-1} \mathbf{A}_{I_{k+1}}^* \mathbf{y}, \quad (8.6.43)$$

$$[\mathbf{x}_{k+1}]_{I_{k+1}^c} = \mathbf{0}. \quad (8.6.44)$$

The name *Orthogonal Matching Pursuit* comes from the observation that the residual

$$\begin{aligned} \mathbf{r}_{k+1} &= \mathbf{Ax}_{k+1} - \mathbf{y} \\ &= \left(\mathbf{A}_{I_{k+1}} (\mathbf{A}_{I_{k+1}}^* \mathbf{A}_{I_{k+1}})^{-1} \mathbf{A}_{I_{k+1}}^* - \mathbf{I} \right) \mathbf{y} \end{aligned} \quad (8.6.45)$$

is orthogonal to the range $\text{range}(\mathbf{A}_{I_{k+1}})$ of the dictionary columns selected through the first $k+1$ iterations.

The overall Orthogonal Matching Pursuit algorithm is given as Algorithm 8.11. This method is sometimes favored by practitioners due to its simplicity, and the fact that it maintains an explicit *active set* I_k . The latter property is useful for problems in which the support of the sparse solution \mathbf{x}_\star is the object of interest.⁸

Although OMP has many variants and extensions, it was originally derived for the specific problem of finding sparse near-solutions to a linear system of equations $\mathbf{Ax} = \mathbf{y}$. Like ℓ^1 minimization, OMP is guaranteed to succeed whenever \mathbf{y} is generated from some sufficiently sparse \mathbf{x}_o and the

⁸For example, in the RF spectrum sensing discussed in Chapter 12, the goal is to determine which bands of the RF spectrum are occupied, in order to avoid interference. The specific energy levels within these bands are of secondary importance.

columns of \mathbf{A} are sufficiently spread in the high-dimensional space \mathbb{R}^m . In particular:

Theorem 8.6.4. *Suppose that $\mathbf{y} = \mathbf{A}\mathbf{x}_o$, with*

$$k = \|\mathbf{x}_o\|_0 \leq \frac{1}{2\mu(\mathbf{A})}. \quad (8.6.46)$$

Then after k iterations, the OMP algorithm terminates with $\mathbf{x}_k = \mathbf{x}_o$ and $I_k = \text{supp}(\mathbf{x}_o)$.

Exercise 8.8 guides the reader through a proof of Theorem 8.6.4. The key message of the proof is that under the conditions of the theorem, at each iteration ℓ the algorithm selects an index i_ℓ that belongs to the true support set $\text{supp}(\mathbf{x}_o)$.

The form of Theorem 8.6.4 can be directly compared to that of Theorem 3.2.3 of Chapter 3. These results imply that *both* OMP and ℓ^1 minimization recover \mathbf{x}_o whenever $\|\mathbf{x}_o\|_0 \leq 1/2\mu(\mathbf{A})$. Hence, at an intuitive level, both methods succeed whenever the target solution is sparse and the matrix \mathbf{A} is “nice”. However, as shown in Chapter 3, the incoherence condition requires \mathbf{x}_o to be extremely sparse. ℓ^1 minimization also recovers denser \mathbf{x}_o under the stronger condition that \mathbf{A} satisfies the Restricted Isometry Property $\delta_{Ck}(\mathbf{A}) < c$. While various improved analyses of OMP are available, the RIP is not sufficient for OMP to succeed. In this sense, convex relaxation achieves a better uniform guarantee. However, OMP can be modified to also guarantee sparse recovery under the RIP. The key ideas are to allow the algorithm to *remove* elements from the active set I_k at each iteration, and to add multiple elements. The resulting method, called COSAMP (Compressed Sampling Matching Pursuit) is described in more detail in Exercise 8.9. The large and varied literature on greedy algorithms also includes greedy methods for more general problems such as low-rank recovery; see the Notes and References to this chapter.

8.7 Notes and References

The name *basis pursuit* was first suggested by Chen and Donoho in their early work on compressive sensing [Chen, 1995, Chen et al., 2001].

More references about the proximity operator [Moreau, 1962, Combettes and Wajs, 2005].

In the literature, ISTA has been also studied under different names, such as forward-backward splitting [Combettes and Wajs, 2005], thresholded Landweber [Daubechies et al., 2004], and separable approximation (SpaRSA) [Wright et al., 2008].

Review the history of ADMM algorithms.

There are several important heuristics in the implementation of convex optimization algorithms, which are significant enough that practitioners shall pay close attention.

8.8 Exercises

8.1 (Proximal Operator for the Nuclear Norm). *Prove the first and the third assertions of Proposition 8.2.3.*

8.2 (Iterative Soft-Thresholding Algorithm for PCP). *Regarding the solving the stable principal component pursuit program using proximal gradient descent:*

- *Apply the proximal gradient method to the PCP program. Based on separability of the two non-smooth terms in the objective function, write down the corresponding proximal operators and the associated \mathbf{w}_1 and \mathbf{w}_2 . Justify the updates for \mathbf{L}_{k+1} and \mathbf{S}_{k+1} in Algorithm 8.7.*
- *Code a MATLAB function that implements the Iterative Soft-Thresholding Algorithm 8.7 for PCP.*

8.3 (Accelerated Proximal Gradient Algorithm for Robust PCA*). *Code a MATLAB function that implement the accelerated proximal gradient method for the stable PCP program.*

8.4 (Dual Feasibility Condition of ALM).

8.5 (Implementation: Augmented Lagrange Multiplier Algorithm for PCP).

8.6 (Unconstrained problems via Frank-Wolfe).

8.7 (Sparse and Low-Rank via Frank-Wolfe).

8.8 (Sparse recovery by Orthogonal Matching Pursuit). *The goal of this exercise is to prove Theorem 8.6.4, which shows that OMP correctly recovers any target sparse solution \mathbf{x}_o with $k = \|\mathbf{x}_o\|_0 \leq \frac{1}{2\mu(\mathbf{A})}$. Let $I = \text{supp}(\mathbf{x}_o)$.*

- *OMP selects a true support index in the first iteration. Let i_{\max} index a maximum-magnitude entry of \mathbf{x}_o , i.e., $\mathbf{x}_o(i_{\max}) = \|\mathbf{x}_o\|_\infty$. Using the incoherence of \mathbf{A} , argue that*

$$|\mathbf{a}_{i_{\max}}^* \mathbf{r}_0| \geq |\mathbf{a}_j^* \mathbf{r}_0| \quad \forall j \in I^c. \quad (8.8.1)$$

- *Argue by induction that OMP selects some $i_\ell \in I$ for every iteration $\ell = 0, \dots, k-1$.*
- *Using the fact that $\mathbf{r}_\ell \perp \text{span}(\mathbf{A}_{I_\ell})$, argue that OMP selects a new index $i_\ell \in I$ at each iteration $\ell = 0, \dots, k-1$. Conclude that OMP terminates with $\mathbf{x}_k = \mathbf{x}_o$ and $I_k = I$, as claimed.*

8.9 (Greedy Methods that Succeed Under RIP). *The Compressive Sampling Matching Pursuit (COSAMP) algorithm modifies OMP by adding and subtracting multiple indices from the active set I_ℓ at each iteration ℓ . This algorithm takes as input a target number of nonzero entries, s , and modifies OMP as follows:*

- *At each iteration ℓ , we let $I_{\ell+1/2} = I_\ell \cup$*
- *We then solve for $\mathbf{x}_{\ell+1/2}$ by*
- *We then*

Implement the COSAMP algorithm, and compare its breakdown in terms of sparsity level to OMP.

8.10 (ADMM for Basis Pursuit).

8.11 (Dual of Principal Component Pursuit). *Show that the dual program to the PCP program is*

$$\max_{\mathbf{\Lambda}} \text{trace}(\mathbf{Y}^* \mathbf{\Lambda}) \quad \text{subj. to} \quad J(\mathbf{\Lambda}) \leq 1,$$

where $\mathbf{\Lambda}$ is the matrix of Lagrange multipliers for the equality constraint $\mathbf{Y} = \mathbf{L} + \mathbf{S}$, and $J(\mathbf{\Lambda}) = \max(\|\mathbf{\Lambda}\|_2, \lambda^{-1} \|\mathbf{\Lambda}\|_\infty)$.