

Lecture notes on regularized, large-scale and distributed optimization

Clément W. Royer

M2 IASD & M2 MASH - 2023/2024

- The last version of these notes can be found at:
<https://www.lamsade.dauphine.fr/~croyer/ensdocs/RLD/LectureNotesOML-RLD.pdf>.
- Comments, typos, etc, can be sent to `clement.royer@lamsade.dauphine.fr`.
- **Major updates of the document**
 - 2023.12.12: Full version.
 - 2023.11.08: First version of the notes.

Foreword

The purpose of these lecture notes is to discuss several characteristics of optimization problems arising in machine learning, that preclude from applying (stochastic) gradient methods as seen in the previous parts of the course. We are particularly interested in the use of regularization to bring structure into an optimization problem, and in techniques that allow to deploy algorithms at scale. More precisely, the learning goals are the following:

- Understand the key principles of proximal algorithms, and why those are useful for regularized problems.
- Understand the purpose of regularization terms in data science problems, with a focus on sparsity-inducing regularizers.
- Understand the main principle behind coordinate descent methods, and their interest in a distributed environment.
- Apply duality to design algorithms for distributed and decentralized optimization.

Notations

- Scalars (i.e. reals) are denoted by lowercase letters: $a, b, c, \alpha, \beta, \gamma$.
- Vectors are denoted by **bold** lowercase letters: $\mathbf{a}, \mathbf{b}, \mathbf{c}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}$.
- Matrices are denoted by **bold** uppercase letters: $\mathbf{A}, \mathbf{B}, \mathbf{C}$.
- Sets are denoted by **bold** uppercase cursive letters : $\mathcal{A}, \mathcal{B}, \mathcal{C}$.
- The set of natural numbers (nonnegative integers) is denoted by \mathbb{N} ; the set of integers is denoted by \mathbb{Z} .
- The set of real numbers is denoted by \mathbb{R} . Our notations for the subset of nonnegative real numbers and the set of positive real numbers are \mathbb{R}_+ and \mathbb{R}_{++} , respectively.
- The notation \mathbb{R}^d is used for the set of vectors with $d \in \mathbb{N}$ real components; although we may not explicitly indicate it in the rest of these notes, we always assume that $d \geq 1$.
- A vector $\mathbf{x} \in \mathbb{R}^d$ is thought as a column vector, with $x_i \in \mathbb{R}$ denoting its i -th coordinate in the canonical basis of \mathbb{R}^d . We thus write $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$, or, in a compact form, $\mathbf{x} = [x_i]_{1 \leq i \leq d}$.
- Given a column vector $\mathbf{x} \in \mathbb{R}^d$, the corresponding row vector is denoted by \mathbf{x}^T , so that $\mathbf{x}^T = [x_1 \ \cdots \ x_d]$ and $[\mathbf{x}^T]^T = \mathbf{x}$. The scalar product between two vectors in \mathbb{R}^d is defined as $\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} = \sum_{i=1}^d x_i y_i$.
- The Euclidean norm of a vector $\mathbf{x} \in \mathbb{R}^d$ is defined by $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$.
- We use $\mathbb{R}^{n \times d}$ to denote the set of real rectangular matrices with n rows and d columns, where n and d will always be assumed to be at least 1. If $n = d$, $\mathbb{R}^{d \times d}$ refers to the set of square matrices of size d .
- We identify a matrix in $\mathbb{R}^{d \times 1}$ with its corresponding column vector in \mathbb{R}^d .
- Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, \mathbf{A}_{ij} refers to the coefficient from the i -th row and the j -th column of \mathbf{A} : the diagonal of \mathbf{A} is given by the coefficients \mathbf{A}_{ii} . Provided this notation is not ambiguous, we use the notations \mathbf{A} , $[\mathbf{A}_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$ and $[\mathbf{A}_{ij}]$ interchangeably.
- For every $d \geq 1$, \mathbf{I}_d refers to the identity matrix in $\mathbb{R}^{d \times d}$ (with 1s on the diagonal and 0s elsewhere).

Contents

1	Proximal methods and regularization	6
1.1	Regularized problems	6
1.2	Proximal methods	7
1.2.1	Proximal operator	7
1.2.2	Proximal point method	8
1.2.3	Proximal gradient methods	9
1.3	The case of ℓ_2 regularization	10
1.3.1	General principles	11
1.3.2	Linear least squares	11
1.4	Conclusion	12
2	Sparse optimization	13
2.1	Sparse regularization terms	13
2.1.1	Motivation : Sparse models and the ℓ_0 norm	13
2.1.2	Convex sparse regularizers	13
2.2	The case of ℓ_1 regularization	14
2.2.1	Proximal gradient and ℓ_1 regularization	14
2.2.2	Subgradient approach	15
2.2.3	The LASSO estimator	17
3	Coordinate descent methods	18
3.1	Basics of coordinate descent	18
3.1.1	Algorithm	18
3.1.2	Coordinate descent and stochastic gradient	19
3.2	Theoretical guarantees of coordinate descent methods	20
3.3	Applications of coordinate descent methods	21
3.4	Conclusion	22
4	Distributed and constrained optimization	23
4.1	Linear constraints and dual problem	23
4.2	Dual algorithms	24
4.2.1	Dual ascent	24
4.2.2	Augmented Lagrangian	24
4.3	Dual methods and decomposition	25
4.3.1	Dual decomposition	25
4.3.2	ADMM	25

4.4	Decentralized optimization	26
4.5	Conclusion	27
5	Exercises	28

Chapter 1

Proximal methods and regularization

1.1 Regularized problems

A common practice in machine learning problems consists in enforcing a specific structure of the machine learning model **through the objective function** rather than by using constraints on the model parameters. Such regularized problems have the following form :

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \quad \underbrace{f(\mathbf{x})}_{\text{data-fitting term}} \quad + \quad \underbrace{\lambda \Omega(\mathbf{x})}_{\text{regularization term}} . \quad (1.1.1)$$

where $\lambda > 0$ is called a regularization parameter. This parameter controls the weight of the regularization term compared to the data-fitting term. It is common for one of the terms in (1.1.1) (and sometimes both) to be nonsmooth, in the sense that the associated functions may not possess a gradient at every point. Such problems are called composite optimization problems in general, and regularized problems in specific settings such as data science and inverse problems.

When $\lambda \rightarrow 0$, the problem becomes essentially equivalent to minimizing the data-fitting term $f(\mathbf{x})$ only, and regularization no longer matters (numerically, this occurs for small values of λ). When $\lambda \rightarrow \infty$, the problem amounts to minimizing the regularization term, without accounting for the data-dependent term. A key issue in regularization consists in finding the right balance between data fitting and regularization.

Example 1.1.1 *The following functions are common choices for regularization terms:*

- ℓ_2 /ridge regularization: $\Omega(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$;
- ℓ_1 /LASSO regularization: $\Omega(\mathbf{x}) = \|\mathbf{x}\|_1$;
- Elastic net: $\Omega(\mathbf{x}) = \frac{\gamma}{2} \|\mathbf{x}\|_2^2 + \|\mathbf{x}\|_1$ for $\gamma > 0$;
- Group LASSO: Given a partition \mathcal{G} of $\{1, \dots, d\}$, set $\Omega(\mathbf{x}) = \sum_{g \in \mathcal{G}} \|\mathbf{x}_{[g]}\|_2$, where $\mathbf{x}_{[g]} \in \mathbb{R}^{|g|}$ is the vector formed by the subset of coordinates of \mathbf{x} corresponding to g .

At a broader level, regularization terms can be used to promote generalization capabilities of the solution, by reducing the sensitivity of the solution with respect to the problem data. They are also used to enforce sparsity in the problem solution. Finally, they can be used to improve the optimization landscape of a problem, for instance by guaranteeing that the regularized problem has a unique solution.

Constrained optimization and regularization The regularization approach bears a connection with the introduction of constraints on the problem variables. In fact, any set of constraints $x \in \mathcal{X}$ can be encoded as a regularization function using the so-called indicator function

$$\Omega(x) = \begin{cases} 0 & \text{if } x \in \mathcal{X} \\ +\infty & \text{otherwise.} \end{cases}$$

Conversely, it is possible to describe the solution of a regularized problem as that of a constrained optimization problem. However, the philosophy behind regularization is to penalize points that do not satisfy a certain property instead of excluding those points like in a constrained formulation. The interest of regularization lies in the possibility of weighing the regularization term relatively to the data-fitting term.

In this chapter, we review the main algorithmic approaches that can be used to tackle problems involving regularized optimization problems.

1.2 Proximal methods

Proximal methods are not restricted to nonsmooth optimization problems, but have proven particularly useful in the presence of nonsmoothness. Those techniques rely on the proximal operator, a mathematical object that is both defined as the solution to an optimization problem and used as a step in other optimization algorithms.

1.2.1 Proximal operator

The proximal operator is a fundamental tool used to both analyze problems with regularization and build algorithms for these problems.

Definition 1.2.1 Let $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ and $x \in \mathbb{R}^d$. The **proximal operator of φ at x** , denoted by $\text{prox}[\varphi][x]$, is given by

$$\text{prox}_\varphi[x] := \underset{u \in \mathbb{R}^d}{\text{argmin}} \left\{ \varphi(u) + \frac{1}{2} \|u - x\|^2 \right\}. \quad (1.2.1)$$

As stated, the proximal operator is a set-valued mapping, i.e. its value is a set of vectors (possibly empty or infinite!). However, for certain classes of functions, the value of the proximal operator is a singleton. It is then convenient to identify it with the corresponding vector. We make this identification below in the case of convex functions, which will be our focus throughout these notes.

Definition 1.2.2 Let $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function. Then the proximal operator of φ is defined as the function $\text{prox}_\varphi[\cdot] : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that

$$\forall x \in \mathbb{R}^d, \text{prox}_\varphi[x] := \mathbf{u}^*, \{\mathbf{u}^*\} = \underset{u \in \mathbb{R}^d}{\text{argmin}} \left\{ \varphi(u) + \frac{1}{2} \|u - x\|^2 \right\}. \quad (1.2.2)$$

Example 1.2.1 For any $x \in \mathbb{R}^d$ and any $\lambda > 0$, the following properties hold:

- $\text{prox}_0[x] = x$.

- $\text{prox}_{\frac{\lambda}{2}\|\cdot\|^2}[\mathbf{x}] = \frac{\mathbf{x}}{1+\lambda}$.
- $\text{prox}_{\lambda\|\cdot\|_1}[\mathbf{x}]$ is defined componentwise by

$$\forall j = 1, \dots, d, \quad \left[\text{prox}_{\lambda\|\cdot\|_1}[\mathbf{x}] \right]_j = \begin{cases} x_j - \lambda & \text{if } x_j > \lambda \\ 0 & \text{if } x_j \in [-\lambda, \lambda] \\ x_j + \lambda & \text{if } x_j < -\lambda \end{cases}$$

- If h is the indicator function of a convex set \mathcal{X} (i.e. $h(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathcal{X}$ and $h(\mathbf{x}) = \infty$ otherwise), then $\text{prox}_h[\mathbf{x}]$ corresponds to the projection on \mathbf{x} onto \mathcal{X} .

In general, proximal operators are useful when they are uniquely defined (which is always the case when the original function is convex) and when they are easy to compute numerically. The next two sections describe two algorithms built on proximal operator calculations.

1.2.2 Proximal point method

In this section, we consider

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} h(\mathbf{x}) \tag{1.2.3}$$

where $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function. Without assumptions on the differentiability of h , we cannot rely on a gradient-type method for optimizing this problem. However, a clever algorithm based on proximal operators was proposed in the 1970s, that builds a sequence of approximate solutions of problem (1.2.3).

Algorithm 1: Proximal point method.

Initialization: $\mathbf{x}_0 \in \mathbb{R}^d$.

for $k = 0, 1, \dots$ **do**

1. Compute a steplength $\alpha_k > 0$.
2. Compute \mathbf{x}_{k+1} such that

$$\mathbf{x}_{k+1} \in \underset{\mathbf{z} \in \mathbb{R}^d}{\text{argmin}} \left\{ h(\mathbf{z}) + \frac{1}{2\alpha_k} \|\mathbf{z} - \mathbf{x}_k\|_2^2 \right\}. \tag{1.2.4}$$

end

The process is given in Algorithm 1. At every iteration, the next iterate \mathbf{x}_{k+1} is computed by solving

$$\underset{\mathbf{z} \in \mathbb{R}^d}{\text{minimize}} \left\{ h(\mathbf{z}) + \frac{1}{2\alpha_k} \|\mathbf{z} - \mathbf{x}_k\|_2^2 \right\}. \tag{1.2.5}$$

At first glance, this problem is not easier than the original one, because it involves the original objective h . However, the presence of the **proximal term** $\frac{1}{2\alpha_k} \|\mathbf{z} - \mathbf{x}_k\|_2^2$ actually makes the problem easier to solve. In particular, the subproblem (1.2.5) is strongly convex (and thus the update (1.2.4) is uniquely defined), and solving this subproblem always improves the objective.

Lemma 1.2.1 *At the k th iteration of Algorithm 1, we have*

$$h(\mathbf{x}_{k+1}) \leq h(\mathbf{x}_k) - \frac{1}{2\alpha_k} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2. \quad (1.2.6)$$

Using the definition of the proximal operator, note that an iteration of Algorithm 1 can be rewritten as

$$\mathbf{x}_{k+1} = \text{prox}_{\alpha_k h}[\mathbf{x}_k].$$

1.2.3 Proximal gradient methods

In this section, we consider a specific class of regularized optimization problems, often termed as **composite optimization** problems.

Definition 1.2.3 (Composite optimization) *A composite optimization problem is of the form:*

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + \lambda \Omega(\mathbf{x}), \quad (1.2.7)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a smooth, \mathcal{C}^1 function, $\lambda > 0$ and $\Omega : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex, nonsmooth regularizer.

Because of the nonsmoothness of Ω , the overall objective of problem (1.2.7) is nonsmooth, and one may then consider applying a subgradient method to tackle this problem. However, knowledge about the particular problem structure, and the fact that f is a smooth function, allows for deriving better methods for this problem, that do not require direct calculations of subgradients.

Proximal gradient techniques, that are described in Algorithm 2, exploit the smoothness of f to iteratively solve a sequence of subproblems approximating the original one.

Algorithm 2: Proximal gradient method.

Initialization: $\mathbf{x}_0 \in \mathbb{R}^d$.

for $k = 0, 1, \dots$ **do**

1. Compute the gradient of the smooth part $\nabla f(\mathbf{x}_k)$.
2. Compute a steplength $\alpha_k > 0$.
3. Compute \mathbf{x}_{k+1} such that

$$\mathbf{x}_{k+1} \in \underset{\mathbf{x} \in \mathbb{R}^d}{\text{argmin}} \left\{ f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2 + \lambda \Omega(\mathbf{x}) \right\}. \quad (1.2.8)$$

end

The cost of an iteration of Algorithm 2 includes a gradient calculation as well as solving an auxiliary optimization problem (1.2.8), called the **proximal subproblem**. Such a method is only practical if the subproblems are easier to solve than the original problem.

Connection with the proximal operator If $\Omega \equiv 0$ (i. e. Ω is the zero function and the problem is un-regularized), one can show that the solution of (1.2.8) is given by

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

We thus recognize the gradient descent iteration. For arbitrary Ω , the following formula applies:

$$\mathbf{x}_{k+1} = \text{prox}_{\alpha_k \lambda \Omega} [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)].$$

As such, proximal gradient bears a close connection with the proximal operator.

However, note that applying the proximal point to problem (1.2.7) yields a different iteration, even in the case $\Omega \equiv 0$! Indeed, in that case, and assuming convexity of f , the proximal point iteration can be written as

$$\mathbf{x}_{k+1} = \text{prox}_{\alpha_k f} [\mathbf{x}_k] = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_{k+1}).$$

The proximal point method is thus an implicit method, in the sense that its iteration is defined implicitly by the above equation to be solved for \mathbf{x}_{k+1} .

Remark 1.2.1 *Proximal gradient methods can be designed using most of the tools that can be applied to gradient descent : this includes stepsize choices, acceleration as well as stochastic aspects. Moreover, theoretical guarantees have been established for both nonconvex and convex f , though the latter has attracted more attention in the optimization literature. Note that the analysis of such methods is performed by looking at the vector*

$$\frac{1}{\alpha_k} (\mathbf{x}_k - \text{prox}_{\alpha_k \lambda \Omega} [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]),$$

that plays a similar role than the gradient in analyzing gradient descent (and reduces to the gradient when $\Omega \equiv 0$).

1.3 The case of ℓ_2 regularization

A **regularized problem with ℓ_2 regularization**, also known as **ridge regularized problem**, has the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + \frac{\lambda}{2} \|\mathbf{x}\|^2. \quad (1.3.1)$$

The ridge regularizer $\mathbf{x} \mapsto \frac{1}{2} \|\mathbf{x}\|^2$ has several interpretations. It effectively penalizes \mathbf{x} s with large components, and can be shown to be equivalent to a constraint on the squared norm $\|\mathbf{w}\|^2$. In addition, a ridge regularizer has the effect to reduce the variance of the problem solution with respect to the data. Finally, when the regularizer $\lambda > 0$ is big enough, this often turns the objective function into a strongly convex one, with the positive implications in terms of convergence speed and uniqueness of the (global) minimum.

1.3.1 General principles

Suppose that the function f in (1.3.1) is \mathcal{C}^1 . Then, the k th iteration of gradient descent applied to (1.3.1) is given by

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k (\nabla f(\mathbf{x}_k) + \lambda \mathbf{x}_k) = (1 - \lambda \alpha_k) \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \quad (1.3.2)$$

where $\alpha_k > 0$ is a positive stepsize. When $\lambda = 0$, we recover the classical gradient descent iteration. When $\lambda > 0$, however, the algorithm will modify the coefficients of the vector \mathbf{x}_k . In this case, one typically chooses α_k so that $1 - \lambda \alpha_k \in (0, 1)$, and thus the coefficients of \mathbf{x}_k are reduced in magnitude at every iteration in the update formula. This process corresponds to **weight decay** in deep learning.

Suppose now that we apply proximal gradient to problem (1.3.1). The corresponding iteration is

$$\mathbf{x}_{k+1} = \text{prox}_{\alpha_k \lambda \Omega} [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)] = \frac{1}{1 + \lambda \alpha_k} [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)] = \frac{1}{1 + \lambda \alpha_k} \mathbf{x}_k - \frac{\alpha_k}{1 + \lambda \alpha_k} \nabla f(\mathbf{x}_k). \quad (1.3.3)$$

As λ increases, this iteration decays both the coordinates of the iterate \mathbf{x}_k and that of the gradient step. In fact, it can be shown that the coordinates of the iterates decrease uniformly in a smooth fashion as $\lambda \rightarrow 0$.

1.3.2 Linear least squares

To illustrate further the use of proximal methods, we consider a simple linear least-squares problem. Given data under the form of a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ and a vector $\mathbf{y} \in \mathbb{R}^n$, we seek to solve

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2n} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{x}\|^2 \quad (1.3.4)$$

for $\lambda \geq 0$.

Gradient descent VS proximal point The iteration of gradient descent applied to this problem gives

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \left[\frac{1}{n} \mathbf{A}^T (\mathbf{A}\mathbf{x}_k - \mathbf{y}) + \lambda \mathbf{x}_k \right], \quad (1.3.5)$$

where $\alpha_k > 0$ and $\frac{1}{n} \mathbf{A}^T (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ is the gradient of the data fitting term at \mathbf{x}_k . On the other hand, the iteration of the proximal point method gives

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \left[\frac{1}{n} \mathbf{A}^T (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}) + \lambda \mathbf{x}_{k+1} \right],$$

which can be rewritten as

$$\mathbf{x}_{k+1} = \left[(1 + \lambda \alpha_k) \mathbf{I} + \frac{1}{n} \mathbf{A}^T \mathbf{A} \right]^{-1} \left(\mathbf{x}_k + \alpha_k \frac{1}{n} \mathbf{A}^T \mathbf{y} \right). \quad (1.3.6)$$

One then sees that each iteration of the proximal point method requires to solve a linear system, which is more expensive than the cost of one gradient descent iteration. However, it can be shown that (1.3.6) has good stability properties. Note also that the linear system solve becomes easier as λ increases.

1.4 Conclusion

Proximal methods are a class of algorithms that rely on proximal operators to perform optimization steps on possibly complex functions. They are particularly useful in the context of regularized optimization, in which

The most classical regularization techniques in optimization and beyond are ℓ_2 and ℓ_1 regularization. In a data-related setting, the former is typically used to reduce the dependency of the solution with respect to the data, while the latter promotes solutions that have zero coordinates.

In both cases, the proximal gradient method can then be written in closed form, leading to algorithms that are implementable quite efficiently. In particular, the case of regularized linear least squares has been widely studied, and the resulting problems are now part of numerous statistics and signal processing toolboxes.

Chapter 2

Sparse optimization

In this chapter, we investigate regularized problems in which the goal is to produce a solution vector with a significant number of zero coefficients that still provide a good model with respect to the data-fitting term in the regularized problem. Such vectors, called **sparse**, are often thought as a way to simplify the original model. For instance, in the case of linear models, zero coefficients means that some features are ignored by the model, and thus learning can be performed using less features. This feature selection process can be done *a posteriori* (i.e. once the learning problem has been solved), but it can also be encoded into the optimization problem through the use of regularization terms.

2.1 Sparse regularization terms

2.1.1 Motivation : Sparse models and the ℓ_0 norm

While computing a model to explain some data, we might want to compute a model that explains the data using as few features as possible¹. Mathematically speaking, if our model is parameterized by a vector $\mathbf{x} \in \mathbb{R}^d$, our goal is to compute a vector that explains the data with as few nonzero coordinates as possible.

There exists a regularizer that penalized vectors with nonzero components (not just large as opposed to the ridge regularizer), called the ℓ_0 norm². An ℓ_0 -regularized problem has the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_0, \quad \|\mathbf{v}\|_0 = |\{i | [\mathbf{v}]_i \neq 0\}|.$$

However, this function is nonsmooth and discontinuous; its combinatorial nature also introduces more complexity to the original problem.

2.1.2 Convex sparse regularizers

Researchers have investigated numerous alternatives to the ℓ_0 . The main surrogate for the ℓ_0 norm is the ℓ_1 norm, defined by

$$\|\mathbf{x}\|_1 = \sum_{i=1}^d |[\mathbf{x}]_i|. \tag{2.1.1}$$

¹The goal of this process is *feature selection*.

²Though technically this function defines a semi-norm.

It can be shown that this function is the closest convex upper bound to the ℓ_0 norm. This function is continuous and convex; moreover, it is actually a norm function, which endows it with many desirable properties. As a result, the class of problems of the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1, \quad (2.1.2)$$

called ℓ_1 **regularized problems** or **lasso regularized problems**, has emerged as a tractable alternative to the ℓ_0 regularized formulation.

Variations on ℓ_1 A number of regularizing terms have been proposed based on the ℓ_1 norm, as well as the ℓ_2 norm

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d [\mathbf{x}]_i^2}.$$

The later is also a convex approximation to the ℓ_0 norm, but it is not as tight as the ℓ_1 norm. Nevertheless, one can combine the ℓ_1 and ℓ_2 norms to build special regularizers, such as the so-called group LASSO regularizer

$$\Omega(\mathbf{x}) = \sum_{g \in \mathcal{G}} \|\mathbf{x}_g\|_2,$$

where \mathcal{G} is a partition of $\{1, \dots, d\}$ and \mathbf{x}_g denotes the vector in $\mathbb{R}^{|g|}$ formed by the coordinates of \mathbf{x} corresponding to g . Other variations on this concept are presented in Figure 2.1 via their level sets.

2.2 The case of ℓ_1 regularization

2.2.1 Proximal gradient and ℓ_1 regularization

A natural way to tackle problems of the form 2.1.2 is through the proximal gradient framework of Algorithm 2. Indeed, unlike for general regularizers, one can obtain a closed-form solution of the subproblem (1.2.8). Indeed, the proximal subproblem, given by

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \left\{ f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\},$$

has a unique solution. To obtain it, one computes the usual gradient step $\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$, then one applies the **soft-thresholding function** $s_{\alpha_k \lambda}(\bullet)$ to each component, where this function is given by

$$\forall \mu > 0, \forall t \in \mathbb{R}, \quad s_\mu(t) = \begin{cases} t + \mu & \text{if } t < -\mu \\ t - \mu & \text{if } t > \mu \\ 0 & \text{otherwise.} \end{cases}$$

(Note that this function encodes the proximal operator for the ℓ_1 norm.) As a result, the solution of the proximal subproblem is defined component-wise according to the components of the gradient step. The resulting update is at the heart of the corresponding proximal algorithm, popularized in signal and image processing under the name ISTA (Iterative Soft-Thresholding Algorithm). A description of ISTA is given in Algorithm 3.

It can be shown that the use of the soft-thresholding function does promote zero components in the new iterates, which results in sparser solutions at the end of the algorithmic run.

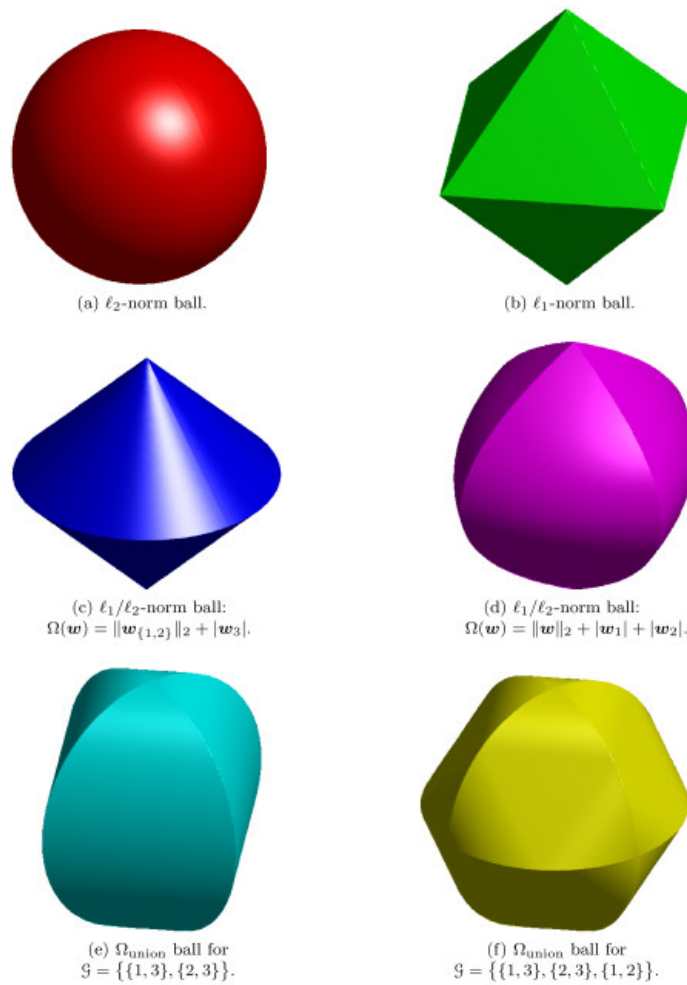


Figure 2.1: Level sets (balls of fixed radius) for sparse regularizers. *Source:[1]*.

Remark 2.2.1 A notable improvement on ISTA was the inclusion of momentum, which resulted in a new algorithm called FISTA (Fast ISTA): this method is now the most widely used instance of ISTA.

2.2.2 Subgradient approach

In smooth optimization, the gradient is a vector that quantifies the rate of local change in the function. Generalized notions of gradient have been defined, especially in the case of convex, nonsmooth functions. We provide below the most common definition of a subgradient.

Definition 2.2.1 (Subgradient and subdifferential) Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function. A vector $\mathbf{g} \in \mathbb{R}^d$ is called a *subgradient* of f at $\mathbf{x} \in \mathbb{R}^d$ if

$$\forall \mathbf{z} \in \mathbb{R}^n, \quad f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}).$$

The set of all subgradients of f at \mathbf{x} is called the *subdifferential* of f at \mathbf{x} , and denoted by $\partial f(\mathbf{x})$.

Algorithm 3: ISTA: Iterative Soft-Thresholding Algorithm.**Initialization:** $\mathbf{x}_0 \in \mathbb{R}^d$.**for** $k = 0, 1, \dots$ **do**

1. Compute the gradient of the smooth part $\nabla f(\mathbf{x}_k)$.
2. Compute a steplength $\alpha_k > 0$.
3. Compute \mathbf{x}_{k+1} component-wise through the following rule

$$[\mathbf{x}_{k+1}]_i = \begin{cases} [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]_i + \alpha_k \lambda & \text{if } [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]_i < -\alpha_k \lambda \\ [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]_i - \alpha_k \lambda & \text{if } [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]_i > \alpha_k \lambda \\ 0 & \text{if } [\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]_i \in [-\alpha_k \lambda, \alpha_k \lambda]. \end{cases} \quad (2.2.1)$$

end

Note that when the function f is differentiable at \mathbf{x} , we have $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$, thus the notion of subdifferential matches that of the gradient for differentiable functions. In addition, the subdifferential provides a characterization of the global minima of f , as shown by the following result.

Theorem 2.2.1 Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, and $\bar{\mathbf{x}} \in \mathbb{R}^d$.

$$\mathbf{0} \in \partial f(\bar{\mathbf{x}}) \iff \bar{\mathbf{x}} \in \underset{\mathbf{x} \in \mathbb{R}^d}{\operatorname{argmin}} f(\mathbf{x}).$$

Again, the result of Theorem 2.2.1 generalizes that of the smooth setting, since in that case the subdifferential consists in a single element.

Example 2.2.1 Let $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = |x|$.

$$\partial f(x) = \begin{cases} -1 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ [-1, 1] & \text{if } x = 0. \end{cases}$$

The set $[-1, 1]$ contains 0, which confirms that $x^* = 0$ is the minimum of f .

Note that the subdifferential of the ℓ_1 norm is defined componentwise according to the example above.

By analogy with gradient descent, we can design a subgradient method, as shown by Algorithm 4.

Such a method offers a flexibility in choosing the subgradient, which can be an issue. Moreover, choosing the stepsize is more difficult than for gradient descent, due to the nonsmooth nature of the problem. In fact, a subgradient can lead to increase in the function value for any stepsize, hence the choice of subgradient is critical to the success of this method. A convenient choice consists in setting \mathbf{g}_k such that

$$\mathbf{g}_k \in \underset{\mathbf{g} \in \mathbb{R}^d}{\operatorname{argmin}} \{ \|\mathbf{g}\| \mid \mathbf{g} \in \partial f(\mathbf{x}_k) \}. \quad (2.2.2)$$

Algorithm 4: Generic subgradient method.

Initialization: $\mathbf{x}_0 \in \mathbb{R}^d$.
for $k = 0, 1, \dots$ **do**
 1. Compute a subgradient $\mathbf{g}_k \in \partial f(\mathbf{x}_k)$.
 2. Compute a steplength $\alpha_k > 0$.
 3. Set $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k$.
end

With this choice, Algorithm 4 will not change the iterate once it reaches a point with $\mathbf{0}_{\mathbb{R}^d}$ in its subdifferential (i.e. a minimum for a convex function). Computing \mathbf{g}_k according to (2.2.2) may however be an expensive procedure as it could require to compute the entire subdifferential.

In the context of ℓ_1 norm regularization, it is possible to define the subdifferential of $f + \lambda \|\cdot\|_1$ explicitly when f is \mathcal{C}^1 . Otherwise, calculus rules about subdifferentials come into play.

2.2.3 The LASSO estimator

As for ℓ_2 regularization, using ℓ_1 regularization together with a linear least-squares loss is an important class of problems. Such problems have the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \lambda \|\mathbf{x}\|_1, \quad (2.2.3)$$

where $\mathbf{A} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$ form the problem data, and $\lambda > 0$. This problem is often referred to as the LASSO (Least Absolute Shrinkage and Selection Operator) problem, though other terminologies such as basis pursuit can be found in the literature.

The properties of the solution of problem (2.2.3) can be described according to the columns of \mathbf{A} . Every column contains all the observed values for a specific feature captured by the linear model.

Theorem 2.2.2 *Let $\mathbf{b}_1, \dots, \mathbf{b}_d$ denote the columns of \mathbf{A} , and let \mathbf{x}^* be a solution of problem (2.2.3). Then, for any $j = 1, \dots, d$, we have*

$$\begin{cases} |\mathbf{b}_j^\top (\mathbf{A}\mathbf{x}^* - \mathbf{y})| \leq \lambda n & \text{if } [\mathbf{x}^*]_j = 0 \\ \mathbf{b}_j^\top (\mathbf{A}\mathbf{x}^* - \mathbf{y}) = -\lambda n \operatorname{sgn}([\mathbf{x}^*]_j) & \text{otherwise,} \end{cases} \quad (2.2.4)$$

where $\operatorname{sgn}(t) = 1$ if $t > 0$ and $\operatorname{sgn}(t) = -1$ if $t < 0$.

The result of Theorem 2.2.2 implies that the larger λ is, the more components of \mathbf{x}^* will be zero. In that sense, the use of ℓ_1 regularization leads to sparser solutions than in absence of regularization.

Chapter 3

Coordinate descent methods

In this chapter, we address the treatment of large-scale optimization problems, where the number of parameters to be optimized over is extremely large. As we witness a growth in both the model complexity (i.e. the number of parameters) and the amount of data available (i.e. the size of the dataset), standard optimization techniques may suffer from the curse of dimensionality and their performance may deteriorate as dimensions grow. As a result, the practical difficulty of the problem increases with the dimension, simply because there are more variables to consider. However, on structured problems such as those arising in data science, there often exists a low-dimensional or separable structure that allows for optimization steps to be taken over a subset of variables. This is the underlying idea of **coordinate descent methods**, that have regained interest in the early 2000s due to their applicability in certain data science settings.

3.1 Basics of coordinate descent

3.1.1 Algorithm

Consider the unconstrained optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}), \quad (3.1.1)$$

where $f \in \mathcal{C}^1(\mathbb{R}^d)$. The idea of coordinate descent methods consist in taking a gradient step with respect to a single decision variable at every iteration. To this end, we observe that for every $\mathbf{x} \in \mathbb{R}^d$, the gradient of f at \mathbf{x} can be decomposed as

$$\nabla f(\mathbf{x}) = \sum_{j=1}^d \nabla_j f(\mathbf{x}) \mathbf{e}_j,$$

where ∇_j denotes the partial derivative with respect to the j -th variable of the function f (that is, the j th coordinate of ∇f) and $\mathbf{e}_j \in \mathbb{R}^d$ is the j th coordinate vector of the canonical basis in \mathbb{R}^d . The coordinate descent approach replaces the full gradient by a step along a coordinate gradient, as formalized in Algorithm 5.

The variants of coordinate descent are mainly identified by the way they select the coordinate sequence $\{j_k\}$. There exist numerous rules for choosing the coordinate index, among which:

Algorithm 5: Coordinate descent method.**Initialization:** $\mathbf{x}_0 \in \mathbb{R}^d$.**for** $k = 0, 1, \dots$ **do**1. Select a coordinate index $j_k \in \{1, \dots, d\}$.2. Compute a steplength $\alpha_k > 0$.

3. Set

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla_{j_k} f(\mathbf{x}_k) \mathbf{e}_{j_k}. \quad (3.1.2)$$

end

- **Cyclic:** Select the indices by cycling over $\{1, \dots, d\}$ in that order. After d iterations, all indices have been selected.
- **Randomized cyclic:** Cycle through a random ordering of $\{1, \dots, d\}$, that changes every d steps.
- **Randomized:** Draw j_k at random in $\{1, \dots, d\}$ at every iteration.

The last two strategies are those for which the strongest results can be obtained.

Block coordinate descent Rather than using a single index, it is possible to select a subset of the coordinates (called “block” in the literature). The k th iteration of such a *block coordinate descent* algorithm thus is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \sum_{j \in \mathcal{B}_k} \nabla_j f(\mathbf{x}_k) \mathbf{e}_j, \quad (3.1.3)$$

where $\mathcal{B}_k \subset \{1, \dots, d\}$.

3.1.2 Coordinate descent and stochastic gradient

Our description of coordinate descent, and particularly the randomized variant, is reminiscent of the stochastic gradient algorithm. In fact, randomized coordinate descent can be viewed as a special case of stochastic gradient, in which the formula

$$\nabla f(\mathbf{x}) = \frac{1}{d} \sum_{j=1}^d \nabla_j f(\mathbf{x}) \mathbf{e}_j$$

is used to define a finite sum with d gradients that can be sampled using any distribution that one would use in a stochastic gradient setting. Note that, unlike in a stochastic gradient framework, any coordinate descent step (even randomized ones) uses a descent direction.

Consider a finite-sum problem of the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad f_i(\mathbf{x}) := \ell_i(\mathbf{a}_i^\top \mathbf{x}), \quad (3.1.4)$$

where $\mathbf{a}_i^T \mathbf{x}$ is a linear model of the data vector \mathbf{a}_i , and $\ell_i : \mathbb{R} \rightarrow \mathbb{R}$ is a convex loss function specific to the i th data point (such as $\ell_i(h) = \frac{1}{2}(h - y_i)^2$ for linear least squares). If the number of data points is large, it is natural to think of applying stochastic gradient to this problem. Another approach consists in considering an equivalent formulation of (3.1.4) through (Fenchel) duality, given by

$$\underset{\mathbf{v} \in \mathbb{R}^n}{\text{maximize}} g(\mathbf{v}) := -\frac{1}{n} \sum_{i=1}^n f_i^*(v_i) \quad (3.1.5)$$

where for any convex function $\phi : \mathbb{R}^m \rightarrow \mathbb{R}$, the convex conjugate function ϕ^* is defined by

$$\phi^*(\mathbf{a}) = \sup_{\mathbf{b} \in \mathbb{R}^m} \{\mathbf{a}^T \mathbf{b} - \phi(\mathbf{b})\}.$$

The so-called dual problem (3.1.5) has a finite-sum, separable form. It can thus be tackled using (dual) *coordinate ascent*, the counterpart of coordinate descent for minimization: the iteration of this method is given by

$$\mathbf{v}_{k+1} = \mathbf{v}_k + \alpha_k \nabla_i g(\mathbf{v}_k), \quad (3.1.6)$$

leading to updating the iterate one coordinate at a time. Under appropriate assumptions on the problem, the iteration (3.1.6) is equivalent to the original stochastic gradient iteration, with $\mathbf{x}_k = \frac{1}{\lambda n} \sum_{i=1}^n [\mathbf{v}_k]_i \mathbf{a}_i$. For this reason, stochastic gradient is sometimes viewed as applying coordinate ascent to the dual problem.

3.2 Theoretical guarantees of coordinate descent methods

A famous 3-dimensional example designed by M. J. D. Powell in 1973 shows that coordinate descent methods do not necessarily converge.

Nevertheless, it is possible to provide guarantees on coordinate descent methods under appropriate assumptions. In particular, a linear rate of convergence can be obtained for coordinate descent methods on strongly convex problems: we provide below the necessary assumptions to arrive at such a result.

Assumption 3.2.1 *The objective function f in (3.1.1) is C^1 and μ -strongly convex, with $f^* = \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$. Moreover, for every $j = 1, \dots, d$, the partial derivative $\nabla_j f$ is L_j -Lipschitz continuous, i.e.*

$$\forall \mathbf{x} \in \mathbb{R}^d, \forall h \in \mathbb{R}, \quad |\nabla_j f(\mathbf{x} + h\mathbf{e}_j) - \nabla_j f(\mathbf{x})| \leq L_j |h|. \quad (3.2.1)$$

We let $L_{\max} = \max_{1 \leq j \leq d} L_j$.

Theorem 3.2.1 *Suppose that Assumption 3.2.1 holds, and that Algorithm 5 is applied to problem (3.1.1) with $\alpha_k = \frac{1}{L_{\max}}$ for all k and j_k being drawn uniformly at random in $\{1, \dots, d\}$. Then, for any $K \in \mathbb{N}$, we have*

$$\mathbb{E}[f(\mathbf{x}_k) - f^*] \leq \left(1 - \frac{\mu}{dL_{\max}}\right)^K (f(\mathbf{x}_0) - f^*). \quad (3.2.2)$$

Other results have been established in the convex and nonconvex settings, under additional assumptions. In all cases, properties on the partial derivatives are required.

Acceleration In a convex optimization setting, it is possible to combine randomized coordinate descent with the accelerated gradient paradigm¹. Starting from $\mathbf{x}_0 \in \mathbb{R}^d$ et $\mathbf{v}_0 = \mathbf{x}_0$, every iteration k draws i_k uniformly at random between 1 and d , then performs the following calculations:

$$\begin{cases} \mathbf{u}_k & := \lambda_k \mathbf{v}_k + (1 - \lambda_k) \mathbf{x}_k \\ \mathbf{x}_{k+1} & := \mathbf{u}_k - \frac{1}{L_{j_k}} \nabla_{j_k} f(\mathbf{u}_k) \mathbf{e}_{j_k} \\ \mathbf{v}_{k+1} & := \mu_k \mathbf{v}_k + (1 - \mu_k) \mathbf{u}_k - \frac{\gamma_k}{L_{j_k}} \nabla_{j_k} f(\mathbf{u}_k) \mathbf{e}_{j_k}, \end{cases}$$

where $\{\lambda_k, \mu_k, \gamma_k\}$ are sequences that depend on the dimension d , and possibly on a strong convexity constant if f happens to be strongly convex. Although this method possesses better complexity guarantees than randomized coordinate descent, it requires to maintain additional vector sequences. In terms of accesses to the coefficients, the cost of this accelerated method is thus higher than that of a basic coordinate descent iteration.

3.3 Applications of coordinate descent methods

Coordinate descent techniques are particularly useful for large-scale sparse optimization. Consider a regularized problem of the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \tilde{f}_i(\mathbf{a}_i^\top \mathbf{x}) + \sum_{j=1}^d \Omega(w_j), \quad (3.3.1)$$

where $\tilde{f}_i : \mathbb{R} \rightarrow \mathbb{R}$ is (possibly) data-dependent, $\mathbf{x}_i \in \mathbb{R}^d$ is a **sparse** data vector, and $\Omega : \mathbb{R} \rightarrow \mathbb{R}$ is a regularization function applied componentwise to the vector \mathbf{x} .

Example 3.3.1 (Regularized least squares with sparse data) Given $\mathbf{X} \in \mathbb{R}^{n \times d}$ with sparse rows and $\mathbf{y} \in \mathbb{R}^n$, consider the problem

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) := \frac{1}{2n} \sum_{i=1}^n \|\mathbf{X} \mathbf{x} - \mathbf{y}\|^2 + \lambda \sum_{j=1}^d w_j^2.$$

Apply Algorithm 5 to this problem. For any iteration k , if we move along the j_k th coordinate, the partial derivative under consideration is

$$\nabla_{j_k} f(\mathbf{x}_k) = \mathbf{x}_{j_k}^\top (\mathbf{X} \mathbf{x}_k - \mathbf{y}) + 2\lambda [\mathbf{x}_k]_{j_k}.$$

By storing the vector $\{\mathbf{X} \mathbf{x}_k\}$ across all iterations, the calculation of $\nabla_{j_k} f(\mathbf{x}_k)$ can be greatly reduced when \mathbf{x}_{j_k} is sparse, to the point that the cost of a coordinate descent iteration will be of the order of the number of nonzero elements in \mathbf{x}_{j_k} .

Coordinate descent techniques are quite prominent in parallel optimization algorithms. In this setting, several cores are cooperating to solve problem (3.1.1): each core can then run *its own coordinate descent method* and all cores update the same shared iterate vector. The most efficient parallel coordinate descent techniques perform these iterations in an asynchronous fashion, which does not prevent from guaranteeing convergence of this framework!

¹See Irène Waldspurger's lecture on this topic.

[Link with stochastic gradient](#)

3.4 Conclusion

Large-scale problems have always pushed optimization algorithms to their limits, and have led to reconsidering certain algorithms in light of their applicability to large-scale settings. Coordinate descent methods are the perfect example of classical techniques that regained popularity because of their efficiency in data science settings. On some instances, randomized coordinate descent techniques bear a close connection with stochastic gradient methods. More globally, coordinate descent methods are quite efficient on large-dimensional problems that have a separable structure. Finally, the use of coordinate descent methods in parallel environments has also contributed to their revival in optimization.

Chapter 4

Distributed and constrained optimization

In this chapter, we describe the theoretical insights behind distributed optimization formulations, in which several agents collaborate to solve an optimization problem. This paradigm can be modeled using a linearly constrained optimization formulation, and handling such formulations requires dedicated algorithms. We first set the mathematical foundations of these methods via a brief introduction to duality, then present our algorithms of interest.

4.1 Linear constraints and dual problem

Consider the following optimization problem with linear equality constraints:

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{b}, \quad (4.1.1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$. For simplicity, we will assume that the feasible set $\{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{Ax} = \mathbf{b}\}$ is not empty.

Duality theory consists in handling constraints formulations by reformulating the problem into an unconstrained optimization problem. We present the theoretical arguments for the special case of problem (4.1.1), which yields a much simpler analysis.

Definition 4.1.1 *The Lagrangian function of problem (4.1.1) is given by*

$$\mathcal{L}(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) + \mathbf{y}^T (\mathbf{Ax} - \mathbf{b}). \quad (4.1.2)$$

The Lagrangian function combines the objective function and the constraints, and allows to restate the original problem as an unconstrained one, called the primal problem:

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \max_{\mathbf{y} \in \mathbb{R}^m} \mathcal{L}(\mathbf{x}, \mathbf{y}). \quad (4.1.3)$$

The solutions of the primal problem are identical to that of problem (4.1.3) in our case. The difficulty of solving problem (4.1.3) lies in the definition of its objective function as the optimal value of a maximization problem.

Definition 4.1.2 The **dual problem** of (4.1.1) is the maximization problem

$$\text{maximize}_{\mathbf{y} \in \mathbb{R}^m} \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{y}), \quad (4.1.4)$$

where the function $\mathbf{y} \mapsto \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{y})$ is called the dual function of the problem.

Unlike the primal problem, the dual problem is always concave (i.e. the opposite of the dual function is convex), which facilitates its resolution by standard optimization techniques. The goal is then to solve the dual problem in order to get the solution of the primal problem, thanks to properties such as the one below.

Assumption 4.1.1 We suppose that **strong duality** holds between problem (4.1.1) and its dual, that is,

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^m} \mathcal{L}(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{y} \in \mathbb{R}^m} \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{y}).$$

Assumption 4.1.1 is typically satisfied when f is convex, but strong duality may hold even on non-convex problems.

4.2 Dual algorithms

We are now concerned with solving the dual problem (4.1.4), and we will present three methods for this purpose.

4.2.1 Dual ascent

The **dual ascent** method is implicitly a subgradient method applied to the dual problem (which we recall is a maximization problem). At every iteration, it starts from a primal-dual pair $(\mathbf{x}_k, \mathbf{y}_k)$ and performs the following iteration:

$$\begin{cases} \mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{y}_k) \\ \mathbf{y}_{k+1} = \mathbf{y}_k + \alpha_k (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}), \end{cases} \quad (4.2.1)$$

where $\alpha_k > 0$ is a stepsize for the dual ascent step, and $\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}$ is a subgradient for the dual function $\mathbf{y} \mapsto \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{y})$ at \mathbf{y}_k .

4.2.2 Augmented Lagrangian

The dual ascent method generally has weak convergence guarantees. For this reason, the optimization literature has introduced other frameworks based on a regularized version of the Lagrangian function.

Definition 4.2.1 The **augmented Lagrangian** of problem (4.1.1) is the function on $\mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}_{++}$ by

$$\mathcal{L}^a(\mathbf{x}, \mathbf{y}; \boldsymbol{\lambda}) := f(\mathbf{x}) + \mathbf{y}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) + \frac{\boldsymbol{\lambda}}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2. \quad (4.2.2)$$

Augmented Lagrangians thus are a family of functions parameterized by $\lambda > 0$, that put more emphasis on the constraint violation as λ grows.

The **augmented Lagrangian** algorithm, also called method of multipliers, performs the following iteration:

$$\begin{cases} \mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}^a(\mathbf{x}, \mathbf{y}_k; \lambda) \\ \mathbf{y}_{k+1} = \mathbf{y}_k + \lambda(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}). \end{cases} \quad (4.2.3)$$

In this algorithm, λ is constant and used as a constant stepsize: many more sophisticated choices of both the augmented Lagrangian function and the stepsizes have been proposed. In general, the advantages of augmented Lagrangian techniques are that the subproblems defining \mathbf{x}_{k+1} become easier to solve (thanks to regularization) and that the overall guarantees on the primal-dual pair are stronger.

4.3 Dual methods and decomposition

A key idea in modern optimization, that has resulted in numerous theoretical and numerical improvements, consists in exploiting the structure of a given problem as much as possible. The underlying idea is that a decomposition of a large, complex problem can lead to many smaller and simpler (sub)problems that will be easier and cheaper to solve than the original one. We describe below how this idea can be carried out in the context of dual algorithms.

4.3.1 Dual decomposition

Suppose that we consider a linearly constrained problem with a separable form:

$$\begin{cases} \operatorname{minimize}_{\mathbf{u} \in \mathbb{R}^{d_1}, \mathbf{v} \in \mathbb{R}^{d_2}} & f(\mathbf{u}) + g(\mathbf{v}) \\ \text{subject to} & \mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} = \mathbf{c}, \end{cases} \quad (4.3.1)$$

where $\mathbf{A} \in \mathbb{R}^{d_1 \times m}$, $\mathbf{B} \in \mathbb{R}^{d_2 \times m}$ and $\mathbf{c} \in \mathbb{R}^m$. Given the particular structure (sometimes called splitting) between the variables \mathbf{u} and \mathbf{v} , one may want to update those variables separately rather than gathering them into a single vector \mathbf{x} and performing a single update.

This idea is precisely that of **dual decomposition**. At iteration k , the dual decomposition method applied to problem (4.3.1) computes

$$\begin{cases} \mathbf{u}_{k+1} \in \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{d_1}} \mathcal{L}(\mathbf{u}, \mathbf{v}_k, \mathbf{y}_k) \\ \mathbf{v}_{k+1} \in \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^{d_2}} \mathcal{L}(\mathbf{u}_k, \mathbf{v}, \mathbf{y}_k) \\ \mathbf{y}_{k+1} = \mathbf{y}_k + \alpha_k(\mathbf{A}\mathbf{u}_{k+1} + \mathbf{B}\mathbf{v}_{k+1} - \mathbf{c}), \end{cases} \quad (4.3.2)$$

where $\alpha_k > 0$. Interestingly, the calculations for \mathbf{u}_{k+1} and \mathbf{v}_{k+1} are completely independent, and can be carried out in parallel. This observation and its practical realization have led to successful applications of dual decomposition in several fields.

4.3.2 ADMM

The **Alternated Direction Method of Multipliers**, or **ADMM**, is an increasingly popular variation on the augmented Lagrangian paradigm that bears some connection with coordinate descent approaches, in that it splits the problem in two sets of variables.

Recall problem (4.3.1) above. For any $\lambda > 0$, the augmented Lagrangian of problem (4.3.1) has the form

$$\mathcal{L}^a(\mathbf{u}, \mathbf{v}, \mathbf{y}; \lambda) = f(\mathbf{u}) + g(\mathbf{v}) + \mathbf{y}^\top (\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} - \mathbf{c}) + \frac{\lambda}{2} \|\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} - \mathbf{c}\|^2.$$

The ADMM iteration exploits the separable nature of the problem by computing the values \mathbf{u} and \mathbf{v} independently. Starting from $(\mathbf{u}_k, \mathbf{v}_k, \mathbf{y}_k)$, the ADMM counterpart to iteration (4.2.3) is

$$\begin{cases} \mathbf{u}_{k+1} \in \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{d_1}} \mathcal{L}^a(\mathbf{u}, \mathbf{v}_k, \mathbf{y}_k; \lambda) \\ \mathbf{v}_{k+1} \in \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^{d_2}} \mathcal{L}^a(\mathbf{u}_{k+1}, \mathbf{v}, \mathbf{y}_k; \lambda) \\ \mathbf{y}_{k+1} = \mathbf{y}_k + \lambda(\mathbf{A}\mathbf{u}_{k+1} + \mathbf{B}\mathbf{v}_{k+1} - \mathbf{c}). \end{cases} \quad (4.3.3)$$

The first two updates of the iteration (4.3.3) cannot be run in parallel, but the philosophy is slightly different from that of the dual decomposition method. Indeed, in ADMM, it is common that solving for a subset of the variables will be much easier than solving for all variables at once (see our example in the next section). The ADMM framework gives the possibility to exploit this property within an augmented Lagrangian method.

Remark 4.3.1 *The idea of splitting the objective and the constraints across two groups of variables can be declined into as many groups of variables as possible, depending on the structure of the problem.*

To end this section, we briefly mention that there exist convergence results for ADMM-type frameworks, typically under convexity assumptions on the problem [3]. A typical result consist in showing that

$$\begin{cases} \|\mathbf{A}\mathbf{u}_k + \mathbf{B}\mathbf{v}_k - \mathbf{c}\| & \rightarrow 0 \\ f(\mathbf{u}_k) + g(\mathbf{v}_k) & \rightarrow \min_{\mathbf{u}, \mathbf{v}} f(\mathbf{u}) + g(\mathbf{v}) \\ \mathbf{y}_k & \rightarrow \mathbf{y}^*, \end{cases}$$

where \mathbf{y}^* is a solution of the dual problem.

4.4 Decentralized optimization

We end this chapter by describing an increasingly common setup in optimization over large datasets, often termed [consensus optimization](#) or [decentralized optimization](#). In this setup, we consider a dataset that is split across m entities called *agents*. Every agent uses its own data to train a certain learning model parameterized by a vector in \mathbb{R}^d . To this end, each agent not only has its own function $f^{(i)}$, but also its own copy of the model parameters $\mathbf{x}^{(i)}$. The optimization problem at hand considers a master iterate \mathbf{x} , and attempts to reach consensus between all the agents. This leads to the following formulation:

$$\begin{aligned} & \text{minimize}_{\mathbf{x}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d} && \sum_{i=1}^m f^{(i)}(\mathbf{x}^{(i)}) \\ & \text{subject to} && \mathbf{x} = \mathbf{x}^{(i)} \quad \forall i = 1, \dots, m. \end{aligned} \quad (4.4.1)$$

This problem is a proxy for $\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \sum_{i=1}^m f^{(i)}(\mathbf{x})$, but the latter problem cannot be solved by a single agent since every agent has exclusive access to its data by design. The formulation (4.4.1)

models the fact that all agents are involved in computing \mathbf{x} by acting on \mathbf{x}_i . It is possible to apply ADMM to problem (4.4.1) by setting

$$\mathbf{u} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \vdots \\ \mathbf{x}^{(m)} \end{bmatrix} \in \mathbb{R}^{md}, \quad \mathbf{v} = \mathbf{x} \in \mathbb{R}^d.$$

Generalization The idea behind the formulation (4.4.1) can be extended to the case of data spread over a network, represented by a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$: every vertex $s \in \mathcal{V}$ of the graph represents an agent, while every edge $(s, s') \in \mathcal{E}$ represents a channel of communication between two agents in the graph. Letting $\mathbf{x}^{(s)} \in \mathbb{R}^d$ and $f^{(s)} : \mathbb{R}^d \rightarrow \mathbb{R}$ represent the parameter copy and objective function for agent $s \in \mathcal{V}$, respectively, the consensus optimization problem can be written as:

$$\begin{aligned} & \text{minimize}_{\{\mathbf{x}^{(s)}\}_{s \in \mathcal{V}} \in (\mathbb{R}^d)^{|\mathcal{V}|}} \sum_{s \in \mathcal{V}} f^{(s)}(\mathbf{x}^{(s)}) \\ & \text{subject to} \quad \mathbf{x}^{(s)} = \mathbf{x}^{(s')} \quad \forall (s, s') \in \mathcal{E}. \end{aligned} \quad (4.4.2)$$

When the graph is fully connected, i.e. all agents communicate, this problem reduces to an unconstrained problem. However, in general, the solutions of this problem are much difficult to identify, and one must work through minimizing the objective and satisfying the so-called consensus constraints.

Decentralized gradient methods To wrap up this chapter, we describe an increasingly popular class of algorithms that extends gradient descent to the decentralized setting. The decentralized gradient framework is designed for problems of the form

$$\text{minimize}_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d} \sum_{i=1}^m f_i(\mathbf{x}^{(i)}),$$

without consensus constraints but with an implicit graph structure $(\mathcal{V}, \mathcal{E})$ connecting the agents. Given a matrix $\mathbf{W} \in \mathbb{R}^{m \times m}$ that is doubly stochastic (i.e. with nonnegative coefficients such that the sum of all rows and all columns are 1) and satisfies $[\mathbf{W}]_{ij} \neq 0$ if and only if $i = j$ or $(i, j) \in \mathcal{E}$, the k th iteration of the decentralized gradient method at agent i reads

$$\mathbf{x}_{k+1}^{(i)} = \sum_{j=1}^m [\mathbf{W}]_{ij} \mathbf{x}_k^{(j)} - \alpha_k \nabla f_i(\mathbf{x}_k^{(i)}). \quad (4.4.3)$$

The iteration (4.4.3) thus combines a gradient step for agent i together with a so-called **mixing step** (or consensus step) in which the current value of the iterate for agent i is combined with that of its neighbors. This framework is steadily gaining popularity in the machine learning community.

4.5 Conclusion

In modern data science tasks, the amount of data available requires distributed storage, and possibly agents cooperating in order to solve the optimization problem at hand. Linearly constrained formulations can capture this behavior, and such constraints can be handled in an efficient manner using dual variables. Augmented Lagrangian techniques are among the most popular methods in this category, and these methods can be further specialized to account for structure in the problem. The ADMM framework has emerged as one of the most interesting formulations used to split the calculation into (presumably) cheaper subproblems. It is also very well suited for operating in a distributed or decentralized environment.

Chapter 5

Exercises

Exercise 1: Proximal gradient and regularization

In this exercise, we revisit the proximal operator and the proximal gradient method on a specific problem. Given a point $w \in \mathbb{R}^d$, we consider

$$\text{minimize}_{x \in \mathbb{R}^d} \|x\|_1 + \frac{1}{2\alpha} \|x - w\|_2^2, \quad (5.0.1)$$

where $\|x\|_1 = \sum_{i=1}^d |[x]_i|$, $\|x\|_2^2 = \sum_{i=1}^d [x]_i^2$ and $\alpha > 0$.

- Using the properties of $x \mapsto \|x\|_1$, explain why the objective function of (5.0.1) cannot be optimized by gradient-type techniques.
- Justify that problem (5.0.1) and

$$\text{minimize}_{x \in \mathbb{R}^d} \alpha \|x\|_1 + \frac{1}{2} \|x - w\|_2^2 \quad (5.0.2)$$

have the same solution set (argmin). Since both functions are strongly convex, what can be said about this solution set?

- Write down an optimality condition for problem (5.0.1).
- In this question, we view problem (5.0.1) as computing a proximal operator.
 - Using the definition of the proximal operator, write down the solution of problem (5.0.1) as the value of a proximal operator of a certain function.
 - By repeatedly solving instances of problem (5.0.1) using the last solution found as w , what algorithm do we obtain?
- In this question, we view problem (5.0.2) in a composite form, where $\frac{1}{2} \|x - w\|_2^2$ is a data-fitting term and $\alpha \|x\|_1$ is a regularization term.
 - What is the purpose of such a regularization term? Why is it computationally worth using?
 - Write down an iteration of proximal gradient applied to this problem with $x_k = w$ and stepsize α . What do you observe then? Is it to be expected?

Exercise 2: Second-difference regularization

In this exercise, we consider a regularized optimization problem of the form

$$\text{minimize}_{x \in \mathbb{R}^d} f(x) + \frac{\lambda}{2} \|Lx\|_2^2, \quad (5.0.3)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable, $\|v\|_2^2 = \sum_{j=1}^d [v]_j^2$ for any $v \in \mathbb{R}^d$ and $L \in \mathbb{R}^{d \times d}$ is the second-difference matrix defined by

$$L_{ij} = \begin{cases} 1 & \text{if } j = i + 1 \text{ or } j = i - 1 \\ -2 & \text{if } j = i \\ 0 & \text{otherwise.} \end{cases}$$

Such problems arise when the vector x represents a discretization of a real-valued function. Using the proposed regularization promotes solutions whose components vary continuously.

a) In this question, we assume that $d \gg 1$ and that the function f is separable, in that it can be written as

$$f(x) = \sum_{j=1}^d f^j([x]_j),$$

where every $f^j : \mathbb{R} \rightarrow \mathbb{R}$ only depends on the j th coordinate of the output vector x .

- i) Justify that the second term in the objective is partially separable (which implies that its gradient also is).
 - ii) Recall that the gradient of the second term is given by λLx for any $x \in \mathbb{R}^d$. Suppose that we apply a basic coordinate descent method to problem (5.0.3). Justify how an iteration of such a method can be considered cheaper than an iteration of gradient descent.
 - iii) Suggest a block variant of that method that makes use of the structure of the second term in the objective.
- b) In this question, we suppose that f is strongly convex. We modify problem (5.0.3) by introducing an auxiliary variable $z \in \mathbb{R}^d$, leading to

$$\begin{aligned} & \text{minimize}_{\substack{x \in \mathbb{R}^d \\ z \in \mathbb{R}^d}} f(x) + \frac{\lambda}{2} \|z\|_2^2 \\ & \text{s.t.} \quad Lx - z = 0. \end{aligned} \quad (5.0.4)$$

- i) Write down the Lagrangian for problem (5.0.4). Using this function, how can we rewrite problem (5.0.4)?
 - ii) What is the difference between a Lagrangian and an augmented Lagrangian?
 - iii) How does the introduction of the variable z allows for applying ADMM to this problem? What is the advantage of such an approach here?
- c) Finally, we suppose that the objective function f can be expressed as a finite sum

$$f(x) = \sum_{i=1}^n f_i(x),$$

where every f_i is strongly convex and continuously differentiable. We consider that all f_i are spread across different agents, but that all agents know the regularization term.

- i) Consider first the formulation (5.0.3). Rewrite this problem under the assumption that every agent has its own copy of the problem variable, and is using its own function f_i instead of f .
- ii) Using the same idea as for obtaining (5.0.4), explain how the problem from the previous question can be reformulated as a linearly-constrained problem in order to apply ADMM.

Bibliography

- [1] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4:1–106, 2012.
- [2] A. Beck. *First-Order Methods in Optimization*. MPS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics, Philadelphia, 2017.
- [3] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3:1–122, 2010.
- [4] N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1:123–231, 2013.
- [5] S. J. Wright. Coordinate descent algorithms. *Math. Program.*, 151:3–34, 2015.
- [6] S. J. Wright and B. Recht. *Optimization for Data Analysis*. Cambridge University Press, 2022.