

Machine Learning Optimization Algorithms & Portfolio Allocation*

Sarah Perrin
Artificial Intelligence & Advanced
Visual Computing Master
Ecole Polytechnique
sarah.perrin@polytechnique.edu

Thierry Roncalli
Quantitative Research
Amundi Asset Management
Paris
thierry.roncalli@amundi.com

June 2019

Abstract

Portfolio optimization emerged with the seminal paper of Markowitz (1952). The original mean-variance framework is appealing because it is very efficient from a computational point of view. However, it also has one well-established failing since it can lead to portfolios that are not optimal from a financial point of view (Michaud, 1989). Nevertheless, very few models have succeeded in providing a real alternative solution to the Markowitz model. The main reason lies in the fact that most academic portfolio optimization models are intractable in real life although they present solid theoretical properties. By intractable we mean that they can be implemented for an investment universe with a small number of assets using a lot of computational resources and skills, but they are unable to manage a universe with dozens or hundreds of assets. However, the emergence and the rapid development of robo-advisors means that we need to re-think portfolio optimization and go beyond the traditional mean-variance optimization approach.

Another industry and branch of science has faced similar issues concerning large-scale optimization problems. Machine learning and applied statistics have long been associated with linear and logistic regression models. Again, the reason was the inability of optimization algorithms to solve high-dimensional industrial problems. Nevertheless, the end of the 1990s marked an important turning point with the development and the rediscovery of several methods that have since produced impressive results. The goal of this paper is to show how portfolio allocation can benefit from the development of these large-scale optimization algorithms. Not all of these algorithms are useful in our case, but four of them are essential when solving complex portfolio optimization problems. These four algorithms are the coordinate descent, the alternating direction method of multipliers, the proximal gradient method and the Dykstra's algorithm. This paper reviews them and shows how they can be implemented in portfolio allocation.

Keywords: Portfolio allocation, mean-variance optimization, risk budgeting optimization, quadratic programming, coordinate descent, alternating direction method of multipliers, proximal gradient method, Dykstra's algorithm.

JEL classification: C61, G11.

*This survey article has been prepared for the book Machine Learning and Asset Management edited by Emmanuel Jurczenko. We would like to thank Mohammed El Mendili, Edmond Lezmi, Lina Mezghani, Jean-Charles Richard, Jules Roche and Jiali Xu for their helpful comments.

1 Introduction

The contribution of Harry Markowitz to economics is considerable. The mean-variance optimization framework marks the beginning of portfolio allocation in finance. In addition to the seminal paper of 1952, Harry Markowitz proposed an algorithm for solving quadratic programming problems in 1956. At that time, very few people were aware of this optimization framework. We can cite Mann (1943) and Martin (1955), but it is widely accepted that Harry Markowitz is the “*father of quadratic programming*” (Cottle and Infanger, 2010). This is not the first time that economists are participating in the development of mathematics¹, but this is certainly the first time that mathematicians will explore a field of research, whose main application during the first years of research is exclusively an economic problem².

The success of mean-variance optimization (MVO) is due to the appealing properties of the quadratic utility function, but it should also be assessed in light of the success of quadratic programming (QP). Because it is easy to solve QP problems and because QP problems are available in mathematical software, solving MVO problems is straightforward and does not require a specific skill. This is why the mean-variance optimization is a universal method which is used by all portfolio managers. However, this approach has been widely criticized by academics and professionals. Indeed, mean-variance optimization is very sensitive to input parameters and produces corner solutions. Moreover, the concept of mean-variance diversification is confused with the concept of hedging (Bourgeron *et al.*, 2018). These different issues make the practice of mean-variance optimization less attractive than the theory (Michaud, 1989). In fact, solving MVO allocation problems requires the right weight constraints to be specified in order to obtain acceptable solutions. It follows that designing the constraints is the most important component of mean-variance optimization. In this case, MVO appears to be a trial-and-error process, not a systematic solution.

The success of the MVO framework is also explained by the fact that there are very few competing portfolio allocation models that can be implemented from an industrial point of view. There are generally two reasons for this. The first one is that some models use input parameters that are difficult to estimate or understand, making these models definitively unusable. The second reason is that other models use a more complex objective function than the simple quadratic utility function. In this case, the computational complexity makes these models less attractive than the standard MVO model. Among these models, some of them are based on the mean-variance objective function, but introduce regularization penalty functions in order to improve the robustness of the portfolio allocation. Again, these models have little chance of being used if they cannot be cast into a QP problem. However, new optimization algorithms have emerged for solving large-scale machine learning problems. The purpose of this article is to present these new mathematical methods and show that they can be easily applied to portfolio allocation in order to go beyond the MVO/QP model.

This survey article is based on several previous research papers ([8], [37], [61] and [62]) and extensively uses four leading references (Beck, 2017; Boyd *et al.*, 2010; Combettes and Pesquet, 2011; Tibshirani, 2017). It is organized as follows. In section two, we present the mean-variance approach and how it is related to the QP framework. The third section is dedicated to large-scale optimization algorithms that have been used in machine learning: coordinate descent, alternating direction method of multipliers, proximal gradient and Dykstra’s algorithm. Section four shows how these algorithms can be implemented in order to solve portfolio optimization problems and build a more robust asset allocation. Finally, section five offers some concluding remarks.

¹For example, Leonid Kantorovich made major contributions to the success of linear programming.

²If we consider the first publications on quadratic programming, most of them were published in *Econometrica* or illustrated the Markowitz problem (see [1], [2], [4], [24], [26], [30], [40], [72] and [73]).

2 The quadratic programming world of portfolio optimization

2.1 Quadratic programming

2.1.1 Primal formulation

A quadratic programming (QP) problem is an optimization problem with a quadratic objective function and linear inequality constraints:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top Q x - x^\top R \\ \text{s.t. } & Sx \leq T \end{aligned} \quad (1)$$

where x is a $n \times 1$ vector, Q is a $n \times n$ matrix and R is a $n \times 1$ vector. We note that the system of constraints $Sx \leq T$ allows us to specify linear equality constraints³ $Ax = B$ or box constraints $x^- \leq x \leq x^+$. Most numerical packages then consider the following formulation:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top Q x - x^\top R \\ \text{s.t. } & \begin{cases} Ax = B \\ Cx \leq D \\ x^- \leq x \leq x^+ \end{cases} \end{aligned} \quad (2)$$

because the problem (2) is equivalent to the canonical problem (1) with the following system of linear inequalities:

$$\begin{bmatrix} -A \\ A \\ C \\ -I_n \\ I_n \end{bmatrix} x \leq \begin{bmatrix} -B \\ B \\ D \\ -x^- \\ x^+ \end{bmatrix}$$

If the space Ω defined by $Sx \leq T$ is non-empty and if Q is a symmetric positive definite matrix, the solution exists because the function $f(x) = \frac{1}{2} x^\top Q x - x^\top R$ is convex. In the general case where Q is a square matrix, the solution may not exist.

2.1.2 Dual formulation

The Lagrange function is equal to:

$$\mathcal{L}(x; \lambda) = \frac{1}{2} x^\top Q x - x^\top R + \lambda^\top (Sx - T)$$

We deduce that the dual problem is defined by:

$$\begin{aligned} \lambda^* &= \arg \max_{\lambda} \left\{ \inf_x \mathcal{L}(x; \lambda) \right\} \\ \text{s.t. } & \lambda \geq 0 \end{aligned}$$

³This is equivalent to impose that $Ax \geq B$ and $Ax \leq B$.

We note that $\partial_x \mathcal{L}(x; \lambda) = Qx - R + S^\top \lambda$. The solution to the equation $\partial_x \mathcal{L}(x; \lambda) = 0$ is then $x = Q^{-1}(R - S^\top \lambda)$. We finally obtain:

$$\begin{aligned} \inf_x \mathcal{L}(x; \lambda) &= \frac{1}{2} (R^\top - \lambda^\top S) Q^{-1} (R - S^\top \lambda) - (R^\top - \lambda^\top S) Q^{-1} R + \\ &\quad \lambda^\top (SQ^{-1} (R - S^\top \lambda) - T) \\ &= \frac{1}{2} R^\top Q^{-1} R - \lambda^\top SQ^{-1} R + \frac{1}{2} \lambda^\top SQ^{-1} S^\top \lambda - R^\top Q^{-1} R + \\ &\quad 2\lambda^\top SQ^{-1} R - \lambda^\top SQ^{-1} S^\top \lambda - \lambda^\top T \\ &= -\frac{1}{2} \lambda^\top SQ^{-1} S^\top \lambda + \lambda^\top (SQ^{-1} R - T) - \frac{1}{2} R^\top Q^{-1} R \end{aligned}$$

We deduce that the dual program is another quadratic programming problem:

$$\begin{aligned} \lambda^* &= \arg \min_{\lambda} \frac{1}{2} \lambda^\top \bar{Q} \lambda - \lambda^\top \bar{R} \\ \text{s.t. } &\lambda \geq 0 \end{aligned} \tag{3}$$

where $\bar{Q} = SQ^{-1}S^\top$ and $\bar{R} = SQ^{-1}R - T$.

Remark 1 *This duality property is very important for some machine learning methods. For example, this is the case of support vector machines and kernel methods that extensively use the duality for defining the solution (Cortes and Vapnik, 1995).*

2.1.3 Numerical algorithms

There is a substantial literature on the methods for solving quadratic programming problems (Gould and Toint, 2000). The research begins in the 1950s with different key contributions: Frank and Wolfe (1956), Markowitz (1956), Beale (1959) and Wolfe (1959). Nowadays, QP problems are generally solved using three approaches: active set methods, gradient projection methods and interior point methods. All these algorithms are implemented in standard mathematical programming languages (Matlab, Mathematica, Python, Gauss, R, etc.). This explains the success of QP problems since 2000s, because they can be easily and rapidly solved.

2.2 Mean-variance optimized portfolios

The concept of portfolio allocation has a long history and dates back to the seminal work of Markowitz (1952). In his paper, Markowitz defined precisely what *portfolio selection* means: “the investor does (or should) consider expected return a desirable thing and variance of return an undesirable thing”. Indeed, Markowitz showed that an efficient portfolio is the portfolio that maximizes the expected return for a given level of risk (corresponding to the variance of portfolio return) or a portfolio that minimizes the risk for a given level of expected return. Even if this framework has been extended to many other allocation problems (index sampling, turnover management, etc.), the mean-variance model remains the optimization approach that is the most widely used in finance.

2.2.1 The Markowitz framework

We consider a universe of n assets. Let $x = (x_1, \dots, x_n)$ be the vector of weights in the portfolio. We assume that the portfolio is fully invested meaning that $\sum_{i=1}^n x_i = \mathbf{1}_n^\top x = 1$. We denote $\mathfrak{R} = (\mathfrak{R}_1, \dots, \mathfrak{R}_n)$ as the vector of asset returns where \mathfrak{R}_i is the return of asset

i. The return of the portfolio is then equal to $\mathfrak{R}(x) = \sum_{i=1}^n x_i \mathfrak{R}_i = x^\top \mathfrak{R}$. Let $\mu = \mathbb{E}[\mathfrak{R}]$ and $\Sigma = \mathbb{E}[(\mathfrak{R} - \mu)(\mathfrak{R} - \mu)^\top]$ be the vector of expected returns and the covariance matrix of asset returns. The expected return of the portfolio is equal to:

$$\mu(x) = \mathbb{E}[\mathfrak{R}(x)] = x^\top \mu$$

whereas its variance is equal to:

$$\sigma^2(x) = \mathbb{E}[(\mathfrak{R}(x) - \mu(x))(\mathfrak{R}(x) - \mu(x))^\top] = x^\top \Sigma x$$

Markowitz (1952) formulated the investor's financial problem as follows:

1. Maximizing the expected return of the portfolio under a volatility constraint (σ -problem):

$$\max \mu(x) \quad \text{s.t.} \quad \sigma(x) \leq \sigma^* \quad (4)$$

2. Or minimizing the volatility of the portfolio under a return constraint (μ -problem):

$$\min \sigma(x) \quad \text{s.t.} \quad \mu(x) \geq \mu^* \quad (5)$$

Markowitz's bright idea was to consider a quadratic utility function:

$$\mathcal{U}(x) = x^\top \mu - \frac{\phi}{2} x^\top \Sigma x$$

where $\phi \geq 0$ is the risk aversion. Since maximizing $\mathcal{U}(x)$ is equivalent to minimizing $-\mathcal{U}(x)$, the Markowitz problems (4) and (5) can be cast into a QP problem⁴:

$$\begin{aligned} x^*(\gamma) &= \arg \min_x \frac{1}{2} x^\top \Sigma x - \gamma x^\top \mu \\ \text{s.t.} \quad &\mathbf{1}_n^\top x = 1 \end{aligned} \quad (6)$$

where $\gamma = \phi^{-1}$. Therefore, solving the μ -problem or the σ -problem is equivalent to finding the optimal value of γ such that $\mu(x^*(\gamma)) = \mu^*$ or $\sigma(x^*(\gamma)) = \sigma^*$. We know that the functions $\mu(x^*(\gamma))$ and $\sigma(x^*(\gamma))$ are increasing with respect to γ and are bounded. The optimal value of γ can then be easily computed using the bisection algorithm. It is obvious that a large part of the success of the Markowitz framework lies on the QP trick. Indeed, Problem (6) corresponds to the QP problem (2) where $Q = \Sigma$, $R = \gamma\mu$, $A = \mathbf{1}_n^\top$ and $B = 1$. Moreover, it is easy to include bounds on the weights, inequalities between asset classes, etc.

2.2.2 Solving complex MVO problems

The previous framework can be extended to other portfolio allocation problems. However, from a numerical point of view, the underlying idea is to always find an equivalent QP formulation (Roncalli, 2013).

Portfolio optimization with a benchmark We now consider a benchmark b . We note $\mu(x | b) = (x - b)^\top \mu$ as the expected excess return and $\sigma(x | b) = \sqrt{(x - b)^\top \Sigma (x - b)}$ as the tracking error volatility of Portfolio x with respect to Benchmark b . The objective

⁴This transformation is called the QP trick.

function corresponds to a trade-off between minimizing the tracking error volatility and maximizing the expected excess return (or the alpha):

$$f(x | b) = \frac{1}{2} \sigma^2(x | b) - \gamma \mu(x | b)$$

We can show that the equivalent QP problem is⁵:

$$x^*(\gamma) = \arg \min_x \frac{1}{2} x^\top \Sigma x - \gamma x^\top \tilde{\mu}$$

where $\tilde{\mu} = \mu + \gamma^{-1} \Sigma b$ is the regularized vector of expected returns. Therefore, portfolio allocation with a benchmark can be viewed as a regularization of the MVO problem and is solved using a QP numerical algorithm.

Index sampling The goal of index sampling is to replicate an index portfolio with a smaller number of assets than the index (or the benchmark) b . From a mathematical point of view, index sampling could be written as follows:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} (x - b)^\top \Sigma (x - b) \\ \text{s.t. } &\begin{cases} \mathbf{1}_n^\top x = 1 \\ x \geq \mathbf{0}_n \\ \sum_{i=1}^n \mathbb{1}\{x_i > 0\} \leq n_x \end{cases} \end{aligned} \quad (7)$$

The idea is to minimize the volatility of the tracking error such that the number of stocks n_x in the portfolio is smaller than the number of stocks n_b in the benchmark. For example, one would like to replicate the S&P 500 index with only 50 stocks and not the entire 500 stocks that compose this index. Professionals generally solve Problem (7) with the following heuristic algorithm:

1. We set $x_{(0)}^+ = \mathbf{1}_n$. At the iteration $k + 1$, we solve the QP problem:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} (x - b)^\top \Sigma (x - b) \\ \text{s.t. } &\begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq x_{(k)}^+ \end{cases} \end{aligned}$$

2. We then update the upper bound $x_{(k)}^+$ of the QP problem by deleting the asset i^* with the lowest non-zero optimized weight⁶:

$$x_{(k+1),i}^+ \leftarrow x_{(k),i}^+ \quad \text{if } i \neq i^* \quad \text{and} \quad x_{(k+1),i^*}^+ \leftarrow 0$$

3. We iterate the two steps until $\sum_{i=1}^n \mathbb{1}\{x_i^* > 0\} = n_x$.

The purpose of the heuristic algorithm is to delete one asset at each iteration in order to obtain an invested portfolio, which is exactly composed of n_x assets and has a low tracking error volatility. Again, we notice that solving the index sampling problem is equivalent to solving $(n_b - n_x)$ QP problems.

⁵See Appendix A.1 on page 51.

⁶We have $i^* = \{i : \arg \inf x_i^* | x_i^* > 0\}$.

Turnover management If we note \bar{x} as the current portfolio and x as the new portfolio, the turnover of Portfolio x with respect to Portfolio \bar{x} is the sum of purchases and sales:

$$\tau(x | \bar{x}) = \sum_{i=1}^n (x_i - \bar{x}_i)^+ + \sum_{i=1}^n (\bar{x}_i - x_i)^+ = \sum_{i=1}^n |x_i - \bar{x}_i|$$

Adding a turnover constraint in long-only MVO portfolios leads to the following problem:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top \Sigma x - \gamma x^\top \mu \\ \text{s.t.} \quad &\begin{cases} \sum_{i=1}^n x_i = 1 \\ \sum_{i=1}^n |x_i - \bar{x}_i| \leq \tau^+ \\ 0 \leq x_i \leq 1 \end{cases} \end{aligned}$$

where τ^+ is the maximum turnover with respect to the current portfolio \bar{x} . Scherer (2007) introduces the additional variables x_i^- and x_i^+ such that:

$$x_i = \bar{x}_i + x_i^+ - x_i^-$$

with $x_i^- \geq 0$ indicates a negative weight change with respect to the initial weight \bar{x}_i and $x_i^+ \geq 0$ indicates a positive weight change. The expression of the turnover becomes:

$$\sum_{i=1}^n |x_i - \bar{x}_i| = \sum_{i=1}^n |x_i^+ - x_i^-| = \sum_{i=1}^n x_i^+ + \sum_{i=1}^n x_i^-$$

because one of the variables x_i^+ or x_i^- is necessarily equal to zero due to the minimization problem. The γ -problem of Markowitz becomes:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top \Sigma x - \gamma x^\top \mu \\ \text{s.t.} \quad &\begin{cases} \sum_{i=1}^n x_i = 1 \\ x_i = \bar{x}_i + x_i^+ - x_i^- \\ \sum_{i=1}^n x_i^+ + \sum_{i=1}^n x_i^- \leq \tau^+ \\ 0 \leq x_i, x_i^-, x_i^+ \leq 1 \end{cases} \end{aligned}$$

We obtain an augmented QP problem of dimension $3n$ (see Appendix A.2 on page 51).

Transaction costs The previous analysis assumes that there is no transaction cost $\mathbf{c}(x | \bar{x})$ when we rebalance the portfolio from the current portfolio \bar{x} to the new optimized portfolio x . If we note c_i^- and c_i^+ as the bid and ask transaction costs, we have:

$$\mathbf{c}(x | \bar{x}) = \sum_{i=1}^n x_i^- c_i^- + \sum_{i=1}^n x_i^+ c_i^+$$

The net expected return of Portfolio x is then equal to $\mu(x) - \mathbf{c}(x | \bar{x})$. It follows that the γ -problem of Markowitz becomes⁷:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top \Sigma x - \gamma \left(\sum_{i=1}^n x_i \mu_i - \sum_{i=1}^n x_i^- c_i^- - \sum_{i=1}^n x_i^+ c_i^+ \right) \\ \text{s.t.} \quad &\begin{cases} \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^- c_i^- + \sum_{i=1}^n x_i^+ c_i^+ = 1 \\ x_i = \bar{x}_i + x_i^+ - x_i^- \\ 0 \leq x_i, x_i^-, x_i^+ \leq 1 \end{cases} \end{aligned}$$

Once again, we obtain a QP problem (see Appendix A.3 on page 52).

⁷The equality constraint $\mathbf{1}_n^\top x = 1$ becomes $\mathbf{1}_n^\top x + \mathbf{c}(x | \bar{x}) = 1$ because the rebalancing process has to be financed.

2.3 Issues with QP optimization

The concurrent model of the Markowitz framework is the risk budgeting approach (Qian, 2005; Maillard *et al.*, 2010; Roncalli, 2013). The goal is to define a convex risk measure $\mathcal{R}(x)$ and to allocate the risk according to some specified risk budgets $\mathcal{RB} = (\mathcal{RB}_1, \dots, \mathcal{RB}_n)$ where $\mathcal{RB}_i > 0$. This approach exploits the Euler decomposition property of the risk measure:

$$\mathcal{R}(x) = \sum_{i=1}^n x_i \frac{\partial \mathcal{R}(x)}{\partial x_i}$$

By noting $\mathcal{RC}_i(x) = x_i \cdot \partial_{x_i} \mathcal{R}(x)$ as the risk contribution of Asset i with respect to portfolio x , the risk budgeting (RB) portfolio is defined by the following set of equations:

$$x^* = \{x \in [0, 1]^n : \mathcal{RC}_i(x) = \mathcal{RB}_i\}$$

Roncalli (2013) showed that it is equivalent to solving the following non-linear optimization problem⁸:

$$\begin{aligned} x^* &= \arg \min_x \mathcal{R}(x) - \lambda \sum_{i=1}^n \mathcal{RB}_i \cdot \ln x_i \\ \text{s.t. } &x_i > 0 \end{aligned} \quad (8)$$

where $\lambda > 0$ is an arbitrary positive constant. Generally, the most frequently used risk measures are the volatility risk measure (Maillard *et al.*, 2010):

$$\mathcal{R}(x) = \sqrt{x^\top \Sigma x}$$

and the standard deviation-based risk measure (Roncalli, 2015):

$$\mathcal{R}(x) = -x^\top (\mu - r) + \xi \sqrt{x^\top \Sigma x}$$

where r is the risk-free rate and ξ is a positive scalar. In particular, this last one encompasses the Gaussian value-at-risk — $\xi = \Phi^{-1}(\alpha)$ — and the Gaussian expected shortfall — $\xi = (1 - \alpha)^{-1} \phi(\Phi^{-1}(\alpha))$.

The risk budgeting approach has displaced the MVO approach in many fields of asset management, in particular in the case of factor investing and alternative risk premia. Nevertheless, we notice that Problem (8) is not a quadratic programming problem, but a logarithmic barrier problem. Therefore, the risk budgeting framework opens a new world of portfolio optimization that is not necessarily QP! That is all the more true since MVO portfolios face robustness issues (Bourgeron *et al.*, 2018). Regularization of portfolio allocation has then become the industry standard. Indeed, it is frequent to add a ℓ_1 -norm or ℓ_2 -norm penalty functions to the MVO objective function. This type of penalty is, however, tractable in a quadratic programming setting. With the development of robo-advisors, non-linear penalty functions have emerged, in particular the logarithmic barrier penalty function. And these regularization techniques result in a non-quadratic programming world of portfolio optimization.

The success of this non-QP financial world will depend on how quickly and easily these complex optimization problems can be solved. Griveau-Billon *et al.* (2013), Bourgeron *et al.* (2018), and Richard and Roncalli (2019) have already proposed numerical algorithms that are doing the work in some special cases. The next section reviews the candidate algorithms that may compete with QP numerical algorithms.

⁸In fact, the solution x^* must be rescaled after the optimization step.

3 Machine learning optimization algorithms

The machine learning industry has experienced a similar trajectory to portfolio optimization. Before the 1990s, statistical learning focused mainly on models that were easy to solve from a numerical point of view. For instance, the linear (and the ridge) regression has an analytical solution, we can solve logistic regression with the Newton-Raphson algorithm whereas supervised and unsupervised classification models⁹ consist in performing a singular value decomposition or a generalized eigenvalue decomposition. The 1990s saw the emergence of three models that have deeply changed the machine learning approach: neural networks, support vector machines and lasso regression.

Neural networks have been extensively studied since the seminal work of Rosenblatt (1958). However, the first industrial application dates back to the publication of LeCun *et al.* (1989) on handwritten zip code recognition. At the beginning of 1990s, a fresh craze then emerged with the writing of many handbooks that were appropriate for students, the most popular of which was Bishop (1995). With neural networks, two main issues arise concerning calibration: the large number of parameters to estimate and the absence of a global maximum. The traditional numerical optimization algorithms¹⁰ that were popular in the 1980s cannot be applied to neural networks. New optimization approaches are then proposed. First, researchers have considered more complex learning rules than the steepest descent (Jacobs, 1988), for example the momentum method of Polyak (1964) or the Nesterov accelerated gradient approach (Nesterov, 1983). Second, the descent method is generally not performed on the full sample of observations, but on a subset of observations that changes at each iteration. This is the underlying idea behind batch gradient descent (BGD), stochastic descent gradient (SGD) and mini-batch gradient descent (MGD). We notice that adaptive learning methods and batch optimization techniques have marked the revival of the gradient descend method.

The development of support vector machines is another important step in the development of machine learning techniques. Like neural networks, they can be seen as an extension of the perceptron. However, they present nice theoretical properties and a strong geometrical framework. Once SVMs have been first developed for linear classification, they have been extended for non-linear classification and regression. A support vector machine consists in separating hyperplanes and finding the optimal separation by maximizing the margin. The original problem called the hard margin classification can be formulated as a quadratic programming problem. However, the dual problem, which is also a QP problem, is generally preferred to the primal problem for solving SVM classification, because of the sparse property of the solution (Cortes and Vapnik, 1995). Over the years, the original hard margin classification has been extended to other problems: soft margin classification with binary hinge loss, soft margin classification with squared hinge loss, least squares SVM regression, ε -SVM regression, kernel machines (Vapnik, 1998). All these statistical problems share the same calibration framework. The primal problem can be cast into a QP problem, implying that the corresponding dual problem is also a QP problem. Again, we notice that the success and the prominence of statistical methods are related to the efficiency of the optimization algorithms, and it is obvious that support vector machines have substantially benefited from the QP formulation. From an industrial point of view, support vector machines present however some limitations. Indeed, if the dimension of the primal QP problem is the number p of features (or parameters), the dimension of the dual QP problem is the number n of obser-

⁹e.g. principal component analysis (PCA), linear/quadratic discriminant analysis (LDA/QDA), Fisher classification method, etc.

¹⁰For example, we can cite the quasi-Newton BFGS (Broyden-Fletcher-Goldfarb-Shanno) and DFP (Davidon-Fletcher-Powell) methods, and the Fletcher-Reeves and Polak-Ribiere conjugate gradient methods.

vations. It is becoming absolutely impossible to solve the dual problem when the number of observations is larger than 100 000 and sometimes as high as several millions. This implies that new algorithms that are more appropriate for large-scale optimization problems need to be developed.

Lasso regression is the third disruptive approach that put machine learning in the spotlight in the 1990s. Like the ridge regression, lasso regression is a regularized linear regression where the ℓ_2 -norm penalty is replaced by the ℓ_1 -norm penalty (Tibshirani, 1996). Since the ℓ_1 regularization forces the solution to be sparse, it has been first largely used for variable selection, and then for pattern recognition and robust estimation of linear models. For finding the lasso solution, the technique of augmented QP problems is widely used since it is easy to implement. The extension of the lasso-ridge regularization to the other ℓ_p norms is straightforward, but these approaches have never been popular. The main reason is that existing numerical algorithms are not sufficient to make these models tractable.

Therefore, the success of a quantitative model may be explained by two conditions. First, the model must be obviously appealing. Second, the model must be solved by an efficient numerical algorithm that is easy to implement or available in mathematical programming software. As shown previously, quadratic programming and gradient descent methods have been key for many statistical and financial models. In what follows, we consider four algorithms and techniques that have been popularized by their use in machine learning: coordinate descent, alternating direction method of multipliers, proximal operators and Dykstra's algorithm. In particular, we illustrate how they can be used for solving complex optimization problems.

3.1 Coordinate descent

3.1.1 Definition

We consider the following unconstrained minimization problem:

$$x^* = \arg \min_x f(x) \quad (9)$$

where $x \in \mathbb{R}^n$ and $f(x)$ is a continuous, smooth and convex function. A popular method to find the solution x^* is to consider the descent algorithm, which is defined by the following rule:

$$x^{(k+1)} = x^{(k)} + \Delta x^{(k)} = x^{(k)} - \eta D^{(k)}$$

where $x^{(k)}$ is the approximated solution of Problem (9) at the k^{th} Iteration, $\eta > 0$ is a scalar that determines the step size and $D^{(k)}$ is the direction. We notice that the current solution $x^{(k)}$ is updated by going in the opposite direction to $D^{(k)}$ in order to obtain $x^{(k+1)}$. In the case of the gradient descent, the direction is equal to the gradient vector of $f(x)$ at the current point: $D^{(k)} = \nabla f(x^{(k)})$. Coordinate descent (CD) is a variant of the gradient descent and minimizes the function along one coordinate at each step:

$$x_i^{(k+1)} = x_i^{(k)} + \Delta x_i^{(k)} = x_i^{(k)} - \eta D_i^{(k)}$$

where $D_i^{(k)} = \nabla_i f(x^{(k)})$ is the i^{th} element of the gradient vector. At each iteration, a coordinate i is then chosen via a certain rule, while the other coordinates are assumed to be fixed. Coordinate descent is an appealing algorithm, because it transforms a vector-valued problem into a scalar-valued problem that is easier to implement. Algorithm (1) summarizes the CD algorithm. The convergence criterion can be a predefined number of iterations or an error rule between two iterations. The step size η can be either a given parameter or computed with a line search, implying that the parameter $\eta^{(k)}$ changes at each iteration.

Algorithm 1 Coordinate descent algorithm (gradient formulation)

The goal is to find the solution $x^* = \arg \min f(x)$
We initialize the vector $x^{(0)}$ and we note η the step size
Set $k \leftarrow 0$
repeat
 Choose a coordinate $i \in \{1, n\}$
 $x_i^{(k+1)} \leftarrow x_i^{(k)} - \eta \nabla_i f(x^{(k)})$
 $x_j^{(k+1)} \leftarrow x_j^{(k)}$ if $j \neq i$
 $k \leftarrow k + 1$
until convergence
return $x^* \leftarrow x^{(k)}$

Another formulation of the coordinate descent method is given in Algorithm (2). The underlying idea is to replace the descent approximation by the exact problem. Indeed, the objective of the descend step is to minimize the scalar-valued problem:

$$x_i^* = \arg \min_{\mathcal{X}} f\left(x_1^{(k)}, \dots, x_{i-1}^{(k)}, \mathcal{X}, x_{i+1}^{(k)}, \dots, x_n^{(k)}\right) \quad (10)$$

Algorithm 2 Coordinate descent algorithm (exact formulation)

The goal is to find the solution $x^* = \arg \min f(x)$
We initialize the vector $x^{(0)}$
Set $k \leftarrow 0$
repeat
 Choose a coordinate $i \in \{1, n\}$
 $x_i^{(k+1)} = \arg \min_{\mathcal{X}} f\left(x_1^{(k)}, \dots, x_{i-1}^{(k)}, \mathcal{X}, x_{i+1}^{(k)}, \dots, x_n^{(k)}\right)$
 $x_j^{(k+1)} \leftarrow x_j^{(k)}$ if $j \neq i$
 $k \leftarrow k + 1$
until convergence
return $x^* \leftarrow x^{(k)}$

Coordinate descent is efficient in large-scale optimization problems, in particular when there is a solution to the scalar-valued problem (10). Furthermore, convergence is guaranteed when $f(x)$ is convex and differentiable (Luo and Tseng, 1992; Luo and Tseng, 1993).

Remark 2 *Coordinate descent methods have been introduced in several handbooks on numerical optimization in the 1980s and 1990s (Wright, 1985). However, the most important step is the contribution of Tseng (2001), who studied the block-coordinate descent method and extended CD algorithms in the case of a non-differentiable and non-convex function $f(x)$.*

3.1.2 Cyclic or random coordinates?

There are several options for choosing the coordinate of the k^{th} iteration. A natural choice could be to choose the coordinate which minimizes the function:

$$i^* = \arg \inf \left\{ f_i^* : i \in \{1, n\}, f_i^* = \min_{\mathcal{X}} f\left((1 - e_i)x^{(k)} + e_i \mathcal{X}\right) \right\}$$

However, it is obvious that choosing the optimal coordinate i^* would require the gradient along each coordinate to be calculated. This causes the coordinate descent to be no longer

efficient, since a classic gradient descent would then be of equivalent cost at each iteration and would converge faster because it requires fewer iterations.

The simplest way to implement the CD algorithm is to consider cyclic coordinates, meaning that we cyclically iterate through the coordinates (Tseng, 2001):

$$i = k \bmod n$$

This ensures that all the coordinates are selected during one cycle $\{k - n + 1, \dots, k\}$ in the same order. This approach, called cyclical coordinate descent (CCD), is the most popular and used method, even if it is difficult to estimate the rate of convergence.

The second way is to consider random coordinates. Let π_i be the probability of choosing the coordinate i at the iteration k . The simplest approach is to consider uniform probabilities: $\pi_i = 1/n$. A better approach consists in pre-specifying probabilities according to the Lipschitz constants¹¹:

$$\pi_i = \frac{\mathfrak{L}_i^\alpha}{\sum_{j=1}^n \mathfrak{L}_j^\alpha} \quad (11)$$

Nesterov (2012) considers three schemes: $\alpha = 0$, $\alpha = 1$ and $\alpha = \infty$ — in this last case, we have $i = \arg \max \{\mathfrak{L}_1, \dots, \mathfrak{L}_n\}$. From a theoretical point of view, the random coordinate descent (RCD) method based on the probability distribution (11) leads to a faster convergence, since coordinates that have a large Lipschitz constant \mathfrak{L}_i are more likely to be chosen. However, it requires additional calculus to compute the Lipschitz constants and CCD is often preferred from a practical point of view. In what follows, we only use the CCD algorithm described below. In Algorithm (3), the variable k represents the number of cycles whereas the number of iterations is equal to $k \cdot n$. For the coordinate i , the lower coordinates $j < i$ correspond to the current cycle $k + 1$ while the upper coordinates $j > i$ correspond to the previous cycle k .

Algorithm 3 Cyclical coordinate descent algorithm

The goal is to find the solution $x^* = \arg \min f(x)$
We initialize the vector $x^{(0)}$
Set $k \leftarrow 0$
repeat
 for $i = 1 : n$ **do**
 $x_i^{(k+1)} = \arg \min_{\mathcal{X}} f\left(x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, \mathcal{X}, x_{i+1}^{(k)}, \dots, x_n^{(k)}\right)$
 end for
 $k \leftarrow k + 1$
until convergence
return $x^* \leftarrow x^{(k)}$

3.1.3 Application to the λ -problem of the lasso regression

We consider the linear regression:

$$Y = X\beta + \varepsilon \quad (12)$$

where Y is the $n \times 1$ vector, X is the $n \times p$ design matrix, β is the $p \times 1$ vector of coefficients and ε is the $n \times 1$ vector of residuals. In this model, n is the number of observations and p

¹¹Nesterov (2012) assumes that $f(x)$ is convex, differentiable and Lipschitz-smooth for each coordinate:

$$\|\nabla_i f(x + e_i h) - \nabla_i f(x)\| \leq \mathfrak{L}_i \|h\|$$

where $h \in \mathbb{R}$.

is the number of parameters (or the number of explanatory variables). The objective of the ordinary least squares is to minimize the residual sum of squares:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \text{RSS}(\beta)$$

where $\text{RSS}(\beta) = \sum_{i=1}^n \varepsilon_i^2$. Since we have:

$$\text{RSS}(\beta) = (Y - X\beta)^\top (Y - X\beta)$$

we obtain:

$$\frac{\partial f(\beta)}{\partial \beta_j} = -x_j^\top (Y - X\beta)$$

where x_j is the $n \times 1$ design vector corresponding to the j^{th} explanatory variable. Because we can write:

$$X\beta = X_{(-j)}\beta_{(-j)} + x_j\beta_j$$

where $X_{(-j)}$ and $\beta_{(-j)}$ are the design matrix and the beta vector by excluding the j^{th} explanatory variable, it follows that:

$$\begin{aligned} \frac{\partial f(\beta)}{\partial \beta_j} &= x_j^\top (X_{(-j)}\beta_{(-j)} + x_j\beta_j - Y) \\ &= x_j^\top X_{(-j)}\beta_{(-j)} + x_j^\top x_j\beta_j - x_j^\top Y \end{aligned}$$

At the optimum, we have $\partial_{\beta_j} f(\beta) = 0$ or:

$$\beta_j = \frac{x_j^\top (Y - X_{(-j)}\beta_{(-j)})}{x_j^\top x_j} \quad (13)$$

The implementation of the coordinate descent algorithm is straightforward. It suffices to iterate Equation (13) through the coordinates.

The lasso regression problem is a variant of the OLS regression by adding a ℓ_1 -norm regularization (Tibshirani, 1996):

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \frac{1}{2} (Y - X\beta)^\top (Y - X\beta) + \lambda \|\beta\|_1 \quad (14)$$

In this formulation, the residual sum of squares of the linear regression is penalized by a term that will force a sparse selection of the coordinates. Since the objective function is the sum of two convex norms, the convergence is guaranteed for the lasso problem. Because $\|\beta\|_1 = \sum_{j=1}^n |\beta_j|$, the first order condition becomes:

$$\begin{aligned} 0 &= \nabla_i f(\beta) \\ &= x_j^\top x_j\beta_j - x_j^\top (Y - X_{(-j)}\beta_{(-j)}) + \lambda \partial |\beta_j| \end{aligned}$$

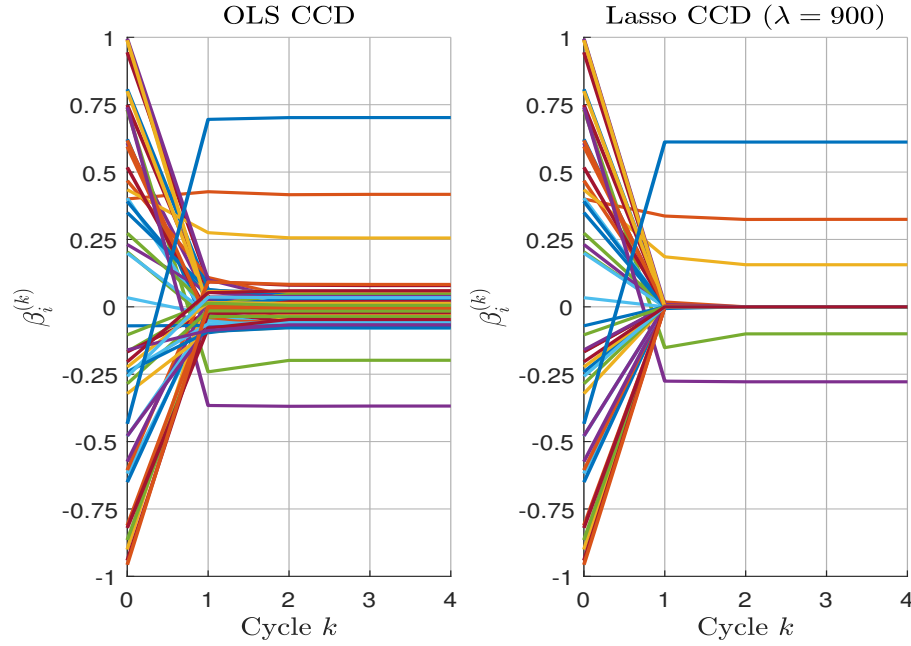
In Appendix A.5, we show that the solution is given by:

$$\beta_j = \frac{\mathcal{S}(x_j^\top (Y - X_{(-j)}\beta_{(-j)}); \lambda)}{x_j^\top x_j} \quad (15)$$

where $\mathcal{S}(v; \lambda)$ is the soft-thresholding operator:

$$\mathcal{S}(v; \lambda) = \text{sign}(v) \cdot (|v| - \lambda)_+$$

Figure 1: CCD algorithm applied to the lasso optimization problem



It follows that the lasso CD algorithm is a variation of the linear regression CD algorithm by applying the soft-threshold operator to the residuals $x_j^\top (Y - X_{(-j)}\beta_{(-j)})$ at each iteration.

Let us consider an experiment with $n = 10\,000$ and $p = 50$. The design matrix X is built using the uniform distribution while the residuals are simulated using a Gaussian distribution and a standard deviation of 20%. The beta coefficients are distributed uniformly between -3 and $+3$ except four coefficients that take a larger value. We then standardize the data of X and Y because the practice of the lasso regression is to consider comparable beta coefficients. By considering uniform numbers between -1 and $+1$ for initializing the coordinates, results of the CCD algorithm are given in Figure 1. We notice that the CCD algorithm converges quickly after three complete cycles. In the case of a large-scale problem when $p \gg 1000$, it has been shown that CCD may be faster for the lasso regression than for the OLS regression because of the soft-thresholding operator. Indeed, we can initialize the algorithm with the null vector $\mathbf{0}_p$. If λ is large, a lot of optimal coordinates are equal to zero and a few cycles are needed to find the optimal values of non-zero coefficients.

3.1.4 Application to the box-constrained QP problem

Coordinate descent can also be applied to the box-constrained QP problem:

$$x^* = \arg \min_x \frac{1}{2} x^\top Q x - x^\top R \quad \text{s.t.} \quad x^- \leq x \leq x^+ \quad (16)$$

In Appendix A.6 on page 53, we show that the coordinate update of the CCD algorithm is equal to:

$$x_i^{(k+1)} = \mathcal{T} \left(\frac{R_i - \frac{1}{2} \sum_{j < i} x_j^{(k+1)} (Q_{i,j} + Q_{j,i}) - \frac{1}{2} \sum_{j > i} x_j^{(k)} (Q_{i,j} + Q_{j,i})}{Q_{i,i}}; x_i^-, x_i^+ \right)$$

where $\mathcal{T}(v; x^-, x^+)$ is the truncation operator:

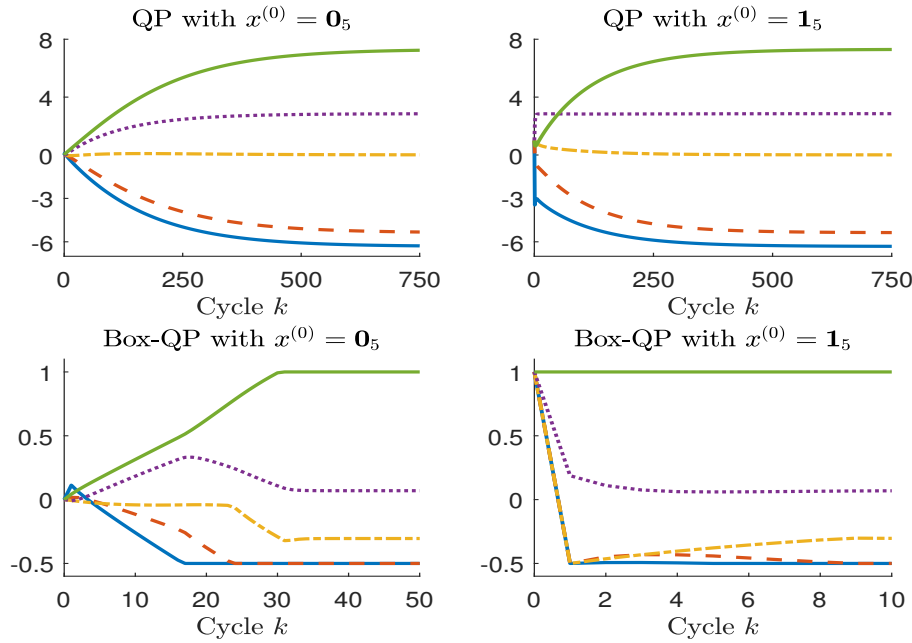
$$\begin{aligned} \mathcal{T}(v; x^-, x^+) &= v \odot \mathbb{1}\{x^- < v < x^+\} + \\ &\quad x^- \odot \mathbb{1}\{v \leq x^-\} + \\ &\quad x^+ \odot \mathbb{1}\{v \geq x^+\} \end{aligned} \quad (17)$$

Generally, we assume that Q is a symmetric matrix, implying that the CCD update reduces to:

$$x_i^{(k+1)} = \mathcal{T}\left(\frac{R_i - \sum_{j < i} x_j^{(k+1)} Q_{i,j} - \sum_{j > i} x_j^{(k)} Q_{i,j}}{Q_{i,i}}; x_i^-, x_i^+\right)$$

Remark 3 CCD can be applied to Problem (16) because the box constraint $x^- \leq x \leq x^+$ is pointwise¹².

Figure 2: CCD algorithm applied to the box-constrained QP problem



We consider the following example:

$$Q = \begin{pmatrix} 5.76 & 5.11 & 3.47 & 5.13 & 6.82 \\ 5.11 & 7.98 & 5.38 & 4.30 & 8.70 \\ 3.47 & 5.38 & 4.01 & 2.83 & 5.91 \\ 5.13 & 4.30 & 2.83 & 4.70 & 5.84 \\ 6.82 & 8.70 & 5.91 & 5.84 & 10.18 \end{pmatrix} \quad \text{and} \quad R = \begin{pmatrix} 0.65 \\ 0.72 \\ 0.46 \\ 0.59 \\ 1.26 \end{pmatrix}$$

In Figure 2, we have reported the solution obtained with the CCD algorithm. The top panels correspond to the QP problem without any constraints, whereas the bottom panel corresponds to the QP problem with the box constraint $-0.5 \leq x_i \leq 1$. We notice that

¹²See the discussion on page 24.

we need more than 500 cycles for the convergence of the CCD algorithm in the case of the unconstrained QP problem, whereas CCD finds the solution of the constrained QP problem using less than 50 cycles. We also observe that the convergence speed is highly dependent on the starting values. In the case of the box-constrained QP problem, we need 40 cyclical iterations if the starting value is the vector $x^{(0)} = \mathbf{0}_5$, whereas less than 10 cyclical iterations are sufficient if we consider the unit vector $x^{(0)} = \mathbf{1}_5$.

3.2 Alternating direction method of multipliers

3.2.1 Definition

The alternating direction method of multipliers (ADMM) is an algorithm introduced by Gabay and Mercier (1976) to solve optimization problems which can be expressed as:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x,y)} f_x(x) + f_y(y) \\ \text{s.t. } Ax + By &= c \end{aligned} \quad (18)$$

where $A \in \mathbb{R}^{p \times n}$, $B \in \mathbb{R}^{p \times m}$, $c \in \mathbb{R}^p$, and the functions $f_x : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ and $f_y : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper closed convex functions. Boyd et al. (2011) show that the ADMM algorithm consists of the following three steps:

1. The x -update is:

$$x^{(k+1)} = \arg \min_x \left\{ f_x(x) + \frac{\varphi}{2} \|Ax + By^{(k)} - c + u^{(k)}\|_2^2 \right\} \quad (19)$$

2. The y -update is:

$$y^{(k+1)} = \arg \min_y \left\{ f_y(y) + \frac{\varphi}{2} \|Ax^{(k+1)} + By - c + u^{(k)}\|_2^2 \right\} \quad (20)$$

3. The u -update is:

$$u^{(k+1)} = u^{(k)} + (Ax^{(k+1)} + By^{(k+1)} - c) \quad (21)$$

In this approach, $u^{(k)}$ is the dual variable of the primal residual $r = Ax + By - c$ and φ is the ℓ_2 -norm penalty variable. The parameter φ can be constant or may change at each iteration¹³. The ADMM algorithm benefits from the dual ascent principle and the method of multipliers. The difference with the latter is that the x - and y -updates are performed in an alternating way. Therefore, it is more flexible because the updates are equivalent to computing proximal operators for f_x and f_y independently. In practice, ADMM may be slow to converge with high accuracy, but is fast to converge if we consider modest accuracy. This is why ADMM is a good candidate for solving large-scale machine learning problems, where high accuracy does not necessarily lead to a better solution.

Remark 4 In this paper, we use the notations $f_x^{(k+1)}(x)$ and $f_y^{(k+1)}(y)$ when referring to the objective functions that are defined in the x - and y -updates. Algorithm (4) summarizes the different ADMM steps.

¹³See Appendix A.7 on page 56 for a discussion about the convergence of the ADMM algorithm.

Algorithm 4 ADMM algorithm

The goal is to compute the solution (x^*, y^*)
 We initialize the vectors $x^{(0)}$ and $y^{(0)}$ and we choose a value for the parameter φ
 We set $u^{(0)} = \mathbf{0}_n$
 $k \leftarrow 0$
repeat
 $x^{(k+1)} = \arg \min_x \left\{ f_x^{(k+1)}(x) = f_x(x) + \frac{\varphi}{2} \|Ax + By^{(k)} - c + u^{(k)}\|_2^2 \right\}$
 $y^{(k+1)} = \arg \min_y \left\{ f_y^{(k+1)}(y) = f_y(y) + \frac{\varphi}{2} \|Ax^{(k+1)} + By - c + u^{(k)}\|_2^2 \right\}$
 $u^{(k+1)} = u^{(k)} + (Ax^{(k+1)} + By^{(k+1)} - c)$
 $k \leftarrow k + 1$
until convergence
return $x^* \leftarrow x^{(k)}$ and $y^* \leftarrow y^{(k)}$

3.2.2 ADMM tricks

The appeal of ADMM is that it can separate a complex problem into two sub-problems that are easier to solve. However, most of the time, the optimization problem is not formulated using a separable objective function. The question is then how to formulate the initial problem as a separable problem. We now list some tricks that show how ADMM may be used in practice.

First trick We consider a problem of the form $x^* = \arg \min_x g(x)$. The idea is then to write $g(x)$ as a separable function $g(x) = g_1(x) + g_2(x)$ and to consider the following equivalent ADMM problem:

$$\begin{aligned}
 \{x^*, y^*\} &= \arg \min_{(x,y)} f_x(x) + f_y(y) \\
 \text{s.t. } &x = y
 \end{aligned} \tag{22}$$

where $f_x(x) = g_1(x)$ and $f_y(y) = g_2(y)$. Usually, the smooth part of $g(x)$ will correspond to $g_1(x)$ while the non-smooth part will be included in $g_2(x)$. The underlying idea is that the x -update is straightforward, whereas the y -update deals with the tricky part of $g(x)$.

Second trick If we want to minimize the function $g(x)$ where $x \in \Omega$ is a set of constraints, the optimization problem can be cast into the ADMM form (22) where $f_x(x) = g(x)$, $f_y(y) = \mathbf{1}_\Omega(y)$ and $\mathbf{1}_\Omega(x)$ is the convex indicator function of Ω :

$$\mathbf{1}_\Omega(x) = \begin{cases} 0 & \text{if } x \in \Omega \\ +\infty & \text{if } x \notin \Omega \end{cases} \tag{23}$$

For example, if we want to solve the QP problem (2) given on page 3, we have:

$$f_x(x) = \frac{1}{2}x^\top Qx - x^\top R$$

and:

$$\Omega = \{x \in \mathbb{R}^n : Ax = B, Cx \leq D, x^- \leq x \leq x^+\}$$

Third trick We can combine the first and second tricks. For instance, if we consider the following optimization problem:

$$\begin{aligned} x^* &= \arg \min_x g_1(x) + g_2(x) \\ \text{s.t. } &x \in \Omega_1 \cap \Omega_2 \end{aligned}$$

the equivalent ADMM form is:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x,y)} \underbrace{(g_1(x) + \mathbb{1}_{\Omega_1}(x))}_{f_x(x)} + \underbrace{(g_2(y) + \mathbb{1}_{\Omega_2}(y))}_{f_y(y)} \\ \text{s.t. } &x = y \end{aligned}$$

Let us consider a variant of the QP problem where we add a non-linear constraint $h(x) = 0$. In this case, we can write the set of constraints as $\Omega = \Omega_1 \cap \Omega_2$ where:

$$\Omega_1 = \{x \in \mathbb{R}^n : Ax = B, Cx \leq D, x^- \leq x \leq x^+\}$$

and:

$$\Omega_2 = \{x \in \mathbb{R}^n : h(x) = 0\}$$

Fourth trick Finally, if we want to minimize the function $g(x) = g(x, Ax + b) = g_1(x) + g_2(Ax + b)$, we can write:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x,y)} g_1(x) + g_2(y) \\ \text{s.t. } &y = Ax + b \end{aligned}$$

For instance, this trick can be used for a QP problem with a non-linear part:

$$g(x) = \frac{1}{2}x^\top Qx - x^\top R + h(x)$$

If we assume that Q is a symmetric positive-definite matrix, we set $x = Ly$ where L is the lower Cholesky matrix such that $LL^\top = Q$. It follows that the ADMM form is equal to¹⁴:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x,y)} \underbrace{\frac{1}{2}x^\top x}_{f_x(x)} + \underbrace{h(y) - y^\top R}_{f_y(y)} \\ \text{s.t. } &x - Ly = \mathbf{0}_n \end{aligned}$$

We notice that the x -update is straightforward because it corresponds to a standard QP problem. If we add a set Ω of constraints, we specify:

$$f_y(y) = h(y) - y^\top R + \mathbb{1}_\Omega(y)$$

Remark 5 In the previous cases, we have seen that when the function $g(x)$ may contain a QP problem, it is convenient to isolate this QP problem into the x -update:

$$x^{(k+1)} = \arg \min_x \left\{ \frac{1}{2}x^\top Qx - x^\top R + \mathbb{1}_\Omega(x) + \frac{\varphi}{2} \|x - y^{(k)} + u^{(k)}\|_2^2 \right\}$$

Since we have:

$$\frac{\varphi}{2} \|x - y^{(k)} + u^{(k)}\|_2^2 = \frac{\varphi}{2}x^\top x - \varphi x^\top (y^{(k)} - u^{(k)}) + \frac{\varphi}{2} (y^{(k)} - u^{(k)})^\top (y^{(k)} - u^{(k)})$$

we deduce that the x -update is a standard QP problem where:

$$f_x^{(k+1)}(x) = \frac{1}{2}x^\top (Q + \varphi I_n)x - x^\top (R + \varphi (y^{(k)} - u^{(k)})) + \mathbb{1}_\Omega(x) \quad (24)$$

¹⁴This Cholesky trick has been used by Gonzalvez et al. (2019) to solve trend-following strategies using the ADMM algorithm in the context of Bayesian learning.

3.2.3 Application to the λ -problem of the lasso regression

The λ -problem of the lasso regression (14) has the following ADMM formulation:

$$\begin{aligned} \{\beta^*, \bar{\beta}^*\} &= \arg \min \frac{1}{2} (Y - X\beta)^\top (Y - X\beta) + \lambda \|\bar{\beta}\|_1 \\ \text{s.t. } &\beta - \bar{\beta} = \mathbf{0}_p \end{aligned}$$

Since the x -step corresponds to a QP problem¹⁵, we use the results given in Remark 5 to find the value of $\beta^{(k+1)}$:

$$\begin{aligned} \beta^{(k+1)} &= (Q + \varphi I_p)^{-1} \left(R + \varphi (\bar{\beta}^{(k)} - u^{(k)}) \right) \\ &= (X^\top X + \varphi I_p)^{-1} \left(X^\top Y + \varphi (\bar{\beta}^{(k)} - u^{(k)}) \right) \end{aligned}$$

The y -step is:

$$\begin{aligned} \bar{\beta}^{(k+1)} &= \arg \min_{\bar{\beta}} \left\{ \lambda \|\bar{\beta}\|_1 + \frac{\varphi}{2} \left\| \beta^{(k+1)} - \bar{\beta} + u^{(k)} \right\|_2^2 \right\} \\ &= \arg \min \left\{ \frac{1}{2} \left\| \bar{\beta} - (\beta^{(k+1)} + u^{(k)}) \right\|_2^2 + \frac{\lambda}{\varphi} \|\bar{\beta}\|_1 \right\} \end{aligned}$$

We recognize the soft-thresholding problem with $v = \beta^{(k+1)} + u^{(k)}$. Finally, the ADMM algorithm is made up of the following steps (Boyd *et al.*, 2011):

$$\begin{cases} \beta^{(k+1)} = (X^\top X + \varphi I_p)^{-1} (X^\top Y + \varphi (\bar{\beta}^{(k)} - u^{(k)})) \\ \bar{\beta}^{(k+1)} = \mathcal{S}(\beta^{(k+1)} + u^{(k)}; \varphi^{-1} \lambda) \\ u^{(k+1)} = u^{(k)} + (\beta^{(k+1)} - \bar{\beta}^{(k+1)}) \end{cases}$$

We consider the example of the lasso regression with $\lambda = 900$ on page 14. By setting $\varphi = \lambda$ and by initialing the algorithm with the OLS estimates, we obtain the convergence given in Figure 3. We notice that the ADMM algorithm converges more slowly than the CCD algorithm for this example. In practice, we generally observe that the convergence is poor for low and very high values of φ . However, finding an optimal value of φ is difficult. A better approach involves using a varying parameter $\varphi^{(k)}$ such as the method described on page 56.

3.3 Proximal operators

The x - and y -update steps of the ADMM algorithm require a ℓ_2 -norm penalized optimization problem to be solved. Proximal operators are special cases of this type of problem when the matrices A or B correspond to the identity matrix I_n or its opposite $-I_n$.

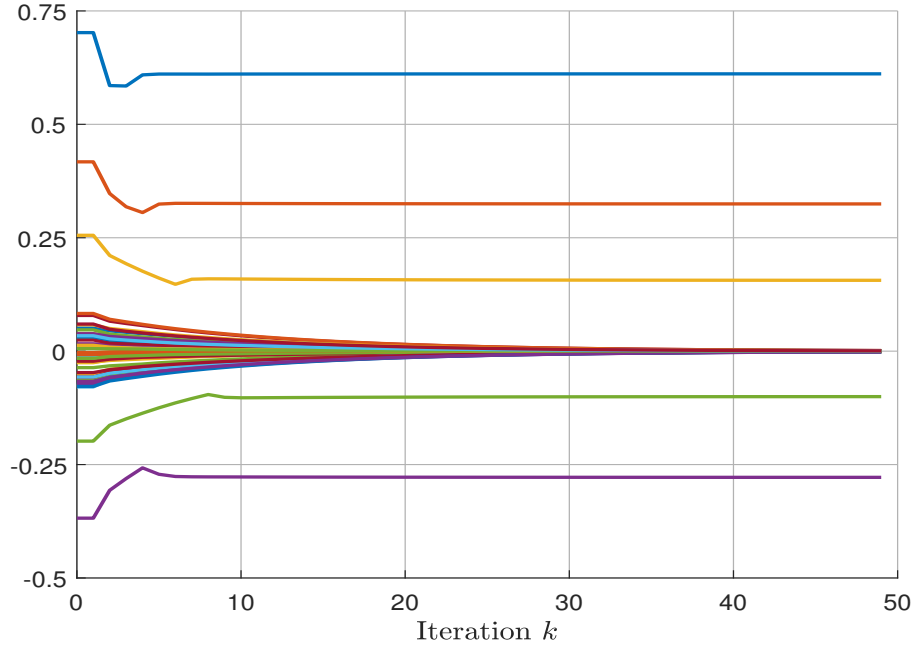
3.3.1 Definition

Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper closed convex function. The proximal operator $\text{prox}_f(v) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined by:

$$\text{prox}_f(v) = x^* = \arg \min_x \left\{ f(x) + \frac{1}{2} \|x - v\|_2^2 \right\} \quad (25)$$

¹⁵We have $Q = X^\top X$ and $R = X^\top Y$.

Figure 3: ADMM algorithm applied to the lasso optimization problem



Since the function $f_v(x) = f(x) + \frac{1}{2} \|x - v\|_2^2$ is strongly convex, it has a unique minimum for every $v \in \mathbb{R}^n$ (Parikh and Boyd, 2014). By construction, the proximal operator defines a point x^* which is a trade-off between minimizing $f(x)$ and being close to v .

In many situations, we need to calculate the proximal of the scaled function $\lambda f(x)$ where $\lambda > 0$. In this case, we use the notation $\mathbf{prox}_{\lambda f}(v)$ and we have:

$$\begin{aligned} \mathbf{prox}_{\lambda f}(v) &= \arg \min_x \left\{ \lambda f(x) + \frac{1}{2} \|x - v\|_2^2 \right\} \\ &= \arg \min_x \left\{ f(x) + \frac{1}{2\lambda} \|x - v\|_2^2 \right\} \end{aligned}$$

For instance, if we consider the y -update of the ADMM algorithm with $B = -I_n$, we have:

$$\begin{aligned} y^{(k+1)} &= \arg \min_y \left\{ f_y(y) + \frac{\varphi}{2} \|y - v_y^{(k+1)}\|_2^2 \right\} \\ &= \arg \min_y \left\{ \varphi^{-1} f_y(y) + \frac{1}{2} \|y - v_y^{(k+1)}\|_2^2 \right\} \\ &= \mathbf{prox}_{\varphi^{-1} f_y}(v_y^{(k+1)}) \end{aligned}$$

where $v_y^{(k)} = Ax^{(k+1)} - c + u^{(k)}$. ADMM is then given by Algorithm (5). The interest of this mathematical formulation is to write the ADMM algorithm in a convenient form such that the x -update corresponds to the tricky part of the optimization while the y -update is reduced to an analytical formula.

Algorithm 5 ADMM algorithm in the case $Ax - y = c$

The goal is to compute the solution (x^*, y^*)
 We initialize the vectors $x^{(0)}$ and $y^{(0)}$ and we choose a value for the parameter φ
 We set $u^{(0)} = \mathbf{0}_n$
 $k \leftarrow 0$
repeat
 $x^{(k+1)} = \arg \min_x \left\{ f_x^{(k+1)}(x) = f_x(x) + \frac{\varphi}{2} \|Ax - y^{(k)} - c + u^{(k)}\|_2^2 \right\}$
 $v_y^{(k+1)} = Ax^{(k+1)} - c + u^{(k)}$
 $y^{(k+1)} = \mathbf{prox}_{\varphi^{-1}f_y} \left(v_y^{(k+1)} \right)$
 $u^{(k+1)} = u^{(k)} + (Ax^{(k+1)} - y^{(k+1)} - c)$
 $k \leftarrow k + 1$
until convergence
return $x^* \leftarrow x^{(k)}$ and $y^* \leftarrow y^{(k)}$

3.3.2 Proximal operators and generalized projections

In the case where $f(x) = \mathbf{1}_\Omega(x)$ is the indicator function, the proximal operator is then the Euclidean projection onto Ω :

$$\begin{aligned}
 \mathbf{prox}_f(v) &= \arg \min_x \left\{ \mathbf{1}_\Omega(x) + \frac{1}{2} \|x - v\|_2^2 \right\} \\
 &= \arg \min_{x \in \Omega} \left\{ \|x - v\|_2^2 \right\} \\
 &= \mathcal{P}_\Omega(v)
 \end{aligned}$$

where $\mathcal{P}_\Omega(v)$ is the standard projection of v onto Ω . Parikh and Boyd (2014) interpret then proximal operators as a generalization of the Euclidean projection.

Let us consider the constrained optimization problem $x^* = \arg \min f(x)$ subject to $x \in \Omega$. Using the second ADMM trick, we have $f_x(x) = f(x)$, $f_y(y) = \mathbf{1}_\Omega(y)$ and $x - y = \mathbf{0}_n$. Therefore, we can use Algorithm (5) since the v - and y -steps become $v_y^{(k+1)} = x^{(k+1)} + u^{(k)}$ and¹⁶ $y^{(k+1)} = \mathcal{P}_\Omega(v_y^{(k+1)})$.

Here, we give the results of Parikh and Boyd (2014) for some simple polyhedra:

Notation	Ω	$\mathcal{P}_\Omega(v)$
$\mathcal{A}_{ffineset}[A, B]$	$Ax = B$	$v - A^\dagger(Av - B)$
$\mathcal{H}_{hyperplane}[a, b]$	$a^\top x = b$	$v - \frac{(a^\top v - b)}{\ a\ _2^2} a$
$\mathcal{H}_{halfspace}[c, d]$	$c^\top x \leq d$	$v - \frac{(c^\top v - d)_+}{\ c\ _2^2} c$
$\mathcal{B}_{box}[x^-, x^+]$	$x^- \leq x \leq x^+$	$\mathcal{T}(v; x^-, x^+)$

where A^\dagger is the Moore-Penrose pseudo-inverse of A , and $\mathcal{T}(v; x^-, x^+)$ is the truncation operator.

¹⁶We notice that the parameter φ has no impact on the y -update because $\varphi^{-1}f_y(y) = f_y(y) = \mathbf{1}_\Omega(y)$. We then deduce that:

$$\mathbf{prox}_{\varphi^{-1}f_y} \left(v_y^{(k+1)} \right) = \mathbf{prox}_{f_y} \left(v_y^{(k+1)} \right) = \mathcal{P}_\Omega \left(v_y^{(k+1)} \right)$$

3.3.3 Main properties

There are many properties that are useful for finding the analytical expression of the proximal operator. In what follows, we consider three main properties, but the reader may refer to Combettes and Pesquet (2011), Parikh and Boyd (2014) and Beck (2017) for a more exhaustive list.

Separable sum Let us assume that $f(x) = \sum_{i=1}^n f_i(x_i)$ is fully separable, then the proximal of $f(v)$ is the vector of the proximal operators applied to each scalar-valued function $f_i(x_i)$:

$$\mathbf{prox}_f(v) = \begin{pmatrix} \mathbf{prox}_{f_1}(v_1) \\ \vdots \\ \mathbf{prox}_{f_n}(v_n) \end{pmatrix}$$

For example, if $f(x) = \lambda \|x\|_1$, we have $f(x) = \lambda \sum_{i=1}^n |x_i|$ and $f_i(x_i) = \lambda |x_i|$. We deduce that the proximal operator of $f(x)$ is the vector formulation of the soft-thresholding operator:

$$\mathbf{prox}_{\lambda \|x\|_1}(v) = \begin{pmatrix} \text{sign}(v_1) \cdot (|v_1| - \lambda)_+ \\ \vdots \\ \text{sign}(v_n) \cdot (|v_n| - \lambda)_+ \end{pmatrix} = \text{sign}(v) \odot (|v| - \lambda \mathbf{1}_n)_+$$

This result has been used to solve the λ -problem of the lasso regression on page 19.

If we consider the scalar-valued logarithmic barrier function $f(x) = -\lambda \ln x$, we have:

$$\begin{aligned} f_v(x) &= -\lambda \ln x + \frac{1}{2}(x - v)^2 \\ &= -\lambda \ln x + \frac{1}{2}x^2 - xv + \frac{1}{2}v^2 \end{aligned}$$

The first-order condition is $-\lambda x^{-1} + x - v = 0$. We obtain two roots with opposite signs:

$$x^* = \frac{v \pm \sqrt{v^2 + 4\lambda}}{2}$$

Since the logarithmic function is defined for $x > 0$, we deduce that the proximal operator is the positive root. In the case of the vector-valued logarithmic barrier $f(x) = -\lambda \sum_{i=1}^n \ln x_i$, it follows that:

$$\mathbf{prox}_f(v) = \frac{v + \sqrt{v \odot v + 4\lambda}}{2}$$

Moreau decomposition An important property of the proximal operator is the Moreau decomposition theorem:

$$\mathbf{prox}_f(v) + \mathbf{prox}_{f^*}(v) = v$$

where f^* is the convex conjugate of f . This result is used extensively to find the proximal of norms, the max function, the sum-of- k -largest-values function, etc. (Beck, 2017).

In the case of the pointwise maximum function $f(x) = \max x$, we can show that:

$$\mathbf{prox}_{\lambda \max x}(v) = \min(v, s^*)$$

where s^* is the solution of the following equation (see Appendix A.8.1 on page 57):

$$s^* = \left\{ s \in \mathbb{R} : \sum_{i=1}^n (v_i - s)_+ = \lambda \right\}$$

If we assume that $f(x) = \|x\|_p$, we obtain:

p	$\mathbf{prox}_{\lambda f}(v)$
$p = 1$	$\mathcal{S}(v; \lambda) = \text{sign}(v) \odot (v - \lambda \mathbf{1}_n)_+$
$p = 2$	$\left(1 - \frac{\lambda}{\max(\lambda, \ v\ _2)}\right) v$
$p = \infty$	$\text{sign}(v) \odot \mathbf{prox}_{\lambda \max x}(v)$

If $f(x)$ is a ℓ_q -norm function, then $f^*(x) = \mathbf{1}_{\mathcal{B}_p}(x)$ where \mathcal{B}_p is the ℓ_p unit ball and $p^{-1} + q^{-1} = 1$. Since we have $\mathbf{prox}_{f^*}(v) = \mathcal{P}_{\mathcal{B}_p}(v)$, we deduce that:

$$\mathbf{prox}_f(v) + \mathcal{P}_{\mathcal{B}_p}(v) = v$$

More generally, we have:

$$\mathbf{prox}_{\lambda f}(v) + \lambda \mathcal{P}_{\mathcal{B}_p}\left(\frac{v}{\lambda}\right) = v$$

It follows that the projection onto the ℓ_p ball can be deduced from the proximal operator of the ℓ_q -norm function. Let $\mathcal{B}_p(c, \lambda) = \{x \in \mathbb{R}^n : \|x - c\|_p \leq \lambda\}$ be the ℓ_p ball with center c and radius λ . We obtain:

p	$\mathcal{P}_{\mathcal{B}_p(\mathbf{0}_n, \lambda)}(v)$	q
$p = 1$	$v - \text{sign}(v) \odot \mathbf{prox}_{\lambda \max x}(v)$	$q = \infty$
$p = 2$	$v - \mathbf{prox}_{\lambda \ x\ _2}(v)$	$q = 2$
$p = \infty$	$\mathcal{T}(v; -\lambda, \lambda)$	$q = 1$

Scaling and translation Let us define $g(x) = f(ax + b)$ where $a \neq 0$. We have¹⁷:

$$\mathbf{prox}_g(v) = \frac{\mathbf{prox}_{a^2 f}(av + b) - b}{a}$$

We can use this property when the center c of the ℓ_p ball is not equal to $\mathbf{0}_n$. Since we have $\mathbf{prox}_g(v) = \mathbf{prox}_f(v - c) + c$ where $g(x) = f(x - c)$ and the equivalence $\mathcal{B}_p(\mathbf{0}_n, \lambda) = \{x \in \mathbb{R}^n : f(x) \leq \lambda\}$ where $f(x) = \|x\|_p$, we deduce that:

$$\mathcal{P}_{\mathcal{B}_p(c, \lambda)}(v) = \mathcal{P}_{\mathcal{B}_p(\mathbf{0}_n, \lambda)}(v - c) + c$$

3.3.4 Application to the τ -problem of the lasso regression

We have previously presented the lasso regression problem by considering the Lagrange formulation (λ -problem). We now consider the original τ -problem:

$$\begin{aligned} \hat{\beta}(\tau) &= \arg \min_{\beta} \frac{1}{2} (Y - X\beta)^\top (Y - X\beta) \\ \text{s.t. } &\|\beta\|_1 \leq \tau \end{aligned}$$

The ADMM formulation is:

$$\begin{aligned} \{\beta^*, \bar{\beta}^*\} &= \arg \min_{(\beta, \bar{\beta})} \frac{1}{2} (Y - X\beta)^\top (Y - X\beta) + \mathbf{1}_\Omega(\bar{\beta}) \\ \text{s.t. } &\beta = \bar{\beta} \end{aligned}$$

¹⁷The proof can be found in Beck (2017) on page 138. We have reported it in Appendix A.8.3 on page 58.

where $\Omega = \mathcal{B}_1(\mathbf{0}_n, \tau)$ is the centered ℓ_1 ball with radius τ . We notice that the x -update is:

$$\begin{aligned}\beta^{(k+1)} &= \arg \min_{\beta} \left\{ \frac{1}{2} (Y - X\beta)^\top (Y - X\beta) + \frac{\varphi}{2} \left\| \beta - \bar{\beta}^{(k)} + u^{(k)} \right\|_2^2 \right\} \\ &= (X^\top X + \varphi I_p)^{-1} \left(X^\top Y + \varphi (\bar{\beta}^{(k)} - u^{(k)}) \right)\end{aligned}$$

where $v_x^{(k+1)} = \bar{\beta}^{(k)} - u^{(k)}$. For the y -update, we deduce that:

$$\begin{aligned}\bar{\beta}^{(k+1)} &= \arg \min_{\bar{\beta}} \left\{ \mathbf{1}_\Omega(\bar{\beta}) + \frac{\varphi}{2} \left\| \beta^{(k+1)} - \bar{\beta} + u^{(k)} \right\|_2^2 \right\} \\ &= \mathbf{prox}_{f_y} \left(\beta^{(k+1)} + u^{(k)} \right) \\ &= \mathcal{P}_\Omega \left(v_y^{(k+1)} \right) \\ &= v_y^{(k+1)} - \text{sign} \left(v_y^{(k+1)} \right) \odot \mathbf{prox}_{\tau \max x} \left(\left\| v_y^{(k+1)} \right\| \right)\end{aligned}$$

where $v_y^{(k+1)} = \beta^{(k+1)} + u^{(k)}$. Finally, the u -update is defined by $u^{(k+1)} = u^{(k)} + \beta^{(k+1)} - \bar{\beta}^{(k+1)}$.

Remark 6 The ADMM algorithm is similar for λ - and τ -problems since the only difference concerns the y -step. For the λ -problem, we apply the soft-thresholding operator while we use the ℓ_1 projection in the case of the τ -problem. However, our experience shows that the τ -problem is easier to solve with the ADMM algorithm from a practical point of view. The reason is that the y -update of the τ -problem is independent of the penalization parameter φ . This is not the case for the λ -problem, because the soft-thresholding depends on the value taken by $\varphi^{-1}\lambda$.

3.3.5 Application to the CD algorithm with pointwise constraints

We consider the following constrained minimization problem:

$$x^* = \arg \min_x f(x) \quad \text{s.t.} \quad x \in \Omega$$

where the set Ω of constraints is fully separable:

$$\mathbf{1}_\Omega(x) = \sum_{i=1}^n \mathbf{1}_{\Omega_i}(x_i)$$

The scalar-valued problem of the CD algorithm becomes:

$$x_i^* = \arg \min_{\varkappa} f(x_1, \dots, x_{i-1}, \varkappa, x_{i+1}, \dots, x_n) + \lambda \sum_{i=1}^n \mathbf{1}_{\Omega_i}(x_i)$$

Nesterov (2012) and Wright (2015) propose the following coordinate update:

$$x_i^* = \arg \min_{\varkappa} (\varkappa - x_i) g_i + \frac{1}{2\eta} (\varkappa - x_i)^2 + \lambda \cdot \mathbf{1}_{\Omega_i}(\varkappa)$$

where $g_i = \nabla_i f(x)$ is the first-derivative of the function with respect to x_i , $\eta > 0$ controls the quadratic penalization term and λ is a positive scalar. The objective function is equivalent

to:

$$\begin{aligned}
 (*) &= (\varkappa - x_i) g_i + \frac{1}{2\eta} (\varkappa - x_i)^2 + \lambda \cdot \mathbb{1}_{\Omega_i}(\varkappa) \\
 &= \frac{1}{2\eta} \left((\varkappa - x_i)^2 + 2(\varkappa - x_i) \eta g_i \right) + \lambda \cdot \mathbb{1}_{\Omega_i}(\varkappa) \\
 &= \frac{1}{2\eta} (\varkappa - x_i + \eta g_i)^2 + \lambda \cdot \mathbb{1}_{\Omega_i}(\varkappa) - \frac{\eta}{2} g_i^2
 \end{aligned}$$

By taking $\lambda = \eta^{-1}$, we deduce that:

$$\begin{aligned}
 x_i^* &= \arg \min_{\varkappa} \mathbb{1}_{\Omega_i}(\varkappa) + \frac{1}{2} \|\varkappa - (x_i - \eta g_i)\|^2 \\
 &= \mathbf{prox}_{\psi}(x_i - \eta g_i) \\
 &= \mathcal{P}_{\Omega_i}(x_i - \eta g_i)
 \end{aligned}$$

where $\psi(\varkappa) = \mathbb{1}_{\Omega_i}(\varkappa)$. Extending the CD algorithm in the case of pointwise constraints is then equivalent to implement a standard CD algorithm and apply the projection onto the i^{th} coordinate at each iteration¹⁸. For instance, this algorithm is particularly efficient when we consider box constraints.

3.4 Dykstra's algorithm

We now consider the proximal optimization problem where the function $f(x)$ is the convex sum of basic functions $f_j(x)$:

$$x^* = \arg \min_x \left\{ \sum_{j=1}^m f_j(x) + \frac{1}{2} \|x - v\|_2^2 \right\}$$

and the proximal of each basic function is known.

3.4.1 The $m = 2$ case

In the previous section, we listed some analytical solutions of the proximal problem when the function $f(x)$ is basic. For instance, we know the proximal solution of the ℓ_1 -norm function $f_1(x) = \lambda_1 \|x\|_1$ or the proximal solution of the logarithmic barrier function $f_2(x) = \lambda_2 \sum_{i=1}^n \ln x_i$. However, we don't know how to compute the proximal operator of $f(x) = f_1(x) + f_2(x)$:

$$\begin{aligned}
 x^* &= \arg \min_x f_1(x) + f_2(x) + \frac{1}{2} \|x - v\|_2^2 \\
 &= \mathbf{prox}_f(v)
 \end{aligned}$$

Nevertheless, an elegant solution is provided by the Dykstra's algorithm (Dykstra, 1983; Bauschke and Borwein, 1994; Combettes and Pesquet, 2011), which is defined by the following iterations:

$$\begin{cases} x^{(k+1)} = \mathbf{prox}_{f_1}(y^{(k)} + p^{(k)}) \\ p^{(k+1)} = y^{(k)} + p^{(k)} - x^{(k+1)} \\ y^{(k+1)} = \mathbf{prox}_{f_2}(x^{(k+1)} + q^{(k)}) \\ q^{(k+1)} = x^{(k+1)} + q^{(k)} - y^{(k+1)} \end{cases} \quad (26)$$

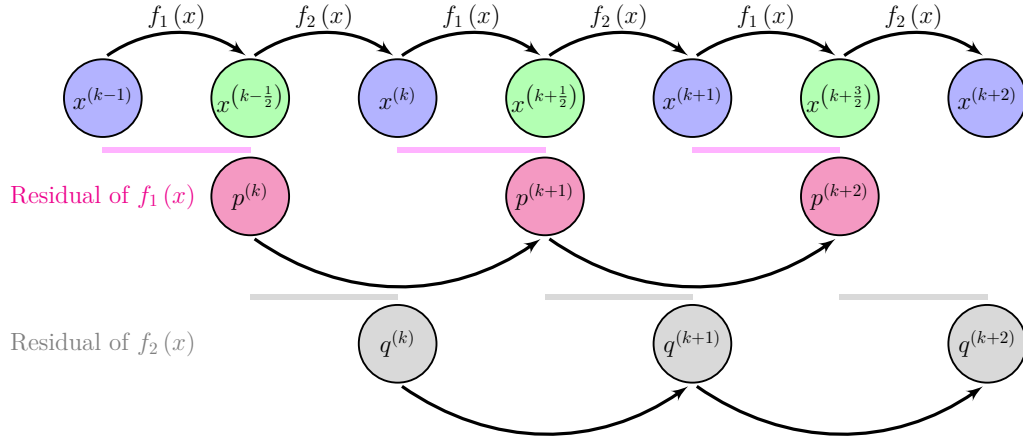
¹⁸This method corresponds to the proximal gradient algorithm.

where $x^{(0)} = y^{(0)} = v$ and $p^{(0)} = q^{(0)} = \mathbf{0}_n$. This algorithm is obviously related to the Douglas-Rachford splitting framework¹⁹ where $x^{(k)}$ and $p^{(k)}$ are the variable and the residual associated to $f_1(x)$, and $y^{(k)}$ and $q^{(k)}$ are the variable and the residual associated to $f_2(x)$. Algorithm (26) can be reformulated by introducing the intermediary step $k + \frac{1}{2}$:

$$\begin{cases} x^{(k+\frac{1}{2})} = \text{prox}_{f_1}(x^{(k)} + p^{(k)}) \\ p^{(k+1)} = p^{(k)} - \Delta_{1/2}x^{(k+\frac{1}{2})} \\ x^{(k+1)} = \text{prox}_{f_2}(x^{(k+\frac{1}{2})} + q^{(k)}) \\ q^{(k+1)} = q^{(k)} - \Delta_{1/2}x^{(k+1)} \end{cases} \quad (27)$$

where $\Delta_h x^{(k)} = x^{(k)} - x^{(k-h)}$. Figure 4 illustrates the splitting method used by the Dykstra's algorithm and clearly shows the relationship with the Douglas-Rachford algorithm.

Figure 4: Splitting method of the Dykstra's algorithm



3.4.2 The $m > 2$ case

The case $m > 2$ is a generalization of the previous algorithm by considering m residuals:

1. The x -update is:

$$x^{(k+1)} = \text{prox}_{f_{j(k)}}(x^{(k)} + z^{(k+1-m)})$$

2. The z -update is:

$$z^{(k+1)} = x^{(k)} + z^{(k+1-m)} - x^{(k+1)}$$

where $x^{(0)} = v$, $z^{(k)} = \mathbf{0}_n$ for $k < 0$ and $j(k) = \text{mod}(k+1, m)$ denotes the modulo operator taking values in $\{1, \dots, m\}$. The variable $x^{(k)}$ is updated at each iteration while the residual $z^{(k)}$ is updated every m iterations. This implies that the basic function $f_j(x)$ is related to the residuals $z^{(j)}$, $z^{(j+m)}$, $z^{(j+2m)}$, etc. Following Tibshirani (2017), it is better to write the Dykstra's algorithm by using two iteration indices k and j . The main index k refers to the cycle²⁰, whereas the sub-index j refers to the constraint number:

¹⁹See Douglas and Rachford (1956), Combettes and Pesquet (2011), and Lindstrom and Sims (2018).

²⁰Exactly like the coordinate descent algorithm.

1. The x -update is:

$$x^{(k+1,j)} = \mathbf{prox}_{f_j} \left(x^{(k+1,j-1)} + z^{(k,j)} \right) \quad (28)$$

2. The z -update is:

$$z^{(k+1,j)} = x^{(k+1,j-1)} + z^{(k,j)} - x^{(k+1,j)} \quad (29)$$

where $x^{(1,0)} = v$, $z^{(k,j)} = \mathbf{0}_n$ for $k = 0$ and $x^{(k+1,0)} = x^{(k,m)}$.

The Dykstra's algorithm is particularly efficient when we consider the projection problem:

$$x^* = \mathcal{P}_\Omega(v)$$

where:

$$\Omega = \Omega_1 \cap \Omega_2 \cap \dots \cap \Omega_m$$

Indeed, the solution is found by replacing Equation (28) with:

$$x^{(k+1,j)} = \mathcal{P}_{\Omega_j} \left(x^{(k+1,j-1)} + z^{(k,j)} \right) \quad (30)$$

3.4.3 Application to general linear constraints

Let us consider the case $\Omega = \{x \in \mathbb{R}^n : Cx \leq D\}$ where the number of inequality constraints is equal to m . We can write:

$$\Omega = \Omega_1 \cap \Omega_2 \cap \dots \cap \Omega_m$$

where $\Omega_j = \{x \in \mathbb{R}^n : c_{(j)}^\top x \leq d_{(j)}\}$, $c_{(j)}^\top$ corresponds to the j^{th} row of C and $d_{(j)}$ is the j^{th} element of D . Since the projection \mathcal{P}_{Ω_j} is known and has been given on page 21, we can find the projection \mathcal{P}_Ω using Algorithm (6).

Algorithm 6 Dykstra's algorithm for solving the proximal problem with linear inequality constraints

The goal is to compute the solution $x^* = \mathbf{prox}_f(v)$ where $f(x) = \mathbf{1}_\Omega(x)$ and $\Omega = \{x \in \mathbb{R}^n : Cx \leq D\}$

We initialize $x^{(0,m)} \leftarrow v$

We set $z^{(0,1)} \leftarrow \mathbf{0}_n, \dots, z^{(0,m)} \leftarrow \mathbf{0}_n$

$k \leftarrow 0$

repeat

$x^{(k+1,0)} \leftarrow x^{(k,m)}$

for $j = 1 : m$ **do**

The x -update is:

$$x^{(k+1,j)} = x^{(k+1,j-1)} + z^{(k,j)} - \frac{\left(c_{(j)}^\top x^{(k+1,j-1)} + c_{(j)}^\top z^{(k,j)} - d_{(j)} \right)}{\|c_{(j)}\|_2^2} c_{(j)}$$

The z -update is:

$$z^{(k+1,j)} = x^{(k+1,j-1)} + z^{(k,j)} - x^{(k+1,j)}$$

end for

$k \leftarrow k + 1$

until Convergence

return $x^* \leftarrow x^{(k,m)}$

If we define Ω as follows:

$$\Omega = \{x \in \mathbb{R}^n : Ax = B, Cx \leq D, x^- \leq x \leq x^+\}$$

we decompose Ω as the intersection of three basic convex sets:

$$\Omega = \Omega_1 \cap \Omega_2 \cap \Omega_3$$

where $\Omega_1 = \{x \in \mathbb{R}^n : Ax = B\}$, $\Omega_2 = \{x \in \mathbb{R}^n : Cx \leq D\}$ and $\Omega_3 = \{x \in \mathbb{R}^n : x^- \leq x \leq x^+\}$. Using Dykstra's algorithm is equivalent to formulating Algorithm (7).

Algorithm 7 Dykstra's algorithm for solving the proximal problem with general linear constraints

The goal is to compute the solution $x^* = \mathbf{prox}_f(v)$ where $f(x) = \mathbb{1}_\Omega(x)$ and $\Omega = \{x \in \mathbb{R}^n : Ax = B, Cx \leq D, x^- \leq x \leq x^+\}$

We initialize $x_m^{(0)} \leftarrow v$

We set $z_1^{(0)} \leftarrow \mathbf{0}_n$, $z_2^{(0)} \leftarrow \mathbf{0}_n$ and $z_3^{(0)} \leftarrow \mathbf{0}_n$

$k \leftarrow 0$

repeat

$$x_0^{(k+1)} \leftarrow x_m^{(k)}$$

$$x_1^{(k+1)} \leftarrow x_0^{(k+1)} + z_1^{(k)} - A^\dagger (Ax_0^{(k+1)} + Az_1^{(k)} - B)$$

$$z_1^{(k+1)} \leftarrow x_0^{(k+1)} + z_1^{(k)} - x_1^{(k+1)}$$

$$x_2^{(k+1)} \leftarrow \mathcal{P}_{\Omega_2} \left(x_1^{(k+1)} + z_2^{(k)} \right)$$

► Algorithm (6)

$$z_2^{(k+1)} \leftarrow x_1^{(k+1)} + z_2^{(k)} - x_2^{(k+1)}$$

$$x_3^{(k+1)} \leftarrow \mathcal{T} \left(x_2^{(k+1)} + z_3^{(k)}; x^-, x^+ \right)$$

$$z_3^{(k+1)} \leftarrow x_2^{(k+1)} + z_3^{(k)} - x_3^{(k+1)}$$

$$k \leftarrow k + 1$$

until Convergence

return $x^* \leftarrow x_3^{(k)}$

Since we have:

$$\frac{1}{2} \|x - v\|_2^2 = \frac{1}{2} x^\top x - x^\top v + \frac{1}{2} v^\top v$$

we deduce that the two previous problems can be cast into a QP problem:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top I_n x - x^\top v \\ \text{s.t. } &x \in \Omega \end{aligned}$$

We can then compare the efficiency of Dykstra's algorithm with the QP algorithm. Let us consider the proximal problem where the vector v is defined by the elements $v_i = \ln(1 + i^2)$ and the set of constraints is:

$$\Omega = \left\{ x \in \mathbb{R}^n : \sum_{i=1}^n x_i \leq \frac{1}{2}, \sum_{i=1}^n e^{-i} x_i \geq 0 \right\}$$

Using a Matlab implementation²¹, we find that the computational time of the Dykstra's algorithm when n is equal to 10 million is equal to the QP algorithm when n is equal to 12 500, meaning that there is a factor of 800 between the two methods!

²¹The QP implementation corresponds to the `quadprog` function.

3.4.4 Application to the ℓ_2 -penalized logarithmic barrier function

We consider the following proximal problem:

$$\begin{aligned} x^* &= \arg \min_x -\lambda \sum_{i=1}^n b_i \ln x_i + \frac{1}{2} \|x - v\|_2^2 \\ \text{s.t. } &\|x - c\|_2 \leq r \end{aligned}$$

In Appendix A.8.6 on page 59, we show that the corresponding Dykstra's algorithm is:

$$\begin{cases} x^{(k+1)} = \frac{y^{(k)} + z_1^{(k)} + \sqrt{(y^{(k)} + z_1^{(k)}) \odot (y^{(k)} + z_1^{(k)}) + 4\lambda b}}{2} \\ z_1^{(k+1)} = y^{(k)} + z_1^{(k)} - x^{(k+1)} \\ y^{(k+1)} = c + \frac{r}{\max(r, \|x^{(k+1)} + z_2^{(k)} - c\|_2)} (x^{(k+1)} + z_2^{(k)} - c) \\ z_2^{(k+1)} = x^{(k+1)} + z_2^{(k)} - y^{(k+1)} \end{cases}$$

4 Applications to portfolio optimization

The development of the previous algorithms will fundamentally change the practice of portfolio optimization. Until now, we have seen that portfolio managers live in a quadratic programming world. With these new optimization algorithms, we can consider more complex portfolio optimization programs with non-quadratic objective function, regularization with penalty functions and non-linear constraints.

Table 1: Some objective functions used in portfolio optimization

Item	Portfolio	$f(x)$	Reference
(1)	MVO	$\frac{1}{2}x^\top \Sigma x - \gamma x^\top \mu$	Markowitz (1952)
(2)	GMV	$\frac{1}{2}x^\top \Sigma x$	Jagganathan and Ma (2003)
(3)	MDP	$\ln(\sqrt{x^\top \Sigma x}) - \ln(x^\top \sigma)$	Choueifaty and Coignard (2008)
(4)	KL	$\sum_{i=1}^n x_i \ln(x_i/\tilde{x}_i)$	Bera and Park (2008)
(5)	ERC	$\frac{1}{2}x^\top \Sigma x - \lambda \sum_{i=1}^n \ln x_i$	Maillard <i>et al.</i> (2010)
(6)	RB	$\mathcal{R}(x) - \lambda \sum_{i=1}^n \mathcal{R}\mathcal{B}_i \cdot \ln x_i$	Roncalli (2015)
(7)	RQE	$\frac{1}{2}x^\top D x$	Carmichael <i>et al.</i> (2018)

We consider a universe of n assets. Let x be the vector of weights in the portfolio. We denote by μ and Σ the vector of expected returns and the covariance matrix of asset returns²². Some models consider also a reference portfolio \tilde{x} . In Table 1, we report the main objective functions that are used by professionals²³. Besides the mean-variance optimized portfolio (MVO) and the global minimum variance portfolio (GMV), we find the equal risk contribution portfolio (ERC), the risk budgeting portfolio (RB) and the most diversified portfolio (MDP). According to Choueifaty and Coignard (2008), the MDP is defined as the portfolio which maximizes the diversification ratio $\mathcal{DR}(x) = \frac{x^\top \sigma}{\sqrt{x^\top \Sigma x}}$. We also include in

²²The vector of volatilities is defined by $\sigma = (\sigma_1, \dots, \sigma_n)$.

²³For each model, we write the optimization problem as a minimization problem.

the list two ‘*academic*’ portfolios, which are based on the Kullback-Leibler (KL) information criteria and the Rao’s quadratic entropy (RQE) measure²⁴.

Table 2: Some regularization penalties used in portfolio optimization

Item	Regularization	$\mathfrak{R}(x)$	Reference
(8)	Ridge	$\lambda \ x - \tilde{x}\ _2^2$	DeMiguel <i>et al.</i> (2009)
(9)	Lasso	$\lambda \ x - \tilde{x}\ _1$	Brodie <i>et al.</i> (2009)
(10)	Log-barrier	$-\sum_{i=1}^n \lambda_i \ln x_i$	Roncalli (2013)
(11)	Shannon’s entropy	$\lambda \sum_{i=1}^n x_i \ln x_i$	Yu <i>et al.</i> (2014)

In a similar way, we list in Table 2 some popular regularization penalty functions that are used in the industry (Bruder *et al.*, 2013; Bourgeron *et al.*, 2018). The ridge and lasso regularization are well-known in statistics and machine learning (Hastie *et al.*, 2009). The log-barrier penalty function comes from the risk budgeting optimization problem, whereas Shannon’s entropy is another approach for imposing a sufficient weight diversification.

Table 3: Some constraints used in portfolio optimization

(12)	No cash and leverage	$\sum_{i=1}^n x_i = 1$
(13)	No short selling	$x_i \geq 0$
(14)	Weight bounds	$x_i^- \leq x_i \leq x_i^+$
(15)	Asset class limits	$c_j^- \leq \sum_{i \in \mathcal{C}_j} x_i \leq c_j^+$
(16)	Turnover	$\sum_{i=1}^n x_i - \tilde{x}_i \leq \tau^+$
(17)	Transaction costs	$\sum_{i=1}^n (c_i^- (\tilde{x}_i - x_i)_+ + c_i^+ (x_i - \tilde{x}_i)_+) \leq \mathbf{c}^+$
(18)	Leverage limit	$\sum_{i=1}^n x_i \leq \mathcal{L}^+$
(19)	Long/short exposure	$-\mathcal{LS}^- \leq \sum_{i=1}^n x_i \leq \mathcal{LS}^+$
(20)	Benchmarking	$\sqrt{(x - \tilde{x})^\top \Sigma (x - \tilde{x})} \leq \sigma^+$
(21)	Tracking error floor	$\sqrt{(x - \tilde{x})^\top \Sigma (x - \tilde{x})} \geq \sigma^-$
(22)	Active share floor	$\frac{1}{2} \sum_{i=1}^n x_i - \tilde{x}_i \geq \mathcal{AS}^-$
(23)	Number of active bets	$(x^\top x)^{-1} \geq \mathcal{N}^-$

Concerning the constraints, the most famous are the no cash/no leverage and no short selling restrictions. Weight bounds and asset class limits are also extensively used by practitioners. Turnover and transaction cost management may be an important topic when rebalancing a current portfolio \tilde{x} . When managing long/short portfolios, we generally impose leverage or long/short exposure limits. In the case of a benchmarked strategy, we might also want to have a tracking error limit with respect to the benchmark \tilde{x} . On the contrary, we can impose a minimum tracking error or active share in the case of active management. Finally, the Herfindahl constraint is used for some smart beta portfolios.

In what follows, we consider several portfolio optimization problems. Most of them are a combination of an objective function, one or two regularization penalty functions and some constraints that have been listed above. From an industrial point of view, it is interesting to implement the proximal operator for each item. In this approach, solving any portfolio optimization problem is equivalent to using CCD, ADMM, Dykstra and the appropriate proximal functions as Lego bricks.

²⁴ D is the dissimilarity matrix satisfying $D_{i,j} \geq 0$ and $D_{i,j} = D_{j,i}$.

4.1 Minimum variance optimization

4.1.1 Managing diversification

The global minimum variance (GMV) portfolio corresponds to the following optimization program:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top \Sigma x \\ \text{s.t. } & \mathbf{1}_n^\top x = 1 \end{aligned}$$

We know that the solution is $x^* = (\mathbf{1}_n^\top \Sigma^{-1} \mathbf{1}_n)^{-1} \Sigma^{-1} \mathbf{1}_n$. In practice, nobody implements the GMV portfolio because it is a long/short portfolio and it is not robust. Most of the time, professionals impose weight bounds: $0 \leq x_i \leq x^+$. However, this approach generally leads to corner solutions, meaning that a large number of optimized weights are equal to zero or the upper bound and very few assets have a weight within the range. With the emergence of smart beta portfolios, the minimum variance portfolio gained popularity among institutional investors. For instance, we can find many passive indices based on this framework. In order to increase the robustness of these portfolios, the first generation of minimum variance strategies has used relative weight bounds with respect to a benchmark b :

$$\delta^- b_i \leq x_i \leq \delta^+ b_i \quad (31)$$

where $0 < \delta^- < 1$ and $\delta^+ > 1$. For instance, the most popular scheme is to take $\delta^- = 0.5$ and $\delta^+ = 2$. Nevertheless, the constraint (31) produces the illusion that the portfolio is diversified, because the optimized weights are different. In fact, portfolio weights are different because benchmark weights are different. The second generation of minimum variance strategies imposes a global diversification constraint. The most popular solution is based on the Herfindahl index $\mathcal{H}(x) = \sum_{i=1}^n x_i^2$. This index takes the value 1 for a pure concentrated portfolio ($\exists i : x_i = 1$) and $1/n$ for an equally-weighted portfolio. Therefore, we can define the number of effective bets as the inverse of the Herfindahl index (Meucci, 2009): $\mathcal{N}(x) = \mathcal{H}(x)^{-1}$. The optimization program becomes:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top \Sigma x \\ \text{s.t. } & \begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq x^+ \\ \mathcal{N}(x) \geq \mathcal{N}^- \end{cases} \end{aligned} \quad (32)$$

where \mathcal{N}^- is the minimum number of effective bets.

The Herfindahl constraint is equivalent to:

$$\begin{aligned} \mathcal{N}(x) \geq \mathcal{N}^- &\Leftrightarrow (x^\top x)^{-1} \geq \mathcal{N}^- \\ &\Leftrightarrow x^\top x \leq \frac{1}{\mathcal{N}^-} \end{aligned}$$

Therefore, a first solution to solve (32) is to consider the following QP problem²⁵:

$$\begin{aligned} x^*(\lambda) &= \arg \min_x \frac{1}{2} x^\top \Sigma x + \lambda x^\top x \\ \text{s.t. } & \begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq x^+ \end{cases} \end{aligned} \quad (33)$$

²⁵The objective function can be written as:

$$\frac{1}{2} x^\top \Sigma x + \lambda x^\top x = \frac{1}{2} x^\top (\Sigma + 2I_n) x$$

where $\lambda \geq 0$ is a scalar. Since $\mathcal{N}(x^*(\infty))$ is equal to the number n of assets and $\mathcal{N}(x^*(\lambda))$ is an increasing function of λ , Problem (33) has a unique solution if $\mathcal{N}^- \in [\mathcal{N}(x^*(0)), n]$. There is an optimal value λ^* such that for each $\lambda \geq \lambda^*$, we have $\mathcal{N}(x^*(\lambda)) \geq \mathcal{N}^-$. Computing the optimal portfolio $x^*(\lambda^*)$ therefore implies finding the solution λ^* of the non-linear equation²⁶ $\mathcal{N}(x^*(\lambda)) = \mathcal{N}^-$.

A second method is to consider the ADMM form:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x, y)} \frac{1}{2} x^\top \Sigma x + \mathbf{1}_{\Omega_1}(x) + \mathbf{1}_{\Omega_2}(y) \\ \text{s.t. } &x = y \end{aligned}$$

where $\Omega_1 = \{x \in \mathbb{R}^n : \mathbf{1}_n^\top x = 1, \mathbf{0}_n \leq x \leq x^+\}$ and $\Omega_2 = \mathcal{B}_2(\mathbf{0}_n, \sqrt{\frac{1}{\mathcal{N}^-}})$. We deduce that the x -update is a QP problem:

$$x^{(k+1)} = \arg \min_x \left\{ \frac{1}{2} x^\top (\Sigma + \varphi I_n) x - \varphi x^\top (y^{(k)} - u^{(k)}) + \mathbf{1}_{\Omega_1}(x) \right\}$$

whereas the y -update is:

$$y^{(k+1)} = \frac{x^{(k+1)} + u^{(k)}}{\max\left(1, \sqrt{\mathcal{N}^-} \|x^{(k+1)} + u^{(k)}\|_2\right)}$$

A better approach is to write the problem as follows:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x, y)} \frac{1}{2} x^\top \Sigma x + \mathbf{1}_{\Omega_3}(x) + \mathbf{1}_{\Omega_4}(y) \\ \text{s.t. } &x = y \end{aligned}$$

where $\Omega_3 = \mathcal{H}_{\text{hyperplane}}[\mathbf{1}_n, 1]$ and $\Omega_4 = \mathcal{B}_{\text{ox}}[\mathbf{0}_n, x^+] \cap \mathcal{B}_2(\mathbf{0}_n, \sqrt{\frac{1}{\mathcal{N}^-}})$. In this case, the x - and y -updates become²⁷:

$$\begin{aligned} x^{(k+1)} &= \arg \min_x \left\{ \frac{1}{2} x^\top (\Sigma + \varphi I_n) x - \varphi x^\top (y^{(k)} - u^{(k)}) + \mathbf{1}_{\Omega_3}(x) \right\} \\ &= (\Sigma + \varphi I_n)^{-1} \left(\varphi (y^{(k)} - u^{(k)}) + \frac{1 - \mathbf{1}_n^\top (\Sigma + \varphi I_n)^{-1} \varphi (y^{(k)} - u^{(k)})}{\mathbf{1}_n^\top (\Sigma + \varphi I_n)^{-1} \mathbf{1}_n} \mathbf{1}_n \right) \end{aligned}$$

and:

$$y^{(k+1)} = \mathcal{P}_{\mathcal{B}_{\text{ox}} - \mathcal{B}_{\text{all}}} \left(x^{(k+1)} + u^{(k)}; \mathbf{0}_n, x^+, \mathbf{0}_n, \sqrt{\frac{1}{\mathcal{N}^-}} \right)$$

where $\mathcal{P}_{\mathcal{B}_{\text{ox}} - \mathcal{B}_{\text{all}}}$ corresponds to the Dykstra's algorithm given in Appendix A.8.8 on page 59.

We consider the parameter set #1 given in Appendix B on page 64. The investment universe is made up of eight stocks. We would like to build a diversified minimum variance long-only portfolio without imposing an upper weight bound²⁸. In Table 4, we report the solutions found by the ADMM algorithm for several values of \mathcal{N}^- . When there is no Herfindahl constraint, the portfolio is fully invested in the seventh stock, meaning that the

²⁶We generally use the bisection algorithm to determine the optimal solution λ^* .

²⁷See Appendix A.4 on page 52 for the derivation of the x -update.

²⁸This means that x^+ is set to $\mathbf{1}_n$.

asset diversification is very poor. Then we increase the number of effective bets. If \mathcal{N}^- is equal to the number n of stocks, we verify that the solution corresponds to the equally-weighted portfolio. Between these two limit cases, we see the impact of the Herfindahl constraint on the portfolio diversification. The parameter set #1 is defined with respect to a capitalization-weighted index, whose weights are equal to 23%, 19%, 17%, 9%, 8%, 6% and 5%. The number of effective bets of this benchmark is equal to 6.435. If we impose that the effective number of bets of the minimum variance portfolio is at least equal to the effective number of bets of the benchmark, we find the following solution: 14.74%, 15.45%, 1.79%, 15.49%, 6.17%, 13.83%, 23.21% and 9.31%.

Table 4: Minimum variance portfolios (in %)

\mathcal{N}^-	1.00	2.00	3.00	4.00	5.00	6.00	6.50	7.00	7.50	8.00
x_1^*	0.00	3.22	9.60	13.83	15.18	15.05	14.69	14.27	13.75	12.50
x_2^*	0.00	12.75	14.14	15.85	16.19	15.89	15.39	14.82	14.13	12.50
x_3^*	0.00	0.00	0.00	0.00	0.00	0.07	2.05	4.21	6.79	12.50
x_4^*	0.00	10.13	15.01	17.38	17.21	16.09	15.40	14.72	13.97	12.50
x_5^*	0.00	0.00	0.00	0.00	0.71	5.10	6.33	7.64	9.17	12.50
x_6^*	0.00	5.36	8.95	12.42	13.68	14.01	13.80	13.56	13.25	12.50
x_7^*	100.00	68.53	52.31	40.01	31.52	25.13	22.92	20.63	18.00	12.50
x_8^*	0.00	0.00	0.00	0.50	5.51	8.66	9.41	10.14	10.95	12.50
λ^* (in %)	0.00	1.59	3.10	5.90	10.38	18.31	23.45	31.73	49.79	∞

As explained before, we can also solve the optimization problem by combining Problem (33) and the bisection algorithm. This is why we have reported the corresponding value λ^* in the last row in Table 4. However, this approach is no longer valid if we consider diversification constraints that are not quadratic. For instance, let us consider the generalized minimum variance problem:

$$\begin{aligned}
 x^* &= \arg \min_x \frac{1}{2} x^\top \Sigma x \\
 \text{s.t.} & \begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq x^+ \\ \mathcal{D}(x) \geq \mathcal{D}^- \end{cases}
 \end{aligned} \tag{34}$$

where $\mathcal{D}(x)$ is a weight diversification measure and \mathcal{D}^- is the minimum acceptable diversification. For example, we can use Shannon's entropy, the Gini index or the diversification ratio. In this case, it is not possible to obtain an equivalent QP problem, whereas the ADMM algorithm is exactly the same as previously except for the y -update:

$$y^{(k+1)} = \mathcal{P}_{\mathcal{B}_{ox}[\mathbf{0}_n, x^+] \cap \mathfrak{D}} \left(x^{(k+1)} + u^{(k)} \right)$$

where $\mathfrak{D} = \{x \in \mathbb{R}^n : \mathcal{D}(x) \geq \mathcal{D}^-\}$. The projection onto \mathfrak{D} can be easily derived from the proximal operator of the dual function (see the *tips and tricks* on page 42).

Remark 7 *If we compare the computational times, we observe that the best method is the second version of the ADMM algorithm. In our example, the computational time is divided by a factor of eight with respect to the bisection approach²⁹. If we consider a large-scale problem with n larger than 1000, the computational time is generally divided by a factor greater than 50!*

²⁹In contrast, the first version of the ADMM algorithm is not efficient since the computational time is multiply by a factor of five with respect to the bisection approach.

4.1.2 Managing the portfolio rebalancing process

Another big challenge of the minimum variance portfolio is the management of the turnover between two rebalancing dates. Let x_t be the current portfolio. The optimization program for calibrating the optimal solution x_{t+1} for the next rebalancing date $t + 1$ may include a penalty function $\mathbf{c}(x | x_t)$ and/or a weight constraint $\mathfrak{C}(x | x_t)$ that are parameterized with respect to the current portfolio x_t :

$$\begin{aligned} x_{t+1} &= \arg \min_x \frac{1}{2} x^\top \Sigma x + \mathbf{c}(x | x_t) \\ \text{s.t. } &\begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq x^+ \\ x \in \mathfrak{C}(x | x_t) \end{cases} \end{aligned} \quad (35)$$

Again, we can solve this problem using the ADMM algorithm. Thanks to the Dykstra's algorithm, the only difficulty is finding the proximal operator of $\mathbf{c}(x | x_t)$ or $\mathfrak{C}(x | x_t)$ when performing the y -update.

Let us define the cost function as:

$$\mathbf{c}(x | x_t) = \lambda \sum_{i=1}^n \left(c_i^- (x_{i,t} - x_i)_+ + c_i^+ (x_i - x_{i,t})_+ \right)$$

where c_i^- and c_i^+ are the bid and ask transaction costs. In Appendix A.8.12 on page 62, we show that the proximal operator is:

$$\text{prox}_{\mathbf{c}(x|x_t)}(v) = x_t + \mathcal{S}(v - x_t; \lambda c^-, \lambda c^+) \quad (36)$$

where $\mathcal{S}(v; \lambda_-, \lambda_+) = (v - \lambda_+)_+ - (v + \lambda_-)_-$ is the two-sided soft-thresholding operator.

If we define the cost constraint $\mathfrak{C}(x | x_t)$ as a turnover constraint:

$$\mathfrak{C}(x | x_t) = \{x \in \mathbb{R}^n : \|x - x_t\|_1 \leq \tau^+\}$$

the proximal operator is:

$$\mathcal{P}_{\mathfrak{C}}(v) = v - \text{sign}(v - x_t) \odot \min(|v - x_t|, s^*) \quad (37)$$

where $s^* = \left\{ s \in \mathbb{R} : \sum_{i=1}^n (|v_i - x_{t,i}| - s)_+ = \tau^+ \right\}$.

Remark 8 *These two examples are very basic and show how we can easily introduce turnover management using the ADMM framework. More sophisticated approaches are presented in Section 4.4 on page 42.*

4.2 Smart beta portfolios

In this section, we consider three main models of smart beta portfolios: the equal risk contribution (ERC) portfolio, the risk budgeting (RB) portfolio and the most diversified portfolio (MDP). Specific algorithms for these portfolios based on the CCD method have already been presented in Griveau-Billion *et al.* (2013) and Richard and Roncalli (2015, 2019). We extend these results to the ADMM algorithm.

4.2.1 The ERC portfolio

The ERC portfolio uses the volatility risk measure $\sigma(x) = \sqrt{x^\top \Sigma x}$ and allocates the weights such that the risk contribution is the same for all the assets of the portfolio (Maillard *et al.*, 2010):

$$\mathcal{RC}_i(x) = x_i \frac{\partial \sigma(x)}{\partial x_i} = x_j \frac{\partial \sigma(x)}{\partial x_j} = \mathcal{RC}_j(x)$$

In this case, we can show that the ERC portfolio is the scaled solution $x^*/(\mathbf{1}_n^\top x^*)$ where x^* is given by:

$$x^* = \arg \min_x \frac{1}{2} x^\top \Sigma x - \lambda \sum_{i=1}^n \ln x_i$$

and λ is any positive scalar. The first-order condition is $(\Sigma x)_i - \lambda x_i^{-1} = 0$. It follows that $x_i (\Sigma x)_i - \lambda = 0$ or:

$$x_i^2 \sigma_i^2 + x_i \sigma_i \sum_{j \neq i} x_j \rho_{i,j} \sigma_j - \lambda = 0$$

We deduce that the solution is the positive root of the second-degree equation. Finally, we obtain the following iteration for the CCD algorithm:

$$x_i^{(k+1)} = \frac{-v_i^{(k+1)} + \sqrt{\left(v_i^{(k+1)}\right)^2 + 4\lambda\sigma_i^2}}{2\sigma_i^2} \quad (38)$$

where:

$$v_i^{(k+1)} = \sigma_i \sum_{j < i} x_j^{(k+1)} \rho_{i,j} \sigma_j + \sigma_i \sum_{j > i} x_j^{(k)} \rho_{i,j} \sigma_j$$

The ADMM algorithm uses the first trick where $f_x(x) = \frac{1}{2} x^\top \Sigma x$ and $f_y(y) = -\lambda \sum_{i=1}^n \ln y_i$. It follows that the x - and y -update steps are:

$$x^{(k+1)} = (\Sigma + \varphi I_n)^{-1} \varphi \left(y^{(k)} - u^{(k)} \right)$$

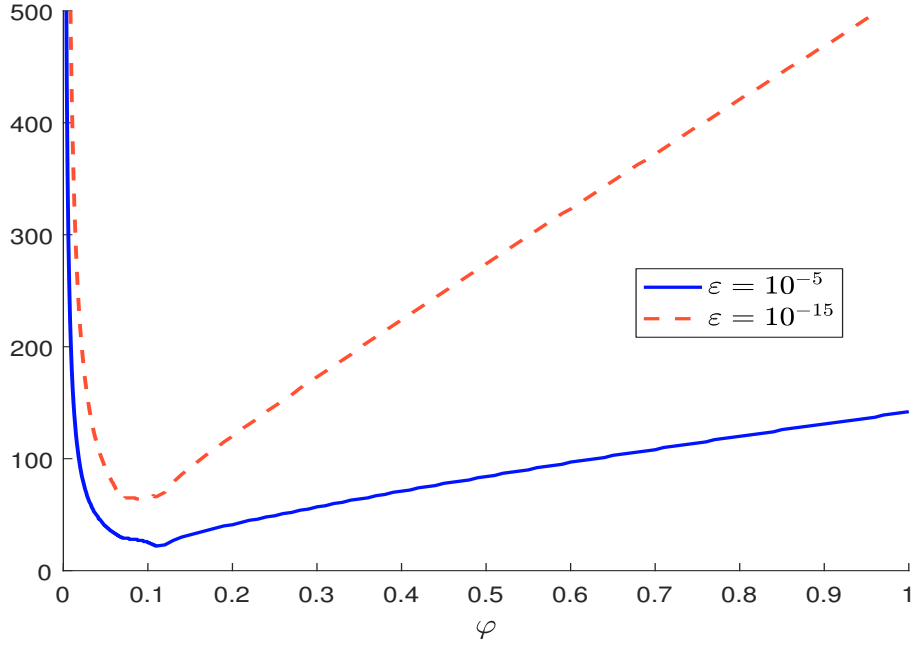
and:

$$y_i^{(k+1)} = \frac{1}{2} \left(\left(x_i^{(k+1)} + u_i^{(k)} \right) + \sqrt{\left(x_i^{(k+1)} + u_i^{(k)} \right)^2 + 4\lambda\varphi^{-1}} \right)$$

We apply the CCD and ADMM algorithms to the parameter set #1. We find that the ERC portfolio is equal to 11.40%, 12.29%, 5.49%, 11.91%, 6.65%, 10.81%, 33.52% and 7.93%. It appears that the CCD algorithm is much more efficient than the ADMM algorithm. For instance, if we set $\lambda = \sqrt{x^{(0)\top \Sigma x^{(0)}}$, $x^{(0)} = n^{-1} \mathbf{1}_n$ and $\varphi = 1$, the CCD algorithm needs six cycles to converge whereas the ADMM algorithm needs 156 iterations if we set the convergence criterion³⁰ $\varepsilon = 10^{-8}$. Whatever the values of λ , $x^{(0)}$ and ε , our experience is that the CCD generally converges within less than 15 cycles even if the number of assets is greater than 1000. The convergence of the ADMM is more of an issue, because it depends on the parameters λ and φ . In Figure 5, we have reported the number of iterations of the ADMM with respect to φ for several values of ε when $\lambda = 1$ and $x^{(0)} = \mathbf{1}_n$. We verify that it is very sensitive to the value taken by φ . Curiously, the parameter λ has little influence, meaning that the convergence issue mainly concerns the x -update step.

³⁰The termination rule is defined as $\max_i |x_i^{(k+1)} - x_i^{(k)}| \leq \varepsilon$.

Figure 5: Number of ADMM iterations for finding the ERC portfolio



4.2.2 Risk budgeting optimization

The ERC portfolio has been extended by Roncalli (2013) when the risk budgets are not equal and when the risk measure $\mathcal{R}(x)$ is convex and coherent:

$$\mathcal{RC}_i(x) = x_i \frac{\partial \mathcal{R}(x)}{\partial x_i} = \mathcal{RB}_i$$

where \mathcal{RB}_i is the risk budget allocated to Asset i . In this case, we can show that the risk budgeting portfolio is the scaled solution of the following optimization problem:

$$x^* = \arg \min_x \mathcal{R}(x) - \lambda \sum_{i=1}^n \mathcal{RB}_i \cdot \ln x_i$$

where λ is any positive scalar. Depending on the risk measure, we can use the CCD or the ADMM algorithm.

For example, Roncalli (2015) proposes using the standard deviation-based risk measure:

$$\mathcal{R}(x) = -x^\top (\mu - r) + \xi \sqrt{x^\top \Sigma x}$$

In this case, the first-order condition for defining the CCD algorithm is:

$$-(\mu_i - r) + \xi \frac{(\Sigma x)_i}{\sqrt{x^\top \Sigma x}} - \lambda \frac{\mathcal{RB}_i}{x_i} = 0$$

It follows that $\xi x_i (\Sigma x)_i - (\mu_i - r) x_i \sigma(x) - \lambda \sigma(x) \cdot \mathcal{RB}_i = 0$ or equivalently:

$$\alpha_i x_i^2 + \beta_i x_i + \gamma_i = 0$$

where $\alpha_i = \xi \sigma_i^2$, $\beta_i = \xi \sigma_i \sum_{j \neq i} x_j \rho_{i,j} \sigma_j - (\mu_i - r) \sigma(x)$ and $\gamma_i = -\lambda \sigma(x) \cdot \mathcal{RB}_i$. We notice that the solution x_i depends on the volatility $\sigma(x)$. Here, we face an endogenous problem, because $\sigma(x)$ depends on x_i . Griveau-Billon *et al.* (2015) notice that this is not an issue, because we may assume that $\sigma(x)$ is almost constant between two coordinate iterations of the CCD algorithm. They deduce that the coordinate solution is then the positive root of the second-degree equation:

$$x_i^{(k+1)} = \frac{-\beta_i^{(k+1)} + \sqrt{\left(\beta_i^{(k+1)}\right)^2 - 4\alpha_i^{(k+1)}\gamma_i^{(k+1)}}}{2\alpha_i^{(k+1)}} \quad (39)$$

where:

$$\begin{cases} \alpha_i^{(k+1)} = \xi \sigma_i^2 \\ \beta_i^{(k+1)} = \xi \sigma_i \left(\sum_{j < i} x_j^{(k+1)} \rho_{i,j} \sigma_j + \sum_{j > i} x_j^{(k)} \rho_{i,j} \sigma_j \right) - (\mu_i - r) \sigma_i^{(k+1)}(x) \\ \gamma_i^{(k+1)} = -\lambda \sigma_i^{(k+1)}(x) \cdot \mathcal{RB}_i \\ \sigma_i^{(k+1)}(x) = \sqrt{\chi^\top \Sigma \chi} \\ \chi = \left(x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i^{(k)}, x_{i+1}^{(k)}, \dots, x_n^{(k)} \right) \end{cases}$$

In the case of the volatility or the standard deviation-based risk measure, we apply the exact formulation of the CCD algorithm because we have an analytical solution of the first-order condition. This is not always the case, especially when we consider skewness-based risk measure (Bruder *et al.*, 2016; Lezmi *et al.*, 2018). In this case, we can use the gradient formulation of the CCD algorithm or the ADMM algorithm, which is defined as follows:

$$\begin{cases} x^{(k+1)} = \mathbf{prox}_{\varphi^{-1}\mathcal{R}(x)} \left(y^{(k)} - u^{(k)} \right) \\ v_y^{(k+1)} = x^{(k+1)} + u^{(k)} \\ y^{(k+1)} = \frac{1}{2} \left(v_y^{(k+1)} + \sqrt{v_y^{(k+1)} \odot v_y^{(k+1)} + 4\lambda \varphi^{-1} \cdot \mathcal{RB}} \right) \\ u^{(k+1)} = u^{(k)} + x^{(k+1)} - y^{(k+1)} \end{cases}$$

4.2.3 The most diversified portfolio

Choueifaty and Coignard (2008) introduce the concept of diversification ratio, which corresponds to the following expression:

$$\mathcal{DR}(x) = \frac{\sum_{i=1}^n x_i \sigma_i}{\sigma(x)} = \frac{x^\top \sigma}{\sqrt{x^\top \Sigma x}}$$

By construction, the diversification ratio of a portfolio fully invested in one asset is equal to one, whereas it is larger than one in the general case. The authors then propose building the most diversified portfolio as the portfolio which maximizes the diversification ratio. It is also the solution to the following minimization problem³¹:

$$\begin{aligned} x^\star &= \arg \min_x \frac{1}{2} \ln(x^\top \Sigma x) - \ln(x^\top \sigma) \\ \text{s.t.} \quad &\begin{cases} \mathbf{1}_n^\top x = 1 \\ x \in \Omega \end{cases} \end{aligned}$$

³¹See Choueifaty *et al.* (2013).

This problem is relatively easy to solve using standard numerical algorithms if Ω corresponds to linear constraints, for example weight constraints. However, the optimal solution may face the same problem as the minimum variance portfolio since most of the time it is concentrated on a small number of assets. This is why it is interesting to add a weight diversification constraint $\mathcal{D}(x) \geq \mathcal{D}^-$. For example, we can assume that the number of effective bets $\mathcal{N}(x)$ is larger than a minimum acceptable value \mathcal{N}^- . Contrary to the minimum variance portfolio, we do not obtain a QP problem and we observe that the optimization problem is tricky in practice. Thanks to the ADMM algorithm, we can however simplify the optimization problem by splitting the constraints and using the same approach that has been already described on page 33. The x -update consists in finding the regularized standard MDP:

$$x^{(k+1)} = \arg \min_x \left\{ \frac{1}{2} \ln(x^\top \Sigma x) - \ln(x^\top \sigma) + \frac{\varphi}{2} \|x - y^{(k)} + u^{(k)}\|_2^2 \quad \text{s.t.} \quad \mathbf{1}_n^\top x = 1 \right\}$$

whereas the y -update corresponds to the projection onto the intersection of Ω and $\mathfrak{D} = \{x \in \mathbb{R}^n : \mathcal{D}(x) \geq \mathcal{D}^-\}$:

$$y^{(k+1)} = \mathcal{P}_{\Omega \cap \mathfrak{D}}(x^{(k+1)} + u^{(k)})$$

We consider the parameter set #2 given in Appendix B on page 64. Results are reported in Table 5. The second column corresponds to the long/short MDP portfolio (or $\Omega = \mathbb{R}^n$). By definition, we cannot compute the number of effective bets because it contains short positions. The other columns correspond to the long-only MDP portfolio (or $\Omega = [0, 1]^n$) when we impose a sufficient number of effective bets \mathcal{N}^- . We notice that the traditional long-only MDP is poorly diversified in terms of weights since we have $\mathcal{N}(x) = 2.30$. As for the minimum variance portfolio, the MDP tends to the equally-weighted portfolio when \mathcal{N}^- tends to the number of assets

Table 5: MDP portfolios (in %)

\mathcal{N}^-	L/S	Long-only					
		0.00	3.00	4.00	5.00	6.00	7.00
x_1^*	41.81	41.04	35.74	30.29	26.08	22.44	18.83
x_2^*	51.88	50.92	43.91	36.68	31.05	26.12	21.19
x_3^*	8.20	8.05	10.12	11.52	12.33	12.80	13.01
x_4^*	-0.43	0.00	2.48	5.12	7.16	8.90	10.51
x_5^*	-0.26	0.00	0.92	2.28	3.60	5.02	6.85
x_6^*	-0.38	0.00	2.03	4.36	6.28	8.02	9.79
x_7^*	-0.51	0.00	3.47	6.68	8.85	10.44	11.65
x_8^*	-0.31	0.00	1.32	3.07	4.65	6.27	8.17
$\mathcal{N}(x)$		2.30	3.00	4.00	5.00	6.00	7.00

4.3 Robo-advisory optimization

Today's financial industry is facing a digital revolution in all areas: payment services, on-line banking, asset management, etc. This is particularly true for the financial advisory industry, which has been impacted in the last few years by the emergence of digitalization and robo-advisors. The demand for robo-advisors is strong, which explains the growth of

this business³². How does one characterize a robo-advisor? This is not simple, but the underlying idea is to build a systematic portfolio allocation in order to provide a customized advisory service. A robo-advisor has two main objectives. The first objective is to know the investor better than a traditional asset manager. Because of this better knowledge, the robo-advisor may propose a more appropriate asset allocation. The second objective is to perform the task in a systematic way and to build an automated rebalancing process. Ultimately, the goal is to offer a customized solution. In fact, the reality is very different. We generally notice that many robo-advisors are more a web or a digital application, but not really a robo-advisor. The reason is that portfolio optimization is a very difficult task. In many robo-advisors, asset allocation is then rather human-based or not completely systematic with the aim to rectify the shortcomings of mean-variance optimization. Over the next five years, the most important challenge for robo-advisors will be to reduce these discretionary decisions and improve the robustness of their systematic asset allocation process. But this means that robo-advisors must give up the quadratic programming world of the portfolio allocation.

4.3.1 Specification of the objective function

In order to make mean-variance optimization more robust, two directions can be followed. The first one has been largely explored and adds a penalty function in order to regularize or sparsify the solution (Brodie *et al.* 2009; DeMiguel *et al.*, 2009; Carrasco and Noumon, 2010; Bruder *et al.*, 2013; Bourgeron *et al.*, 2018). The second one consists in changing the objective function and considering risk budgeting portfolios instead of mean-variance optimized portfolios (Maillard *et al.*, 2010; Roncalli, 2013). Even if this second direction has encountered great success, it presents a solution that is not sufficiently flexible in terms of active management. Nevertheless, these two directions are not so divergent. Indeed, Roncalli (2013) shows that the risk budgeting optimization can be viewed as a non-linear shrinkage approach of the minimum variance optimization. Richard and Roncalli (2015) propose then a unified approach of smart beta portfolios by considering alternative allocation models as penalty functions of the minimum variance optimization. In particular, they use the logarithmic barrier function in order to regularize minimum variance portfolios. This idea has also been reiterated by de Jong (2018), who considers a mean-variance framework.

Therefore, we propose defining the robo-advisor optimization problem as follows:

$$\begin{aligned} x_{t+1}^* &= \arg \min_x f_{\mathcal{R}obo}(x) \\ \text{s.t.} \quad &\begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq \mathbf{1}_n \\ x \in \Omega \end{cases} \end{aligned} \quad (40)$$

where:

$$\begin{aligned} f_{\mathcal{R}obo}(x) &= \frac{1}{2} (x - b)^\top \Sigma_t (x - b) - \gamma (x - b)^\top \mu_t + \\ &\quad \varrho_1 \|\Gamma_1 (x - x_t)\|_1 + \frac{1}{2} \varrho_2 \|\Gamma_2 (x - x_t)\|_2^2 + \\ &\quad \tilde{\varrho}_1 \|\tilde{\Gamma}_1 (x - \tilde{x})\|_1 + \frac{1}{2} \tilde{\varrho}_2 \|\tilde{\Gamma}_2 (x - \tilde{x})\|_2^2 - \lambda \sum_{i=1}^n \mathcal{RB}_i \cdot \ln x_i \end{aligned} \quad (41)$$

³²For instance, the growth was 60% per year in the US over the last five years. In Europe, the growth is also impressive, even though the market is smaller. In the last two years, assets under management have increased 14-fold.

b is the benchmark portfolio, \tilde{x} is the reference portfolio and x_t is the current portfolio.

This specification is sufficiently broad that it encompasses most models used by the industry. We notice that the objective function is made up of three parts. The first part corresponds to the MVO objective function with a benchmark. If we set b equal to $\mathbf{0}_n$, we obtain the Markowitz utility function. The second part contains ℓ_1 - and ℓ_2 -norm penalty functions. The regularization can be done with respect to the current allocation x_t in order to control the rebalancing process and the smoothness of the dynamic allocation. The regularization can also be done with respect to a reference portfolio, which is generally the strategic asset allocation of the fund. The idea is to control the profile of the fund. For example, if the strategic asset allocation is an 80/20 asset mix policy, we do not want the portfolio to present a defensive or balanced risk profile. Finally, the third part of the objective function corresponds to the logarithmic barrier function, where the parameter λ controls the trade-off between MVO optimization and RB optimization. This last part is very important in order to make the dynamic asset allocation more robust. The hyperparameters of the objective function are ϱ_1 , ϱ_2 , $\tilde{\varrho}_1$, $\tilde{\varrho}_2$ and λ . They are all positive and can also be set to zero in order to deactivate a penalty function. For instance, ϱ_2 and $\tilde{\varrho}_2$ are equal to zero if we don't want to have a shrinkage of the covariance matrix Σ_t . The hyperparameters ϱ_1 and $\tilde{\varrho}_1$ can also be equal to zero because the ℓ_1 regularization is generally introduced when specifying the additional constraints Ω . The parameter γ is not really a hyperparameter, because it is generally calibrated to target volatility or an expected return. We also notice that this model encompasses the Black-Litterman model thanks to the specification of μ_t (Bourgeron *et al.*, 2018). Another important component of this framework is the specification of the set $x \in \Omega$. It may include traditional constraints such as weight bounds and/or asset class limits, but we can also add non-linear constraints such as a turnover limit, an active share floor or a weight diversification constraint.

4.3.2 Derivation of the general algorithm

Problem (40) is equivalent to solving:

$$x_{t+1}^* = \arg \min_x f_{\mathcal{R}obo}^+(x)$$

where the objective function can be broken down as follows:

$$\begin{aligned} f_{\mathcal{R}obo}^+(x) &= f_{\text{MVO}}(x) + f_{\ell_1}(x) + f_{\ell_2}(x) + f_{\text{RB}}(x) + \\ &\quad + \mathbf{1}_{\Omega_0}(x) + \mathbf{1}_{\Omega}(x) \end{aligned}$$

where:

$$\begin{aligned} f_{\text{MVO}}(x) &= \frac{1}{2} (x - b)^\top \Sigma_t (x - b) - \gamma (x - b)^\top \mu_t \\ f_{\ell_1}(x) &= \varrho_1 \|\Gamma_1(x - x_t)\|_1 + \tilde{\varrho}_1 \|\tilde{\Gamma}_1(x - \tilde{x})\|_1 \\ f_{\ell_2}(x) &= \frac{1}{2} \varrho_2 \|\Gamma_2(x - x_t)\|_2^2 + \frac{1}{2} \tilde{\varrho}_2 \|\tilde{\Gamma}_2(x - \tilde{x})\|_2^2 \\ f_{\text{RB}}(x) &= -\lambda \sum_{i=1}^n \mathcal{R}\mathcal{B}_i \cdot \ln x_i \end{aligned}$$

and $\Omega_0 = \{x \in [0, 1]^n : \mathbf{1}_n^\top x = 1\}$. The ADMM algorithm is implemented as follows:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min f_x(x) + f_y(y) \\ \text{s.t. } &x - y = \mathbf{0}_n \end{aligned}$$

This is the general approach for solving the robo-advisor problem.

The main task is then to split the function $f_{\mathcal{R}obo}^+$ into f_x and f_y . However, in order to be efficient, the x - and y -update steps of the ADMM algorithm must be easy to compute. Therefore, we impose that the x -step is solved using QP or CCD methods while the y -step is solved using the Dykstra's algorithm, where each component corresponds to an analytical proximal operator. Moreover, we also split the set of constraints Ω into a set of linear constraints $\Omega_{\mathcal{L}inear}$ and a set of non-linear constraints $\Omega_{\mathcal{N}onlinear}$. This lead defining $f_x(x)$ and $f_y(y)$ as follows:

$$\begin{cases} f_x(x) = f_{\text{MVO}}(x) + f_{\ell_2}(x) + \mathbf{1}_{\Omega_0}(x) + \mathbf{1}_{\Omega_{\mathcal{L}inear}}(x) \\ f_y(y) = f_{\ell_1}(y) + f_{\text{RB}}(x) + \mathbf{1}_{\Omega_{\mathcal{N}onlinear}}(x) \end{cases} \quad (42)$$

We notice that $f_x(x)$ has a quadratic form, implying that the x -step may be solved using a QP algorithm. Another formulation is:

$$\begin{cases} f_x(x) = f_{\text{MVO}}(x) + f_{\ell_2}(x) + f_{\text{RB}}(x) \\ f_y(y) = f_{\ell_1}(y) + \mathbf{1}_{\Omega_0}(x) + \mathbf{1}_{\Omega_{\mathcal{L}inear}}(x) + \mathbf{1}_{\Omega_{\mathcal{N}onlinear}}(x) \end{cases} \quad (43)$$

In this case, the x -step is solved using the CCD algorithm.

4.3.3 Specific algorithms

The ADMM-QP formulation If we consider Formulation (42), we have:

$$\begin{aligned} f_{\text{QP}}(x) &= f_{\text{MVO}}(x) + f_{\ell_2}(x) \\ &= \frac{1}{2}(x-b)^\top \Sigma_t (x-b) - \gamma(x-b)^\top \mu_t + \frac{1}{2}\varrho_2 \|\Gamma_2(x-x_t)\|_2^2 + \frac{1}{2}\tilde{\varrho}_2 \|\tilde{\Gamma}_2(x-\tilde{x})\|_2^2 \\ &= \frac{1}{2}x^\top Qx - x^\top R + C \end{aligned}$$

where $Q = \Sigma_t + \varrho_2 \Gamma_2^\top \Gamma_2 + \tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2$, $R = \gamma \mu_t + \Sigma_t b + \varrho_2 \Gamma_2^\top \Gamma_2 x_t + \tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2 \tilde{x}$ and C is a constant³³. Using the fourth ADMM trick, we deduce that $x^{(k+1)}$ is the solution of the following QP problem:

$$\begin{aligned} x^{(k+1)} &= \arg \min_x \frac{1}{2}x^\top (Q + \varphi I_n) x - x^\top (R + \varphi (y^{(k)} - u^{(k)})) \\ \text{s.t.} \quad &\begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq \mathbf{1}_n \end{cases} \end{aligned}$$

Since the proximal operators of f_{ℓ_1} and f_{RB} have been already computed, finding $y^{(k+1)}$ is straightforward with the Dykstra's algorithm as long as the proximal of each non-linear constraint is known.

The ADMM-CCD formulation If we consider Formulation (43), we have:

$$f_x(x) = f_{\text{QP}}(x) - \lambda \sum_{i=1}^n \mathcal{R}B_i \cdot \ln x_i$$

³³The expression of $f_{\text{QP}}(x)$ is computed in Appendix A.9 on page 63.

Using Appendix A.10 on page 63, the CCD algorithm applied to x -update is:

$$x_i^{(k+1)} = \frac{R_i - \sum_{j < i} x_j^{(k+1)} Q_{i,j} - \sum_{j > i} x_j^{(k)} Q_{i,j}}{2Q_{i,i}} + \frac{\sqrt{\left(\sum_{j < i} x_j^{(k+1)} Q_{i,j} + \sum_{j > i} x_j^{(k)} Q_{i,j} - R_i\right)^2 + 4\lambda_i Q_{i,i}}}{2Q_{i,i}}$$

where the matrices Q and R are defined as:

$$Q = \Sigma_t + \varrho_2 \Gamma_2^\top \Gamma_2 + \tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2 + \varphi I_n$$

and:

$$R = \gamma \mu_t + \Sigma_t b + \varrho_2 \Gamma_2^\top \Gamma_2 x_t + \tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2 \tilde{x} + \varphi \left(y^{(k)} - u^{(k)} \right)$$

and $\lambda_i = \lambda \cdot \mathcal{RB}_i$. Like the ADMM-QP formulation, the y -update step does not pose any particular difficulties.

4.4 Tips and tricks

If we consider the different portfolio optimization approaches presented in Table 1, we have shown how to solve MVO (1), GMV (2), MDP (3), ERC (4) and RB (5) models. The RQE (7) model is equivalent to the GMV (2) model by replacing the covariance matrix Σ by the dissimilarity matrix D . Below, we implement the Kullback-Leibler model (4) of Bera and Park (2008) using the ADMM framework. Concerning the regularization problems in Table 2, ridge (8), lasso (9) and log-barrier (10) penalty functions have been already covered. Indeed, ridge and lasso penalizations correspond to the proximal operator of ℓ_1 - and ℓ_2 -norm functions by applying the translation $g(x) = x - \tilde{x}$. Shannon's entropy (11) penalization is discussed below. For the constraints that are considered in Table 3, imposing no cash and leverage (12) is done with the proximal of the hyperplane $\mathcal{H}_{hyperlane}[\mathbf{1}_n, 1]$. No short selling (13) and weight bounds (14) are equivalent to considering the box projections $\mathcal{B}_{ox}[\mathbf{0}_n, \infty]$ and $\mathcal{B}_{ox}[x^-, x^+]$. Asset class limits can be implemented using the projection onto the intersection of two half-spaces $\mathcal{H}_{halfspace}[\mathbf{1}_{i \in C_j}, c_j^+]$ and $\mathcal{H}_{halfspace}[-\mathbf{1}_{i \in C_j}, -c_j^-]$. The proximal of the turnover (16) had been already given in Equation (37) on page 34. If we want to impose an upper limit on transaction costs (17), we use the Moreau decomposition and Equation (36). Finally, Section 4.1.1 on page 31 dealt with the weight diversification problem of the number of active bets. Therefore, it remains to solve leverage limits (18), long/short exposure (19) restrictions and active management constraints: benchmarking (20), tracking error floor (21) and active share floor (22).

4.4.1 Volatility and return targeting

We first consider the μ -problem and the σ -problem. Targeting a minimum expected return $\mu(x) \geq \mu^*$ can be implemented in the ADMM framework using the proximal operator of the hyperplane³⁴ $\mathcal{H}_{hyperlane}[-\mu, -\mu^*]$. In the case of the σ -problem $\sigma(x) \leq \sigma^*$, we use the fourth ADMM trick. Let L be the lower Cholesky decomposition of Σ , we have:

$$\begin{aligned} \sigma(x) \leq \sigma^* &\Leftrightarrow \sqrt{x^\top \Sigma x} \leq \sigma^* \\ &\Leftrightarrow \sqrt{x^\top (LL^\top) x} \leq \sigma^* \\ &\Leftrightarrow \|y^\top y\|_2 \leq \sigma^* \end{aligned}$$

³⁴We have $\mu(x) \geq \mu^* \Leftrightarrow x^\top \mu \geq \mu^* \Leftrightarrow -\mu^\top x \leq -\mu^*$.

where $y = L^\top x$. It follows that the proximal of the y -update is the projection onto the ℓ_2 ball $\mathcal{B}_2(\mathbf{0}_n, \sigma^*)$.

4.4.2 Leverage management

If we impose a leverage limit $\sum_{i=1}^n |x_i| \leq \mathcal{L}^+$, we have $\|x\|_1 \leq \mathcal{L}^+$ and the proximal of the y -update is the projection onto the ℓ_1 ball $\mathcal{B}_1(\mathbf{0}_n, \mathcal{L}^+)$. If the leverage constraint concerns the long/short limits $-\mathcal{L}\mathcal{S}^- \leq \sum_{i=1}^n x_i \leq \mathcal{L}\mathcal{S}^+$, we consider the intersection of the two half-spaces $\mathcal{H}_{\text{alfspace}}[\mathbf{1}_n, \mathcal{L}\mathcal{S}^+]$ and $\mathcal{H}_{\text{alfspace}}[-\mathbf{1}_n, \mathcal{L}\mathcal{S}^-]$. If we consider an absolute leverage $|\sum_{i=1}^n x_i| \leq \mathcal{L}^+$, we obtain the previous case with $\mathcal{L}\mathcal{S}^- = \mathcal{L}\mathcal{S}^+ = \mathcal{L}^+$. Portfolio managers can also use another constraint concerning the sum of the k largest values³⁵:

$$f(x) = \sum_{i=n-k+1}^n x_{(i:n)} = x_{(n:n)} + \dots + x_{(n-k+1:n)}$$

where $x_{(i:n)}$ is the order statistics of x : $x_{(1:n)} \leq x_{(2:n)} \leq \dots \leq x_{(n:n)}$. Beck (2017) shows that:

$$\text{prox}_{\lambda f(x)}(v) = v - \lambda \mathcal{P}_\Omega\left(\frac{v}{\lambda}\right)$$

where:

$$\Omega = \{x \in [0, 1]^n : \mathbf{1}_n^\top x = k\} = \mathcal{B}_{\text{ox}}[\mathbf{0}_n, \mathbf{1}_n] \cap \mathcal{H}_{\text{hyperlane}}[\mathbf{1}_n, k]$$

4.4.3 Entropy portfolio and diversification measure

Bera and Park (2008) propose using a cross-entropy measure as the objective function:

$$\begin{aligned} x^\star &= \arg \min_x \text{KL}(x \mid \tilde{x}) \\ \text{s.t.} &\begin{cases} \mathbf{1}_n^\top x = 1 \\ \mathbf{0}_n \leq x \leq \mathbf{1}_n \\ \mu(x) \geq \mu^\star \\ \sigma(x) \leq \sigma^\star \end{cases} \end{aligned}$$

where $\text{KL}(x \mid \tilde{x}) = \sum_{i=1}^n x_i \ln(x_i/\tilde{x}_i)$ and \tilde{x} is a reference portfolio, which is well-diversified (e.g. the EW³⁶ or ERC portfolio). In Appendix A.8.9 on page 60, we show that the proximal operator of $\lambda \text{KL}(x \mid \tilde{x})$ is equal to:

$$\text{prox}_{\lambda \text{KL}(x \mid \tilde{x})}(v) = \lambda \begin{pmatrix} W\left(\lambda^{-1} \tilde{x}_1 e^{\lambda^{-1} v_1 - \tilde{x}_1^{-1}}\right) \\ \vdots \\ W\left(\lambda^{-1} \tilde{x}_n e^{\lambda^{-1} v_n - \tilde{x}_n^{-1}}\right) \end{pmatrix}$$

where $W(x)$ is the Lambert W function.

Remark 9 Using the previous result and the fact that $\text{SE}(x) = -\text{KL}(x \mid \mathbf{1}_n)$, we can use Shannon's entropy to define the diversification measure $\mathcal{D}(x) = \text{SE}(x)$. Therefore, solving Problem (34) is straightforward when we consider the following diversification set:

$$\mathcal{D} = \left\{x \in [0, 1]^n : -\sum_{i=1}^n x_i \ln x_i \geq \text{SE}^-\right\}$$

³⁵An example is the 5/10/40 UCITS rule: A UCITS fund may invest no more than 10% of its net assets in transferable securities or money market instruments issued by the same body, with a further aggregate limitation of 40% of net assets on exposures of greater than 5% to single issuers.

³⁶In this case, it is equivalent to maximize Shannon's entropy because $\tilde{x} = \mathbf{1}_n$.

4.4.4 Passive and active management

In the case of the active share, we use the translation property:

$$\begin{aligned}\mathcal{AS}(x \mid \tilde{x}) &= \frac{1}{2} \sum_{i=1}^n |x_i - \tilde{x}_i| \\ &= \frac{1}{2} \|x - \tilde{x}\|_1\end{aligned}$$

The proximal operator is given in Appendix A.8.11 on page 62. It is interesting to notice that this type of problem cannot be solved using an augmented QP algorithm since it involves the complement of the ℓ_1 ball and not directly the ℓ_1 ball itself. In this case, we face a maximization problem and not a minimization problem, and the technique of augmented variables does not work.

For tracking error volatility, again we use the fourth ADMM trick:

$$\begin{aligned}\sigma(x \mid \tilde{x}) &= \sqrt{(x - \tilde{x})^\top \Sigma (x - \tilde{x})} \\ &= \|y\|_2\end{aligned}$$

where $y = L^\top x - L^\top \tilde{x}$. Using our ADMM notations, we have $Ax + By = c$ where $A = L^\top$, $B = -I_n$ and $c = L^\top \tilde{x}$.

4.4.5 Index sampling

Index sampling is based on the cardinality constraint $\sum_{i=1}^n \mathbb{1}\{x_i > 0\} \leq n_x$. It is closed to the ℓ_0 -norm function $\|x\|_0 = \sum_{i=1}^n \mathbb{1}\{x_i \neq 0\}$. Beck (2017) derives the proximal of $\lambda \|x\|_0$ on pages 137-138 of his monograph. However, it does not help to solve the index sampling problem, because we are interested in computing the projection onto the ℓ_0 ball and not the proximal of the ℓ_0 -norm function³⁷. This is why index sampling remains an open problem using the ADMM framework.

5 Conclusion

The aim of this paper is to propose an alternative solution to the quadratic programming algorithm in the context of portfolio allocation. In numerical analysis, the quadratic programming model is a powerful optimization tool, which is computationally very efficient. In portfolio management, the mean-variance optimization model is exactly a quadratic programming model, meaning that it benefits from its computational power. Therefore, the success of the Markowitz allocation model is explained by these two factors: the quadratic utility function and the quadratic programming setup. A lot of academics and professionals have proposed an alternative approach to the MVO framework, but very few of these models are used in practice. The main reason is that these competing models focus on the objective function and not on the numerical implementation. However, we believe that any model which is not tractable will have little success with portfolio managers. The analogy is obvious if we consider the theory of options. The success of the Black-Scholes model lies in the Black-Scholes analytical formula. Over the last thirty years, many models have been created (e.g. local volatility and stochastic volatility models), but only one can really

³⁷We cannot use the Moreau decomposition, because the dual of $\lambda \|x\|_p$ is not necessarily the ball $\mathcal{B}_p(\mathbf{0}_n, \lambda)$. For example, the dual of the ℓ_2 -norm function is the ℓ_2 ball, but the dual of the ℓ_1 -norm function is the ℓ_∞ ball.

compete with the Black-Scholes model. This is the SABR model, and the main reason is that it has an analytical formula for implied volatility.

This paper focuses then on a general approach for numerically solving non-QP portfolio allocation models. For that, we consider some algorithms that have been successfully applied to machine learning and large-scale optimization. For instance, the coordinate descent algorithm is the fastest method for performing high-dimensional lasso regression, while the Dykstra's algorithm has been created to find the solution of restricted least squares regression. Since there is a strong link between MVO and linear regression (Scherer, 2007), this is not a surprise if these algorithms can help solve regularized MVO allocation models. However, these two algorithms are not sufficient for defining a general framework. For that, we need to use the alternating direction method of multipliers and proximal gradient methods. Finally, the combination of these four algorithms (CD, ADMM, PO and Dykstra) allows us to consider allocation models that cannot be cast into a QP form.

In this paper, we have first considered allocation models with non-quadratic objective functions. For example, we have used models based on the diversification ratio, Shannon's entropy or the Kullback-Leibler divergence. Second, we have solved regularized MVO models with non-linear penalty functions such as the ℓ_p -norm penalty or the logarithmic barrier. Third, we have discussed how to handle non-linear constraints. For instance, we have imposed constraints on active share, volatility targeting, leverage limits, transaction costs, etc. Most importantly, these three non-QP extensions can be combined.

With the development of quantitative strategies (smart beta, factor investing, alternative risk premia, systematic strategies, robo-advisors, etc.), the asset management industry has dramatically changed over the last five years. This is just the beginning and we think that alternative data, machine learning methods and artificial intelligence will massively shape investment processes in the future. This paper is an illustration of this trend and shows how machine learning optimization algorithms allow to move away from the traditional QP world of portfolio management.

References

- [1] BARANKIN, E.W., and DORFMAN, R. (1956), A Method for Quadratic Programming, *Econometrica*, 24, pp. 340.
- [2] BARANKIN, E.W., and DORFMAN, R. (1958), On Quadratic Programming, *University of California Publications in Statistics*, 2(13), pp. 285-318.
- [3] BAUSCHKE, H.H., and BORWEIN, J.M. (1994), Dykstra's Alternating Projection Algorithm for Two Sets, *Journal of Approximation Theory*, 79(3), pp. 418-443.
- [4] BEALE, E.M.L. (1959), On Quadratic Programming, *Naval Research Logistics Quarterly*, 6(3), pp. 227-243.
- [5] BECK, A. (2017), *First-Order Methods in Optimization*, MOS-SIAM Series on Optimization, 25, SIAM.
- [6] BERA, A.K., and PARK, S.Y. (2008), Optimal Portfolio Diversification Using the Maximum Entropy Principle, *Econometric Reviews*, 27(4-6), pp. 484-512.
- [7] BISHOP, C.M. (1995), *Neural Networks for Pattern Recognition*, Oxford University Press.
- [8] BOURGERON, T., LEZMI, E., and RONCALLI, T. (2018), Robust Asset Allocation for Robo-Advisors, *arXiv*, 1902.05710.
- [9] BRODIE, J., DAUBECHIES, I., DE MOL, C., GIANNONE, D., and LORIS, I. (2009), Sparse and Stable Markowitz Portfolios, *Proceedings of the National Academy of Sciences*, 106(30), pp. 12267-12272.
- [10] BRUDER, B., GAUSSEL, N., RICHARD, J-C., and RONCALLI, T. (2013), Regularization of Portfolio Allocation, *SSRN*, www.ssrn.com/abstract=2767358.
- [11] BRUDER, B., KOSTYUCHYK, N., and RONCALLI, T. (2016), Risk Parity Portfolios with Skewness Risk: An Application to Factor Investing and Alternative Risk Premia, *SSRN*, www.ssrn.com/abstract=2813384.
- [12] BOYD, S., PARIKH, N., CHU, E., PELEATO, B., and ECKSTEIN, J. (2010), Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers, *Foundations and Trends® in Machine learning*, 3(1), pp. 1-122.
- [13] CANDELON, B., HURLIN, C., and TOKPAVI, S. (2012), Sampling Error and Double Shrinkage Estimation of Minimum Variance Portfolios, *Journal of Empirical Finance*, 19(4), pp. 511-527.
- [14] CARMICHAEL, B., KOUMOU, G.B., and MORAN, K. (2018), Rao's Quadratic Entropy and Maximum Diversification Indexation, *Quantitative Finance*, 18(6), pp. 1017-1031.
- [15] CARRASCO, M., and NOUMON, N. (2010), Optimal Portfolio Selection Using Regularization, University of Montréal, *Discussion paper*.
- [16] CHAUX, C., COMBETTES, P.L., PESQUET, J.C., and WAJS, V.R. (2007), A Variational Formulation for Frame-based Inverse Problems, *Inverse Problems*, 23(4), pp. 1495-1518.
- [17] CHOUEIFATY, Y. and COIGNARD, Y. (2008), Toward Maximum Diversification, *Journal of Portfolio Management*, 35(1), pp. 40-51.

- [18] CHOUEIFATY, Y., FROIDURE, T. and REYNIER, J. (2013), Properties of the Most Diversified Portfolio, *Journal of investment strategies*, 2(2), pp. 49-70.
 - [19] COMBETTES, P.L., and MÜLLER, C.L. (2018), Perspective Functions: Proximal Calculus and Applications in High-dimensional Statistics, *Journal of Mathematical Analysis and Applications*, 457(2), pp. 1283-1306.
 - [20] COMBETTES, P.L., and PESQUET, J.C. (2011), Proximal Splitting Methods in Signal Processing, in Bauschke, H.H., Burachik, R.S., Combettes, P.L., Elser, V., Luke, D.R., and Wolkowicz, H. (Eds), *Fixed-point Algorithms for Inverse Problems in Science and Engineering*, Springer Optimization and Its Applications, 48, pp. 185-212, Springer.
 - [21] CORLESS, R.M., GONNET, G.H., HARE, D.E., JEFFREY, D.J., and KNUTH, D.E. (1996), On the Lambert W Function, *Advances in Computational Mathematics*, 5(1), pp. 329-359.
 - [22] CORTES, C., and VAPNIK, V. (1995), Support-vector Networks, *Machine Learning*, 20(3), pp. 273-297.
 - [23] COTTLE, R.W. and INFANGER, G. (2010), Harry Markowitz and the Early History of Quadratic Programming, in Guerard, J.B. (Eds), *Handbook of Portfolio Construction*, Springer, pp. 179-211.
 - [24] DANTZIG, G.B. (1961), Quadratic Programming: A Variant of the Wolfe-Markowitz algorithms, *Operations Research Center*, University of California-Berkeley, Research Report 2.
 - [25] de JONG, M. (2018), Portfolio Optimisation in an Uncertain World, *Journal of Asset Management*, 19(4), pp. 216-221.
 - [26] DEBREU, G. (1952), Definite and Semidefinite Quadratic Forms, *Econometrica*, 20(2), pp. 295-300.
 - [27] DEMIGUEL, V., GARLAPPI, L., NOGALES, F.J., and UPPAL, R. (2009), A Generalized Approach to Portfolio Optimization: Improving Performance by Constraining Portfolio Norms, *Management Science*, 55(5), pp. 798-812.
 - [28] DOUGLAS, J., and RACHFORD, H.H. (1956), On the Numerical Solution of Heat Conduction Problems in Two and Three Space Variables, *Transactions of the American mathematical Society*, 82(2), pp. 421-439.
 - [29] DYKSTRA, R.L. (1983), An Algorithm for Restricted Least Squares Regression, *Journal of the American Statistical Association*, 78(384), pp. 837-842.
 - [30] FRANK, M., and WOLFE, P. (1956), An Algorithm for Quadratic Programming, *Naval Research Logistics Quarterly*, 3, pp. 95-110.
 - [31] FRIEDMAN, J., HASTIE, T., and TIBSHIRANI, R. (2010), Regularization Paths for Generalized Linear Models via Coordinate Descent, *Journal of Statistical Software*, 33(1), pp. 1-22.
 - [32] GABAY, D., and MERCIER, B. (1976), A Dual Algorithm for the Solution of Nonlinear Variational Problems via Finite Element Approximation, *Computers & Mathematics with Applications*, 2(1), pp. 17-40.
-

- [33] GHADIMI, E., TEIXEIRA, A., SHAMES, I., and JOHANSSON, M. (2015), Optimal Parameter Selection for the Alternating Direction Method of Multipliers (ADMM): Quadratic Problems, *IEEE Transactions on Automatic Control*, 60(3), pp. 644-658.
- [34] GISELSSON, P., and BOYD, S. (2017), Linear Convergence and Metric Selection for Douglas-Rachford Splitting and ADMM, *IEEE Transactions on Automatic Control*, 62(2), pp. 532-544.
- [35] GONZALVEZ, J., LEZMI, E., RONCALLI, T., and XU, J. (2019), Financial Applications of Gaussian Processes and Bayesian Optimization, *arXiv*, arxiv.org/abs/1903.04841.
- [36] GOULD, N.I.M., and TOINT, P.L. (2000), A Quadratic Programming Bibliography, *Numerical Analysis Group Internal Report*, 1, 142 pages.
- [37] GRIVEAU-BILLION, T., RICHARD, J-C., and RONCALLI, T. (2013), A Fast Algorithm for Computing High-dimensional Risk Parity Portfolios, *SSRN*, www.ssrn.com/abstract=2325255.
- [38] HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. (2009), *The Elements of Statistical Learning*, Second edition, Springer.
- [39] HE, B.S., YANG, H., and WANG, S.L. (2000), Alternating Direction Method with Self-Adaptive Penalty Parameters for Monotone Variational Inequalities, *Journal of Optimization Theory and applications*, 106(2), pp. 337-356.
- [40] HILDRETH, C. (1957), A Quadratic Programming Procedure, *Naval Research Logistics Quarterly*, 4, pp. 79-85.
- [41] JACOBS, R.A. (1988), Increased Rates of Convergence Through Learning Rate Adaptation, *Neural Networks*, 1(4), pp. 295-307.
- [42] JAGANNATHAN, R., and MA, T. (2003), Risk Reduction in Large Portfolios: Why Imposing the Wrong Constraints Helps, *Journal of Finance*, 58(4), pp. 1651-1684.
- [43] LECUN, Y., BOSER, B., DENKER, J.S., HENDERSON, D., HOWARD, R.E., HUBBARD, W., and JACKEL, L.D. (1989), Backpropagation Applied to Handwritten Zip Code Recognition, *Neural Computation*, 1(4), pp. 541-551.
- [44] LEZMI, E., MALONGO, H., RONCALLI, T., and SOBOTKA, R. (2018), Portfolio Allocation with Skewness Risk: A Practical Guide, *SSRN*, www.ssrn.com/abstract=3201319.
- [45] LINDSTROM, S.B., and SIMS, B. (2018), Survey: Sixty Years of Douglas-Rachford, *arXiv*, 1809.07181.
- [46] LUO, Z.Q., and TSENG, P. (1992), On the Convergence of the Coordinate Descent Method for Convex Differentiable Minimization, *Journal of Optimization Theory and Applications*, 72(1), pp. 7-35.
- [47] LUO, Z.Q., and TSENG, P. (1993), Error Bounds and Convergence Analysis of Feasible Descent Methods: A General Approach, *Annals of Operations Research*, 46(1), pp. 157-178.
- [48] MAILLARD, S., RONCALLI, T. and TEÏLETCHÉ, J. (2010), The Properties of Equally Weighted Risk Contribution Portfolios, *Journal of Portfolio Management*, 36(4), pp. 60-70.

- [49] MANN, H.B. (1943), Quadratic Forms with Linear Constraints, *American Mathematical Monthly*, 50, pp. 430-433.
- [50] MARKOWITZ, H. (1952), Portfolio Selection, *Journal of Finance*, 7(1), pp. 77-91.
- [51] MARKOWITZ, H. (1956), The Optimization of a Quadratic Function Subject to Linear Constraints, *Naval Research Logistics Quarterly*, 3(1-2), pp. 111-133.
- [52] MARTIN, A.D. (1955), Mathematical Programming of Portfolio Selections, *Management Science*, 1(2), pp. 152-166.
- [53] MEUCCI, A. (2009), Managing Diversification, *Risk*, 22(5), pp. 74-79.
- [54] MICHAUD, R.O. (1989), The Markowitz Optimization Enigma: Is ‘Optimized’ Optimal?, *Financial Analysts Journal*, 45(1), pp. 31-42.
- [55] NESTEROV, Y.E. (1983), A Method for Solving the Convex Programming Problem with Convergence Rate $O(k^{-2})$, *Doklady Akademii Nauk SSSR*, 269, pp. 543-547.
- [56] NESTEROV, Y. (2004), *Introductory Lectures on Convex Optimization: A Basic Course*, Applied Optimization, 87, Kluwer Academic Publishers.
- [57] NESTEROV, Y. (2012), Efficiency of Coordinate Descent Methods on Huge-scale Optimization Problems, *SIAM Journal on Optimization*, 22(2), pp. 341-362.
- [58] PARIKH, N., and BOYD, S. (2014), Proximal Algorithms, *Foundations and Trends® in Optimization*, 1(3), pp. 127-239.
- [59] POLYAK, B.T. (1964), Some Methods of Speeding Up the Convergence of Iteration Methods, *USSR Computational Mathematics and Mathematical Physics*, 4(5), pp. 1-17.
- [60] QIAN, E. (2005), Risk Parity Portfolios: Efficient Portfolios Through True Diversification, *Panagora Asset Management*, September.
- [61] RICHARD, J-C., and RONCALLI, T. (2015), Smart Beta: Managing Diversification of Minimum Variance Portfolios, in Jurczenko, E. (Ed.), *Risk-based and Factor Investing*, ISTE Press – Elsevier.
- [62] RICHARD, J-C., and RONCALLI, T. (2019), Constrained Risk Budgeting Portfolios: Theory, Algorithms, Applications & Puzzles, *arXiv*, 1902.05710.
- [63] RONCALLI, T. (2013), *Introduction to Risk Parity and Budgeting*, Chapman & Hall/CRC Financial Mathematics Series.
- [64] RONCALLI, T. (2015), Introducing Expected Returns into Risk Parity Portfolios: A New Framework for Asset Allocation, *Bankers, Markets & Investors*, 138, pp. 18-28.
- [65] SCHERER, B. (2007), *Portfolio Construction & Risk Budgeting*, Third edition, Risk Books.
- [66] TIBSHIRANI, R. (1996), Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society B*, 58(1), pp. 267-288.
- [67] TIBSHIRANI, R.J. (2017), Dykstra’s Algorithm, ADMM, and Coordinate Descent: Connections, Insights, and Extensions, in Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (Eds), *Advances in Neural Information Processing Systems*, 30, pp. 517-528.

- [68] TSENG, P. (1990), Dual Ascent Methods for Problems with Strictly Convex Costs and Linear Constraints: A Unified Approach, *SIAM Journal on Control and Optimization*, 28(1), pp. 214-242.
- [69] TSENG, P. (2001), Convergence of a Block Coordinate Descent Method for Nondifferentiable Minimization, *Journal of Optimization Theory and Applications*, 109(3), pp. 475-494.
- [70] VAPNIK, V. (1998), *Statistical Learning Theory*, John Wiley & Sons.
- [71] WANG, S.L., and LIAO, L.Z. (2001), Decomposition Method with a Variable Parameter for a Class of Monotone Variational Inequality Problems, *Journal of Optimization Theory and Applications*, 109(2), pp. 415-429.
- [72] WESTON, J.F. and BARENEK, W. (1955), Programming Investment Portfolio Construction, *Analysts Journal*, 11(2), pp. 51-55.
- [73] WOLFE, P. (1959), The Simplex Method for Quadratic Programming, *Econometrica*, 27(3), pp. 382-398.
- [74] WRIGHT, S.J. (2015), Coordinate Descent Algorithms, *Mathematical Programming*, 151(1), pp. 3-34
- [75] YU, J.R., LEE, W.Y., and Chiou, W.J.P. (2014), Diversified Portfolios with Different Entropy Measures, *Applied Mathematics and Computation*, 241, pp. 47-63.

Appendix

A Mathematical results

A.1 QP problem when there is a benchmark

Following Roncalli (2013), the excess return $\mathfrak{R}(x | b)$ of Portfolio x with respect to Benchmark b is the difference between the return of the portfolio and the return of the benchmark:

$$\mathfrak{R}(x | b) = \mathfrak{R}(x) - \mathfrak{R}(b) = (x - b)^\top \mathfrak{R}$$

It is easy to show that the expected excess return is equal to:

$$\mu(x | b) = \mathbb{E}[\mathfrak{R}(x | b)] = (x - b)^\top \mu$$

whereas the volatility of the tracking error is given by:

$$\sigma(x | b) = \sigma(\mathfrak{R}(x | b)) = \sqrt{(x - b)^\top \Sigma (x - b)}$$

The objective function is then:

$$\begin{aligned} f(x | b) &= \frac{1}{2} (x - b)^\top \Sigma (x - b) - \gamma (x - b)^\top \mu \\ &= \frac{1}{2} x^\top \Sigma x - x^\top (\gamma \mu + \Sigma b) + \left(\frac{1}{2} b^\top \Sigma b + \gamma b^\top \mu \right) \\ &= \frac{1}{2} x^\top Q x - x^\top R + C \end{aligned}$$

where C is a constant which does not depend on Portfolio x . We recognize a QP problem where $Q = \Sigma$ and $R = \gamma \mu + \Sigma b$.

A.2 Augmented QP formulation of the turnover management problem

The augmented QP problem is defined by:

$$\begin{aligned} X^* &= \arg \min_X \frac{1}{2} X^\top Q X - X^\top R \\ \text{s.t.} \quad &\begin{cases} AX = B \\ CX \leq D \\ \mathbf{0}_{3n} \leq X \leq \mathbf{1}_{3n} \end{cases} \end{aligned}$$

where $X = (x_1, \dots, x_n, x_1^-, \dots, x_n^-, x_1^+, \dots, x_n^+)$ is a $3n \times 1$ vector, Q is a $3n \times 3n$ matrix:

$$Q = \begin{pmatrix} \Sigma & \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} \end{pmatrix}$$

$R = (\gamma \mu, \mathbf{0}_n, \mathbf{0}_n)$ is a $3n \times 1$ vector, A is a $(n + 1) \times 3n$ matrix:

$$A = \begin{pmatrix} \mathbf{1}_n^\top & \mathbf{0}_n^\top & \mathbf{0}_n^\top \\ I_n & I_n & -I_n \end{pmatrix}$$

$B = (1, \bar{x})$ is a $(n + 1) \times 1$ vector, $C = (\mathbf{0}_n^\top \quad \mathbf{1}_n^\top \quad \mathbf{1}_n^\top)$ is a $1 \times 3n$ matrix and $D = \tau^+$.

A.3 Augmented QP formulation of the MVO problem with transaction costs

The augmented QP problem of dimension $3n$ is defined by:

$$\begin{aligned} X^* &= \arg \min_X \frac{1}{2} X^\top Q X - X^\top R \\ \text{s.t. } &\begin{cases} AX = B \\ \mathbf{0}_{3n} \leq X \leq \mathbf{1}_{3n} \end{cases} \end{aligned}$$

where $X = (x_1, \dots, x_n, x_1^-, \dots, x_n^-, x_1^+, \dots, x_n^+)$ is a $3n \times 1$ vector, Q is a $3n \times 3n$ matrix:

$$Q = \begin{pmatrix} \Sigma & \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n} \end{pmatrix}$$

$R = (\gamma\mu, -c^-, -c^+)$ is a $3n \times 1$ vector, A is a $(n+1) \times 3n$ matrix:

$$A = \begin{pmatrix} \mathbf{1}_n^\top & (c^-)^\top & (c^+)^\top \\ I_n & I_n & -I_n \end{pmatrix}$$

and $B = (1, \bar{x})$ is a $(n+1) \times 1$ vector.

A.4 QP problem with a hyperplane constraint

We consider the following QP problem:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top Q x - x^\top R \\ \text{s.t. } &a^\top x = b \end{aligned}$$

The associated Lagrange function is:

$$\mathcal{L}(x; \lambda) = \frac{1}{2} x^\top Q x - x^\top R + \lambda (a^\top x - b)$$

The first-order conditions are then:

$$\begin{cases} \partial_x \mathcal{L}(x; \lambda) = Qx - R + \lambda a = \mathbf{0}_n \\ \partial_\lambda \mathcal{L}(x; \lambda) = a^\top x - b = 0 \end{cases}$$

We obtain $x = Q^{-1}(R + \lambda a)$. Because $a^\top x - b = 0$, we have $a^\top Q^{-1}R + \lambda a^\top Q^{-1}a = b$ and:

$$\lambda^* = \frac{b - a^\top Q^{-1}R}{a^\top Q^{-1}a}$$

The optimal solution is then:

$$x^* = Q^{-1} \left(R + \frac{b - a^\top Q^{-1}R}{a^\top Q^{-1}a} a \right)$$

A.5 Derivation of the soft-thresholding operator

We consider the following equation:

$$cx - v + \lambda \partial |x| \in 0$$

where $c > 0$ and $\lambda > 0$. Since we have $\partial |x| = \text{sign}(x)$, we deduce that:

$$x^* = \begin{cases} c^{-1}(v + \lambda) & \text{if } x^* < 0 \\ 0 & \text{if } x^* = 0 \\ c^{-1}(v - \lambda) & \text{if } x^* > 0 \end{cases}$$

If $x^* < 0$ or $x^* > 0$, then we have $v + \lambda < 0$ or $v - \lambda > 0$. This is equivalent to set $|v| > \lambda > 0$. The case $x^* = 0$ implies that $|v| \leq \lambda$. We deduce that:

$$x^* = c^{-1} \cdot \mathcal{S}(v; \lambda)$$

where $\mathcal{S}(v; \lambda)$ is the soft-thresholding operator:

$$\begin{aligned} \mathcal{S}(v; \lambda) &= \begin{cases} 0 & \text{if } |v| \leq \lambda \\ v - \lambda \text{sign}(v) & \text{otherwise} \end{cases} \\ &= \text{sign}(v) \cdot (|v| - \lambda)_+ \end{aligned}$$

In Figure 6, we have represented the function $\mathcal{S}(v; \lambda)$ when λ is respectively equal to 1 and 2.

Remark 10 The soft-thresholding operator is the proximal operator of the ℓ_1 -norm $f(x) = \|x\|_1$. Indeed, we have $\text{prox}_f(v) = \mathcal{S}(v; 1)$ and $\text{prox}_{\lambda f}(v) = \mathcal{S}(v; \lambda)$.

A.6 The box-constrained QP problem

We consider the box-constrained QP problem:

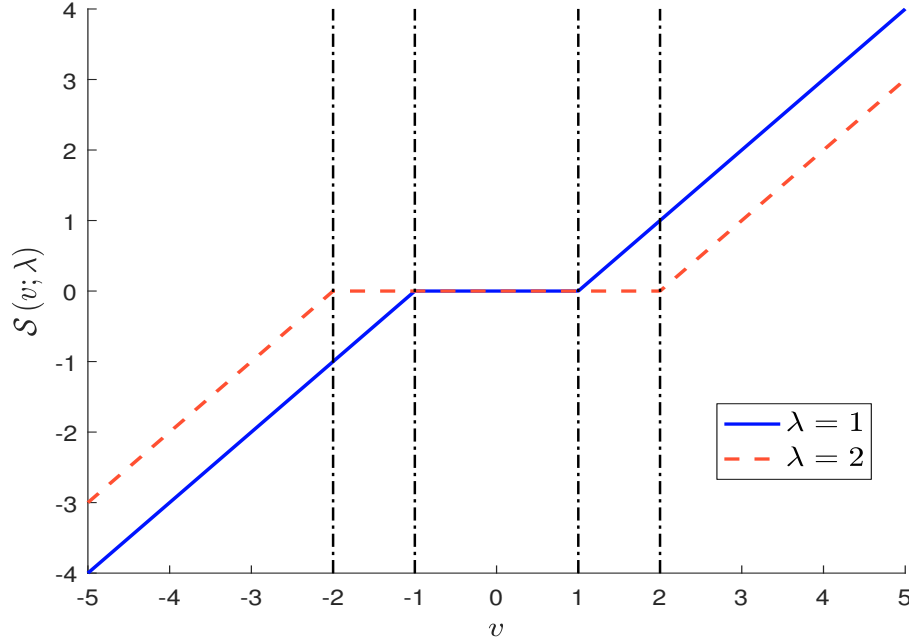
$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} x^\top Q x - x^\top R \\ \text{s.t. } & x^- \leq x \leq x^+ \end{aligned} \tag{44}$$

The objective function is equal to:

$$\begin{aligned} f(x) &= \frac{1}{2} x^\top Q x - x^\top R \\ &= \frac{1}{2} \sum_{i=1}^n x_i \sum_{j=1}^n Q_{i,j} x_j - \sum_{i=1}^n x_i R_i \\ &= \frac{1}{2} \sum_{i=1}^n x_i \left(Q_{i,i} x_i + \sum_{j \neq i} Q_{i,j} x_j \right) - \sum_{i=1}^n x_i R_i \end{aligned}$$

We deduce that:

$$\frac{\partial f(x)}{\partial x_i} = \frac{1}{2} \left(2x_i Q_{i,i} + \sum_{j \neq i} x_j (Q_{i,j} + Q_{j,i}) \right) - R_i$$

Figure 6: Soft-thresholding operator $\mathcal{S}(v; \lambda)$


We notice that:

$$\frac{\partial f(x)}{\partial x_i} = 0 \Leftrightarrow x_i = \frac{1}{Q_{i,i}} \left(R_i - \sum_{j \neq i} x_j \left(\frac{Q_{i,j} + Q_{j,i}}{2} \right) \right)$$

The Lagrange function associated to Problem (44) is equal to:

$$\mathcal{L}(x; \lambda^-, \lambda^+) = f(x) - \sum_{i=1}^n \lambda_i^- (x_i - x_i^-) - \sum_{i=1}^n \lambda_i^+ (x_i^+ - x_i)$$

The first-order condition is then:

$$\frac{\partial \mathcal{L}(x; \lambda^-, \lambda^+)}{\partial x_i} = \frac{\partial f(x)}{\partial x_i} - \lambda_i^- + \lambda_i^+ = 0$$

Since the Kuhn-Tucker conditions are:

$$\begin{cases} \min(\lambda_i^-, x_i - x_i^-) = 0 \\ \min(\lambda_i^+, x_i^+ - x_i) = 0 \end{cases}$$

we obtain three cases:

1. If no bound is reached, we have $\lambda_i^- = \lambda_i^+ = 0$ and the solution is equal to:

$$x_i^* = \frac{1}{Q_{i,i}} \left(R_i - \sum_{j \neq i} x_j \left(\frac{Q_{i,j} + Q_{j,i}}{2} \right) \right)$$

2. If the lower bound is reached, we have $\lambda_i^- > 0$, $\lambda_i^+ = 0$ and $x_i^* = x_i^-$.

3. If the upper bound is reached, we have $\lambda_i^- = 0$, $\lambda_i^+ > 0$ and $x_i^* = x_i^+$.

A.7 ADMM algorithm

The optimization problem is defined as:

$$\begin{aligned} \{x^*, y^*\} &= \arg \min_{(x, y)} f_x(x) + f_y(y) \\ \text{s.t. } &Ax + By = c \end{aligned} \quad (45)$$

The derivation of the algorithm is fully explained in Boyd *et al.* (2011). For that, they consider the augmented Lagrange function:

$$\mathcal{L}(x, y; \lambda, \varphi) = f_x(x) + f_y(y) + \lambda^\top (Ax + By - c) + \frac{\varphi}{2} \|Ax + By - c\|_2^2 \quad (46)$$

where $\varphi > 0$. According to Boyd *et al.* (2011), the ℓ_2 -norm penalty adds robustness to the dual ascent method and accelerates its convergence. The ADMM algorithm uses the property that the objective function is separable, and consists of the following iterations:

$$\begin{aligned} x^{(k+1)} &= \arg \min_x \mathcal{L}(x, y^{(k)}; \lambda^{(k)}, \varphi) \\ &= \arg \min_x \left\{ f_x(x) + \lambda^{(k)\top} (Ax + By^{(k)} - c) + \frac{\varphi}{2} \|Ax + By^{(k)} - c\|_2^2 \right\} \end{aligned}$$

and:

$$\begin{aligned} y^{(k+1)} &= \arg \min_y \mathcal{L}(x^{(k+1)}, y; \lambda^{(k)}, \varphi) \\ &= \arg \min_y \left\{ f_y(y) + \lambda^{(k)\top} (Ax^{(k+1)} + By - c) + \frac{\varphi}{2} \|Ax^{(k+1)} + By - c\|_2^2 \right\} \end{aligned}$$

The update for the dual variable λ is then:

$$\lambda^{(k+1)} = \lambda^{(k)} + \varphi (Ax^{(k+1)} + By^{(k+1)} - c)$$

We repeat the iterations until convergence.

Boyd *et al.* (2011) notice that the previous algorithm can be simplified. Let $r = Ax + By - c$ be the (primal) residual. By combining linear and quadratic terms, we have:

$$\lambda^\top r + \frac{\varphi}{2} \|r\|_2^2 = \frac{\varphi}{2} \|r + u\|_2^2 - \frac{\varphi}{2} \|u\|_2^2$$

where $u = \varphi^{-1}\lambda$ is the *scaled* dual variable. We can then write the Lagrange function (46) as follows:

$$\begin{aligned} \mathcal{L}(x, y; u, \varphi) &= f_x(x) + f_y(y) + \frac{\varphi}{2} \|Ax + By - c + u\|_2^2 - \frac{\varphi}{2} \|u\|_2^2 \\ &= f_x(x) + f_y(y) + \frac{\varphi}{2} \|Ax + By - c + u\|_2^2 - \frac{1}{2\varphi} \|\lambda\|_2^2 \end{aligned}$$

Since the last term is a constant, we deduce that the x - and y -updates become:

$$\begin{aligned} x^{(k+1)} &= \arg \min_x \mathcal{L}(x, y^{(k)}; u^{(k)}, \varphi) \\ &= \arg \min_x \left\{ f_x(x) + \frac{\varphi}{2} \|Ax + By^{(k)} - c + u^{(k)}\|_2^2 \right\} \end{aligned} \quad (47)$$

and:

$$\begin{aligned} y^{(k+1)} &= \arg \min_y \mathcal{L} \left(x^{(k+1)}, y; u^{(k)}, \varphi \right) \\ &= \arg \min_y \left\{ f_y(y) + \frac{\varphi}{2} \|Ax^{(k+1)} + By - c + u^{(k)}\|_2^2 \right\} \end{aligned} \quad (48)$$

For the scaled dual variable u , we have:

$$\begin{aligned} u^{(k+1)} &= u^{(k)} + r^{(k+1)} \\ &= u^{(k)} + \left(Ax^{(k+1)} + By^{(k+1)} - c \right) \end{aligned} \quad (49)$$

where $r^{(k+1)} = Ax^{(k+1)} + By^{(k+1)} - c$ is the primal residual at iteration $k + 1$. Boyd *et al.* (2011) also define the variable $s^{(k+1)} = \varphi A^\top B (y^{(k+1)} - y^{(k)})$ and refer to $s^{(k+1)}$ as the dual residual at iteration $k + 1$.

Under the assumption that the traditional Lagrange function $\mathcal{L}(x, y; \lambda, 0)$ has a saddle point, one can prove that the residual $r^{(k)}$ converges to zero, the objective function $f_x(x^{(k)}) + f_y(y^{(k)})$ converges to the optimal value $f_x(x^*) + f_y(y^*)$, and the dual variable $\lambda^{(k)} = \varphi u^{(k)}$ converges to a dual optimal point. However, the rate of convergence is not known and the primal variables $x^{(k)}$ and $y^{(k)}$ do not necessarily converge to the optimal values x^* and y^* . In general, the stopping criterion is defined with respect to the residuals:

$$\begin{cases} \|r^{(k+1)}\|_2 \leq \varepsilon \\ \|s^{(k+1)}\|_2 \leq \varepsilon' \end{cases}$$

Typical values when implementing this stopping criterion are $\varepsilon = \varepsilon' = 10^{-15}$ (Bourgeron *et al.*, 2018).

From a theoretical point of view, the convergence holds regardless of the choice of the penalization parameter $\varphi > 0$. But the choice of φ affects the convergence rate (Ghadimi *et al.*, 2015; Giselsson and Boyd, 2017). In practice, the penalization parameter φ may be changed at each iteration, implying that φ is replaced by $\varphi^{(k)}$ and the scaled dual variable u^k is equal to $\lambda^{(k)}/\varphi^{(k)}$. This may improve the convergence and make the performance independent of the initial choice $\varphi^{(0)}$. To update $\varphi^{(k)}$ in practice, He *et al.* (2000) and Wang and Liao (2001) provide a simple and efficient scheme. On the one hand, the x - and y -updates in ADMM essentially come from placing a penalty on $\|r^{(k)}\|_2^2$. As a consequence, if $\varphi^{(k)}$ is large, $\|r^{(k)}\|_2^2$ tends to be small. On the other hand, $s^{(k)}$ depends linearly on φ . As a consequence, if $\varphi^{(k)}$ is small, $\|s^{(k)}\|_2^2$ is small. To keep $\|r^{(k)}\|_2^2$ and $\|s^{(k)}\|_2^2$ within a factor μ , one may consider:

$$\varphi^{(k+1)} = \begin{cases} \tau \varphi^{(k)} & \text{if } \|r^{(k)}\|_2^2 > \mu \|s^{(k)}\|_2^2 \\ \varphi^{(k)}/\tau' & \text{if } \|s^{(k)}\|_2^2 > \mu \|r^{(k)}\|_2^2 \\ \varphi^{(k)} & \text{otherwise} \end{cases}$$

where μ , τ and τ' are parameters that are greater than one. In practice, we use $\varphi^{(0)} = 1$, $u^{(0)} = \mathbf{0}_p$, $\mu = 10^3$ and $\tau = \tau' = 2$.

Remark 11 The constant case $\varphi^{(k+1)} = \varphi^{(k)} = \varphi^{(0)}$ is obtained by setting $\tau = \tau' = 1$.

A.8 Proximal operators

A.8.1 Pointwise maximum function

The unit simplex is the generalization of the triangle:

$$\mathbb{S}_n = \left\{ x \in [0, 1]^n, \theta_i \geq 0 : x = \sum_{i=0}^n \theta_i e_i, \sum_{i=0}^n \theta_i = 1, \mathbf{1}_n^\top x \leq 1 \right\}$$

where e_0 is the zero vector and e_i are the unit vectors for $i \geq 1$. Beck (2017) shows that $\mathcal{P}_{\mathbb{S}_n}(v) = (v - \mu^* \mathbf{1}_n)_+$ where μ^* is the root of the equation $\mathbf{1}_n^\top (v - \mu^* \mathbf{1}_n)_+ = 1$. In the case of the pointwise maximum function $f(x) = \max x$, the Moreau decomposition gives:

$$\begin{aligned} \mathbf{prox}_{\lambda \max x}(v) &= v - \lambda \mathcal{P}_{\mathbb{S}_n} \left(\frac{v}{\lambda} \right) \\ &= v - \lambda \left(\frac{v}{\lambda} - \mu^* \mathbf{1}_n \right)_+ \end{aligned}$$

where μ^* is the root of the equation:

$$\begin{aligned} \mathbf{1}_n^\top \left(\frac{v}{\lambda} - \mu^* \mathbf{1}_n \right)_+ = 1 &\Leftrightarrow \sum_{i=1}^n \left(\frac{v_i}{\lambda} - \mu^* \right)_+ = 1 \\ &\Leftrightarrow \sum_{i=1}^n (v_i - s^*)_+ = \lambda \end{aligned}$$

where $s^* = \lambda \mu^*$. It follows that:

$$\begin{aligned} \mathbf{prox}_{\lambda \max x}(v) &= v - (v - s^* \mathbf{1}_n)_+ \\ &= \min(v, s^*) \end{aligned}$$

A.8.2 ℓ_2 -norm function

The projection onto the unit ball $\mathcal{B}_2(\mathbf{0}, 1) = \{x \in \mathbb{R}^n : \|x\|_2 \leq 1\}$ is equal to:

$$\mathcal{P}_{\mathcal{B}_2(\mathbf{0}, 1)}(v) = \begin{cases} v & \text{if } \|v\|_2 \leq 1 \\ \frac{v}{\|v\|_2} & \text{if } \|v\|_2 > 1 \end{cases}$$

Since we have:

$$\mathbf{prox}_{\lambda \|x\|_2}(v) + \lambda \mathcal{P}_{\mathcal{B}_2(\mathbf{0}, 1)} \left(\frac{v}{\lambda} \right) = v$$

we deduce that:

$$\begin{aligned} \mathbf{prox}_{\lambda \|x\|_2}(v) &= v - \lambda \mathcal{P}_{\mathcal{B}_2(\mathbf{0}, 1)} \left(\frac{v}{\lambda} \right) \\ &= \begin{cases} v - \lambda \left(\frac{v}{\lambda} \right) & \text{if } \|v\|_2 \leq \lambda \\ v - \lambda \frac{\frac{v}{\lambda}}{\left\| \frac{v}{\lambda} \right\|_2} & \text{if } \|v\|_2 > \lambda \end{cases} \\ &= \begin{cases} 0 & \text{if } \|v\|_2 \leq \lambda \\ v - \frac{\lambda v}{\|v\|_2} & \text{if } \|v\|_2 > \lambda \end{cases} \\ &= \left(1 - \frac{\lambda}{\max(\lambda, \|v\|_2)} \right) v \end{aligned}$$

A.8.3 Scaling and translation

Let $g(x) = f(ax + b)$ where $a \neq 0$. Using the change of variable $y = ax + b$, we have:

$$\begin{aligned} \mathbf{prox}_g(v) &= \arg \min_x \left\{ g(x) + \frac{1}{2} \|x - v\|_2^2 \right\} \\ &= \arg \min_x \left\{ f(ax + b) + \frac{1}{2} \|x - v\|_2^2 \right\} \\ &= \arg \min_y \left\{ f(y) + \frac{1}{2} \left\| \frac{y - b}{a} - v \right\|_2^2 \right\} \end{aligned}$$

We deduce that:

$$\begin{aligned} f(y) + \frac{1}{2} \left\| \frac{y - b}{a} - v \right\|_2^2 &= f(y) + \frac{1}{2a^2} \|y - b - av\|_2^2 \\ &= \frac{1}{a^2} \left(a^2 f(y) + \frac{1}{2} \|y - (av + b)\|_2^2 \right) \end{aligned}$$

We conclude that $y^* = \mathbf{prox}_{a^2 f}(av + b)$ and:

$$\begin{aligned} \mathbf{prox}_g(v) &= \frac{y^* - b}{a} \\ &= \frac{\mathbf{prox}_{a^2 f}(av + b) - b}{a} \end{aligned}$$

A.8.4 Projection onto the ℓ_1 ball

We have:

$$\begin{aligned} x &= \mathcal{P}_{\mathcal{B}_1(c, r)}(v) \\ &= \mathcal{P}_{\mathcal{B}_1(\mathbf{0}_n, r)}(v - c) + c \\ &= (v - c) - \text{sign}(v - c) \odot \mathbf{prox}_{r \max x}(|v - c|) + c \\ &= v - \text{sign}(v - c) \odot \min(|v - c|, s^*) \end{aligned}$$

where s^* is the solution of the following equation:

$$s^* = \left\{ s \in \mathbb{R} : \sum_{i=1}^n (|v_i - c_i| - s)_+ = r \right\}$$

Remark 12 $\mathcal{P}_{\mathcal{B}_1(c, r)}(v)$ is sometimes expressed using the soft-thresholding operator (Beck, 2017, page 151), but the two formulas are equivalent.

A.8.5 Projection onto the ℓ_2 ball

We have:

$$\begin{aligned} x &= \mathcal{P}_{\mathcal{B}_2(c, r)}(v) \\ &= \mathcal{P}_{\mathcal{B}_2(\mathbf{0}_n, r)}(v - c) + c \\ &= (v - c) - \mathbf{prox}_{r\|x\|_2}(v - c) + c \\ &= v - \left(1 - \frac{r}{\max(r, \|v - c\|_2)} \right) (v - c) \\ &= c + \frac{r}{\max(r, \|v - c\|_2)} (v - c) \end{aligned} \tag{50}$$

A.8.6 ℓ_2 -penalized logarithmic barrier function

We note $f_1(x) = -\lambda \sum_{i=1}^n b_i \ln x_i$ and $f_2(x) = \mathbb{1}_\Omega(x)$ where $\Omega = \mathcal{B}_2(c, r)$. The Dykstra's algorithm becomes:

$$\begin{cases} v_x^{(k)} = y^{(k)} + z_1^{(k)} \\ x^{(k+1)} = \text{prox}_{f_1}(v_x^{(k)}) \\ z_1^{(k+1)} = y^{(k)} + z_1^{(k)} - x^{(k+1)} \\ v_y^{(k)} = x^{(k+1)} + z_2^{(k)} \\ y^{(k+1)} = \text{prox}_{f_2}(v_y^{(k)}) \\ z_2^{(k+1)} = x^{(k+1)} + z_2^{(k)} - y^{(k+1)} \end{cases}$$

It follows that:

$$x^{(k+1)} = \frac{v_x^{(k)} + \sqrt{v_x^{(k)} \odot v_x^{(k)} + 4\lambda b}}{2}$$

and:

$$\begin{aligned} y^{(k+1)} &= \mathcal{P}_{\mathcal{B}_2(c, r)}(v_y^{(k)}) \\ &= c + \frac{r}{\max(r, \|v_y^{(k)} - c\|_2)} (v_y^{(k)} - c) \end{aligned}$$

A.8.7 Quadratic function

Let $f(x) = \frac{1}{2}x^\top Qx - x^\top R$. We have:

$$\begin{aligned} \text{prox}_f(v) &= \arg \min_x \left\{ \frac{1}{2}x^\top Qx - x^\top R + \frac{1}{2}\|x - v\|_2^2 \right\} \\ &= \arg \min_x \left\{ \frac{1}{2}x^\top (Q + I_n)x - x^\top (R + v) + \frac{1}{2}v^\top v \right\} \\ &= (Q + I_n)^{-1} (R + v) \end{aligned}$$

A.8.8 Projection onto the intersection of a ℓ_2 ball and a box

We note $f_1(x) = \mathbb{1}_{\Omega_1}(x)$ and $f_2(x) = \mathbb{1}_{\Omega_2}(x)$ where $\Omega_1 = \mathcal{B}_2(c, r)$ and $\Omega_2 = \mathcal{B}_{\text{ox}}[x^-, x^+] = \{x \in \mathbb{R}^n : x^- \leq x \leq x^+\}$. The Dykstra's algorithm becomes:

$$\begin{cases} x^{(k+1)} = c + \frac{r}{\max(r, \|y^{(k)} + z_1^{(k)} - c\|_2)} (y^{(k)} + z_1^{(k)} - c) \\ z_1^{(k+1)} = y^{(k)} + z_1^{(k)} - x^{(k+1)} \\ y^{(k+1)} = \mathcal{T}(x^{(k+1)} + z_2^{(k)}; x^-, x^+) \\ z_2^{(k+1)} = x^{(k+1)} + z_2^{(k)} - y^{(k+1)} \end{cases}$$

This algorithm is denoted by $\mathcal{P}_{\mathcal{B}_{\text{ox}}-\mathcal{B}_{\text{all}}}(v; x^-, x^+, c, r)$.

A.8.9 Shannon's entropy and Kullback-Leibler divergence

If we consider the scalar function $f(x) = \lambda x \ln(x/\tilde{x})$ where \tilde{x} is a constant, we have:

$$\begin{aligned} \lambda f(x) + \frac{1}{2} \|x - v\|_2^2 &= \lambda x \ln \frac{x}{\tilde{x}} + \frac{1}{2} (x - v)^2 \\ &= \lambda x \ln \frac{x}{\tilde{x}} + \frac{1}{2} x^2 - xv + \frac{1}{2} v^2 \end{aligned}$$

The first-order condition is:

$$\begin{aligned} \lambda \frac{1}{\tilde{x}} + \lambda \ln \frac{x}{\tilde{x}} + x - v &= 0 \Leftrightarrow \ln x + \lambda^{-1} x = \ln \tilde{x} + \lambda^{-1} v - \frac{1}{\tilde{x}} \\ &\Leftrightarrow e^{\ln x + \lambda^{-1} x} = e^{\lambda^{-1} v - \frac{1}{\tilde{x}} + \ln \tilde{x}} \\ &\Leftrightarrow x e^{\lambda^{-1} x} = \tilde{x} e^{\lambda^{-1} v - \frac{1}{\tilde{x}}} \\ &\Leftrightarrow (\lambda^{-1} x) e^{(\lambda^{-1} x)} = \lambda^{-1} \tilde{x} e^{\lambda^{-1} v - \frac{1}{\tilde{x}}} \end{aligned}$$

We deduce that the root is equal to:

$$x^* = \lambda W \left(\frac{\tilde{x} e^{\lambda^{-1} v - \frac{1}{\tilde{x}}}}{\lambda} \right)$$

where $W(x)$ is the Lambert W function satisfying $W(x) e^{W(x)} = x$ (Corless *et al.*, 1996). In the case of the Kullback-Liebler divergence $\text{KL}(x) = \sum_{i=1}^n x_i \ln(x_i/\tilde{x}_i)$, it follows that:

$$\mathbf{prox}_{\lambda \text{KL}(v|\tilde{x})}(v) = \lambda \begin{pmatrix} W \left(\lambda^{-1} \tilde{x}_1 e^{\lambda^{-1} v_1 - \tilde{x}_1^{-1}} \right) \\ \vdots \\ W \left(\lambda^{-1} \tilde{x}_n e^{\lambda^{-1} v_n - \tilde{x}_n^{-1}} \right) \end{pmatrix}$$

Remark 13 The proximal of Shannon's entropy $\text{SE}(x) = -\sum_{i=1}^n x_i \ln x_i$ is a special case of the previous result³⁸ with $\tilde{x}_i = 1$:

$$\mathbf{prox}_{\lambda \text{SE}(x)}(v) = \lambda \begin{pmatrix} W \left(\lambda^{-1} e^{\lambda^{-1} v_1 - 1} \right) \\ \vdots \\ W \left(\lambda^{-1} e^{\lambda^{-1} v_n - 1} \right) \end{pmatrix}$$

This result has been first obtained by Chaux *et al.* (2007).

A.8.10 Projection onto the complement $\bar{\mathcal{B}}_2(c, r)$ of the ℓ_2 ball

We consider the following proximal problem:

$$x^* = \arg \min_x \left\{ \mathbf{1}_{\Omega}(x) + \frac{1}{2} \|x - v\|_2^2 \right\}$$

where:

$$\Omega = \{x \in \mathbb{R}^n : \|x - c\|_2 \geq r\}$$

³⁸We use the fact that $\max \text{SE}(x) = \min - \text{SE}(x)$.

This problem is equivalent to:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} (x - v)^\top (x - v) \\ \text{s.t. } & (x - c)^\top (x - c) - r^2 \geq 0 \end{aligned}$$

We deduce that the Lagrange function is equal to:

$$\mathcal{L}(x; \lambda) = \frac{1}{2} (x - v)^\top (x - v) - \lambda \left((x - c)^\top (x - c) - r^2 \right)$$

The first-order condition is:

$$\frac{\partial \mathcal{L}(x; \lambda)}{\partial x} = x - v - 2\lambda (x - c) = \mathbf{0}_n$$

whereas the KKT condition is $\min \left(\lambda, (x - c)^\top (x - c) - r^2 \right) = 0$. We distinguish two cases:

1. If $\lambda = 0$, this means that $x^* = v$ and $(x - c)^\top (x - c) - r^2 > 0$.
2. If $\lambda > 0$, we have $(x - c)^\top (x - c) = r^2$. Then we obtain the following system:

$$\begin{cases} x - v - 2\lambda (x - c) = \mathbf{0}_n \\ (x - c)^\top (x - c) = r^2 \end{cases}$$

We deduce that:

$$(x - c) - (v - c) - 2\lambda (x - c) = \mathbf{0}_n \quad (51)$$

and:

$$(x - c)^\top (x - c) - (x - c)^\top (v - c) - 2\lambda (x - c)^\top (x - c) = 0$$

It follows that $r^2 - (x - c)^\top (v - c) - 2\lambda r^2 = 0$, meaning that:

$$\lambda^* = \frac{r^2 - (x - c)^\top (v - c)}{2r^2}$$

We notice that:

$$\begin{aligned} (51) &\Leftrightarrow (x - c) - (v - c) - 2 \frac{r^2 - (x - c)^\top (v - c)}{2r^2} (x - c) = \mathbf{0}_n \\ &\Leftrightarrow -r^2 (v - c) + (x - c)^\top (v - c) (x - c) = \mathbf{0}_n \\ &\Leftrightarrow (x - c)^\top (v - c) (x - c) = r^2 (v - c) \end{aligned}$$

Because $(x - c)^\top (v - c)$ is a scalar, we deduce that $x - c$ and $v - c$ are two collinear vectors:

$$x - c = r \frac{(v - c)}{\|v - c\|_2}$$

The optimal solution is:

$$x^* = c + r \frac{(v - c)}{\|v - c\|_2}$$

Combining the two cases gives:

$$\mathbf{prox}_{\mathbf{I}_\Omega(x)}(v) = c + \frac{r}{\min(r, \|v - c\|_2)} (v - c)$$

This is the formula of the projection onto the ℓ_2 ball, but the minimum function has replaced the maximum function³⁹.

³⁹See Equation (50) on page 58.

A.8.11 Projection onto the complement $\bar{\mathcal{B}}_1(c, r)$ of the ℓ_1 ball

This proximal problem associated with $\bar{\mathcal{B}}_1(\mathbf{0}_n, r)$ is:

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} (x - v)^\top (x - v) \\ \text{s.t. } &\|x\|_1 \geq r \end{aligned}$$

We deduce that the Lagrange function is equal to:

$$\mathcal{L}(x; \lambda) = \frac{1}{2} (x - v)^\top (x - v) - \lambda (\|x\|_1 - r)$$

The first-order condition is:

$$\frac{\partial \mathcal{L}(x; \lambda)}{\partial x} = x - v - \lambda \text{sign}(x) = \mathbf{0}_n$$

whereas the KKT condition is $\min(\lambda, \|x\|_1 - r) = 0$. We distinguish two cases:

1. If $\lambda = 0$, this means that $x^* = v$ and $\|x\|_1 \geq r$.
2. If $\lambda > 0$, we have $\|x\|_1 = r$. Then we obtain the following system of equations:

$$\begin{cases} x - v = \lambda \text{sign}(x) \\ \|x\|_1 = r \end{cases}$$

The first condition gives that $x - v$ is a vector whose elements are $+\lambda$ and/or $-\lambda$, whereas the second condition shows that x is on the surface of the ℓ_1 ball. Unfortunately, there is no unique solution. This is why we assume that $\text{sign}(x) = \text{sign}(v)$ and we modify the sign function: $\text{sign}(a) = 1$ if $a \geq 0$ and $\text{sign}(a) = -1$ if $a < 0$. In this case, there is a unique solution $x^* = v + \lambda^* \text{sign}(v)$ where $\lambda^* = n^{-1}(r - \|v\|_1)$ because $|v + \lambda \text{sign}(v)| = |v| + \lambda$.

Combining the two cases implies that:

$$\mathcal{P}_{\bar{\mathcal{B}}_1(\mathbf{0}_n, r)}(v) = \begin{cases} v & \text{if } \|v\|_1 \geq r \\ v + \text{sign}(v) \odot n^{-1}(r - \|v\|_1) & \text{if } \|v\|_1 < r \end{cases}$$

Using the translation property, we deduce that:

$$\mathcal{P}_{\bar{\mathcal{B}}_1(c, r)}(v) = v + \text{sign}(v - c) \odot \frac{\max(r - \|v - c\|_1, 0)}{n}$$

A.8.12 The bid-ask linear cost function

If we consider the scalar function:

$$f(x) = \alpha(\gamma - x)_+ + \beta(x - \gamma)_+$$

we have:

$$\begin{aligned} f_v(x) &= \lambda f(x) + \frac{1}{2} \|x - v\|_2^2 \\ &= \lambda \alpha (\gamma - x)_+ + \lambda \beta (x - \gamma)_+ + \frac{1}{2} x^2 - xv + \frac{1}{2} v^2 \end{aligned}$$

Following Beck (2017), we distinguish three cases:

1. If $f_v(x) = \lambda\alpha(\gamma - x) + \frac{1}{2}x^2 - xv + \frac{1}{2}v^2$, then $f'_v(x) = -\lambda\alpha + x - v$ and $x^* = v + \lambda\alpha$. This implies that $\gamma - x^* > 0$ or $v < \gamma - \lambda\alpha$.
2. If $f_v(x) = \lambda\beta(x - \gamma)_+ + \frac{1}{2}x^2 - xv + \frac{1}{2}v^2$, then $f'_v(x) = \lambda\beta + x - v$ and $x^* = v - \lambda\beta$. This implies that $x^* - \gamma > 0$ or $v < \gamma + \lambda\beta$.
3. If $v \in [\gamma - \lambda\alpha, \gamma + \lambda\beta]$, the minimum is not obtained at a point of differentiability. Since γ is the only point of non-differentiability, we obtain $x^* = \gamma$.

Therefore, we can write the proximal operator in the following compact form:

$$x^* = \gamma + (v - \gamma - \lambda\beta)_+ - (v - \gamma + \lambda\alpha)_-$$

where x_- and x_+ are the negative part and the positive part of x . If we consider the vector-value function $f(x) = \sum_{i=1}^n \alpha_i(\gamma_i - x_i)_+ + \beta_i(x_i - \gamma_i)_+$, we deduce that:

$$\text{prox}_{\lambda f(x)}(v) = \gamma + \mathcal{S}(v - \gamma; \lambda\alpha, \lambda\beta)$$

where $\mathcal{S}(v; \lambda_-, \lambda_+) = (v - \lambda_+)_+ - (v + \lambda_-)_-$ is the two-sided soft-thresholding operator.

A.9 The QP form of the ADMM-QP problem

We have:

$$\begin{aligned} f_{\text{QP}}(x) &= f_{\text{MVO}}(x) + f_{\ell_2}(x) \\ &= \frac{1}{2}(x - b)^\top \Sigma_t (x - b) - \gamma(x - b)^\top \mu_t + \frac{1}{2}\varrho_2 \|\Gamma_2(x - x_t)\|_2^2 + \frac{1}{2}\tilde{\varrho}_2 \|\tilde{\Gamma}_2(x - \tilde{x})\|_2^2 \\ &= \frac{1}{2}x^\top \Sigma_t x - x^\top \Sigma_t b + \frac{1}{2}b^\top \Sigma_t b - \gamma x^\top \mu_t + \gamma b^\top \mu_t + \\ &\quad \frac{1}{2}x^\top (\varrho_2 \Gamma_2^\top \Gamma_2) x - x^\top (\varrho_2 \Gamma_2^\top \Gamma_2) x_t + \frac{1}{2}x_t^\top (\varrho_2 \Gamma_2^\top \Gamma_2) x_t + \\ &\quad \frac{1}{2}x^\top (\tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2) x - x^\top (\tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2) \tilde{x} + \frac{1}{2}\tilde{x}^\top (\tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2) \tilde{x} \\ &= \frac{1}{2}x^\top (\Sigma_t + \varrho_2 \Gamma_2^\top \Gamma_2 + \tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2) x - x^\top (\gamma \mu_t + \Sigma_t b + \varrho_2 \Gamma_2^\top \Gamma_2 x_t + \tilde{\varrho}_2 \tilde{\Gamma}_2^\top \tilde{\Gamma}_2 \tilde{x}) + \\ &\quad \gamma b^\top \mu_t + \frac{1}{2}(b^\top \Sigma_t b + \varrho_2 x_t^\top \Gamma_2^\top \Gamma_2 x_t + \tilde{\varrho}_2 \tilde{x}^\top \tilde{\Gamma}_2^\top \tilde{\Gamma}_2 \tilde{x}) \end{aligned}$$

A.10 The CCD algorithm of a QP form with a logarithmic barrier

We consider the following optimization problem:

$$x^* = \arg \min_x \frac{1}{2}x^\top Qx - x^\top R - \sum_{i=1}^n \lambda_i \ln x_i$$

where Q is a positive-definite matrix and $\lambda_i > 0$. The first-order condition with respect to coordinate x_i is:

$$(Qx)_i - R_i - \frac{\lambda_i}{x_i} = 0$$

It follows that $x_i(Qx)_i - R_i x_i - \lambda_i = 0$ or equivalently:

$$Q_{i,i}x_i^2 + \left(\sum_{j \neq i} x_j Q_{i,j} - R_i \right) x_i - \lambda_i = 0$$

The polynomial function is convex because we have $Q_{i,i} > 0$. Since the product of the roots is negative⁴⁰, we have two solutions with opposite signs. We deduce that the solution is the positive root of the second-degree equation:

$$x_i^* = \frac{R_i - \sum_{j \neq i} x_j Q_{i,j} + \sqrt{\left(\sum_{j \neq i} x_j Q_{i,j} - R_i\right)^2 + 4\lambda_i Q_{i,i}}}{2Q_{i,i}}$$

It follows that CCD algorithm is:

$$x_i^{(k+1)} = \frac{R_i - \sum_{j < i} x_j^{(k+1)} Q_{i,j} - \sum_{j > i} x_j^{(k)} Q_{i,j}}{2Q_{i,i}} + \frac{\sqrt{\left(\sum_{j < i} x_j^{(k+1)} Q_{i,j} + \sum_{j > i} x_j^{(k)} Q_{i,j} - R_i\right)^2 + 4\lambda_i Q_{i,i}}}{2Q_{i,i}}$$

B Data

Parameter set #1 We consider a capitalization-weighted stock index, which is composed of eight stocks. The weights of this benchmark are equal to 23%, 19%, 17%, 9%, 8%, 6% and 5%. We assume that their volatilities are 21%, 20%, 40%, 18%, 35%, 23%, 7% and 29%. The correlation matrix is defined as follows:

$$\rho = \begin{pmatrix} 100\% & & & & & & & \\ 80\% & 100\% & & & & & & \\ 70\% & 75\% & 100\% & & & & & \\ 60\% & 65\% & 90\% & 100\% & & & & \\ 70\% & 50\% & 70\% & 85\% & 100\% & & & \\ 50\% & 60\% & 70\% & 80\% & 60\% & 100\% & & \\ 70\% & 50\% & 70\% & 75\% & 80\% & 50\% & 100\% & \\ 60\% & 65\% & 70\% & 75\% & 65\% & 70\% & 80\% & 100\% \end{pmatrix}$$

Parameter set #2 We consider a universe of eight stocks. We assume that their volatilities are 25%, 20%, 15%, 18%, 30%, 20%, 15% and 35%. The correlation matrix is defined as follows:

$$\rho = \begin{pmatrix} 100\% & & & & & & & \\ 20\% & 100\% & & & & & & \\ 55\% & 60\% & 100\% & & & & & \\ 60\% & 60\% & 60\% & 100\% & & & & \\ 60\% & 60\% & 60\% & 60\% & 100\% & & & \\ 60\% & 60\% & 60\% & 60\% & 60\% & 100\% & & \\ 60\% & 60\% & 60\% & 60\% & 60\% & 60\% & 100\% & \\ 60\% & 60\% & 60\% & 60\% & 60\% & 60\% & 60\% & 100\% \end{pmatrix}$$

C Notations

- $\mu = (\mu_1, \dots, \mu_n)$ is the vector of expected return.

⁴⁰We have $-Q_{i,i}\lambda_i < 0$.

- $\Sigma = [\rho_{i,j}\sigma_i\sigma_j]_{i,j=1}^{i,j=1}$ is the covariance matrix where σ_i is the volatility of Asset i and $\rho_{i,j}$ is the correlation between Asset i and Asset j .
- b is the vector of benchmark weights.
- $\mu(x) = x^\top \mu$ is the expected return of Portfolio x .
- $\sigma(x) = \sqrt{x^\top \Sigma x}$ is the volatility of Portfolio x .
- $\mu(x|b) = (x-b)^\top \mu$ is the expected excess return of Portfolio x with respect to Benchmark b .
- $\sigma(x|b) = \sqrt{(x-b)^\top \Sigma (x-b)}$ is the tracking error volatility of Portfolio x with respect to Benchmark b .
- $\mathcal{R}(x)$ is a convex risk measure.
- $\mathcal{RB} = (\mathcal{RB}_1, \dots, \mathcal{RB}_n)$ is the vector of risk budgets.
- $\mathcal{RC}_i(x)$ is the risk contribution of Asset i with respect to Portfolio x .
- $\tau(x|\tilde{x}) = \sum_{i=1}^n |x_i - \tilde{x}_i|$ is the turnover between Portfolios x and \tilde{x} . The maximum acceptable turnover is denoted by τ^+ .
- $\mathbf{c}(x|\tilde{x})$ is the cost function when rebalancing Portfolio x from Portfolio \tilde{x} . The maximum acceptable cost is denoted by \mathbf{c}^+ .
- $\mathcal{AS}(x|b) = \frac{1}{2} \sum_{i=1}^n |x_i - b_i|$ is the active share of Portfolio x with respect to Benchmark b . \mathcal{AS}^- is the minimum acceptable active share.
- $\mathcal{H}(x) = \sum_{i=1}^n x_i^2$ is the Herfindahl index.
- $\mathcal{N}(x) = 1/\mathcal{H}(x)$ is the number of effective bets. \mathcal{N}^- corresponds to the minimum acceptable number of effective bets.
- $\mathcal{DR}(x) = (x^\top \sigma) / \sqrt{x^\top \Sigma x}$ is the diversification ratio of Portfolio x .
- $\mathcal{LS}(x) = |\sum_{i=1}^n x_i|$ is the long/short exposure of Portfolio x .
- $\mathcal{L}(x) = \sum_{i=1}^n |x_i|$ is the leverage of Portfolio x .
- $\text{SE}(x) = -\sum_{i=1}^n x_i \ln x_i$ is Shannon's entropy of x .
- $\text{KL}(x) = \sum_{i=1}^n x_i \ln(x_i/\tilde{x}_i)$ is the Kullback-Leibler divergence between x and \tilde{x} .
- $W(x)$ is the Lambert W function satisfying $W(x)e^{W(x)} = x$.
- $\mathbf{0}_n$ is the vector of zeros.
- $\mathbf{1}_n$ is the vector of ones.
- e_i is the unit vector, i.e. $[e_i]_i = 1$ and $[e_i]_j = 0$ for all $j \neq i$.
- $x_- = \max(-x, 0) = -\min(x, 0)$ is the negative part of x .
- $x_+ = \max(x, 0)$ is the positive part of x .
- $\mathbb{1}_\Omega(x)$ is the convex indicator function of Ω : $\mathbb{1}_\Omega(x) = 0$ for $x \in \Omega$ and $\mathbb{1}_\Omega(x) = +\infty$ for $x \notin \Omega$.

- A^\dagger is the Moore-Penrose pseudo-inverse matrix of A .
- $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ is the ℓ_p norm.
- $\|x\|_A = (x^\top A x)^{1/2}$ is the weighted ℓ_2 norm.
- $x \odot y$ is the Hadamard element-wise product: $[x \odot y]_{i,j} = [x]_{i,j} [y]_{i,j}$.
- $\mathbf{prox}_f(v)$ is the proximal operator of $f(x)$: $\mathbf{prox}_f(v) = \arg \min_x \left\{ f(x) + \frac{1}{2} \|x - v\|_2^2 \right\}$.
- $\mathcal{S}(v; \lambda) = \text{sign}(v) \cdot (|v| - \lambda)_+$ is the soft-thresholding operator.
- $\mathcal{S}(v; \lambda_-, \lambda_+) = (v - \lambda_+)_+ - (v + \lambda_-)_-$ is the two-sided soft-thresholding operator. We have the following property: $\mathcal{S}(v; \lambda) = \mathcal{S}(v; \lambda, \lambda)$.
- $\mathcal{T}(v, x^-, x^+) = \max(x^-, \min(x, x^+))$ is the truncation operator.
- $\mathcal{P}_\Omega(v)$ is the projection of v onto the set Ω : $\mathcal{P}_\Omega(v) = \arg \min_{x \in \Omega} \frac{1}{2} \|x - v\|_2^2 = \mathbf{prox}_{\mathbb{1}_\Omega(x)}(v)$.
- \mathbb{S}_n is the unit simplex with dimension n .
- $\mathcal{A}_{\text{ffineset}}[A, B]$ is the affine set $\{x \in \mathbb{R}^n : Ax = B\}$.
- $\mathcal{H}_{\text{hyperplane}}[a, b]$ is the hyperplane $\{x \in \mathbb{R}^n : a^\top x = b\}$.
- $\mathcal{H}_{\text{halfspace}}[c, d]$ is the half-space $\{x \in \mathbb{R}^n : c^\top x \leq d\}$.
- $\mathcal{B}_{\text{ox}}[x^-, x^+]$ is the box $\{x \in \mathbb{R}^n : x^- \leq x \leq x^+\}$.
- $\mathcal{B}_p(c, r)$ is the ℓ_p -ball $\{x \in \mathbb{R}^n : \|x - c\|_p \leq r\}$.
- \mathfrak{D} is the weight diversification set $\{x \in \mathbb{R}^n : \mathcal{D}(x) \geq \mathcal{D}^-\}$ where $\mathcal{D}(x)$ is the diversification measure and \mathcal{D}^- is the minimum acceptable diversification.