

Optimization for large scale machine learning

Mathurin Massias
mathurin.massias@gmail.com

Last updated: October 1, 2024

Contents

0	Reminders	5
0.1	Calculus	5
0.2	Linear algebra	6
0.3	Exercises	9
1	Motivation: gradient descent on ordinary least squares	10
1.1	Why does statistical learning need optimization?	10
1.2	Solving least squares with gradient descent	13
1.3	Numerical puzzles	14
1.4	Exercises	15
2	Reminder on convex analysis	17
2.1	Convexity and minimizers	17
2.2	Exercises	18
2.2.1	Convexity	18
2.2.2	Gradient	19
3	Gradient descent on convex functions	20
3.1	Convergence rates for Lipschitz convex functions	20
3.2	Convergence rates for convex L -smooth functions	21
3.3	Convergence rate for smooth and strongly convex functions	23
3.4	Exercises	25
3.4.1	Convexity inequalities	25
3.4.2	Gradient descent	25
3.5	Appendix	26
4	Line search	27
5	Inertial acceleration: Nesterov and Heavy-ball	28
5.1	Polyak's heavy ball method	28
5.2	Nesterov acceleration	29

6	Lower bounds for first-order methods	33
7	Non-smooth convex optimization	36
7.1	Constrained optimization	36
7.2	Extended value functions	37
7.3	Subdifferential of convex functions	38
7.4	The proximal operator	39
7.5	The proximal point algorithm	41
7.6	The proximal gradient descent algorithm	41
7.7	Exercises	43
7.7.1	Subdifferential	43
7.7.2	Constrained optimization	43
8	Duality	45
8.1	Motivation	45
8.2	The Fenchel transform	46
8.3	Fenchel-Rockafellar duality	50
8.4	A calculus rule for subdifferentials	53
8.5	Primal-dual algorithms	53
8.6	Lagrangian duality and KKT conditions	55
8.7	Exercises	58
9	Convergence through fixed-point iterations	60
9.1	Picard iterations	60
9.2	Monotone operators	63
9.3	Applications	64
9.4	Exercises	64
10	Second order optimization: Newton method	65
10.1	Newton method	65
10.2	Quasi Newton methods	67
10.2.1	The Broyden/SR1 formula	67
10.2.2	BFGS	68
10.2.3	DFP	72
10.3	Exercises	72
11	The finite sum structure: stochastic methods	73
11.1	Stochastic gradient descent	73
11.2	Variance reduction	76
12	Mirror descent	79
12.1	The mirror descent algorithm	81
12.2	One proof to rule them all	83
12.3	Online learning and mirror descent	84

13 Continuous view: gradient flows, mirror flows, and Nesterov	86
13.1 Gradient flow	86
13.2 Mirror descent flow	88
13.3 An ODE modeling Nesterov acceleration	89
13.4 Subgradient flow	90
13.5 Exercises	90
14 Bilevel optimization	92
14.1 The implicit approach	92
14.2 The value function approach	94
14.3 See also	94
15 Implicit regularization	95
15.1 Implicit bias of GD on least squares	95
15.2 “Implicit” bias of mirror descent	95
15.3 Logistic regression on separable data	97
15.4 Matrix factorization: can everything be explained in terms of norms?	97
16 Automatic differentiation	98
16.1 Forward and backward automatic differentiation	98
17 Variational inequalities	99
17.1 Solving VIs: the extragradient method	101
18 Additional bibliographic resources	103
A The singular value decomposition	106
B A coercive l.s.c. function without minimizer	107

Goals of the class

In this class, we will study the resolution of optimization problems arising in Machine Learning. These problems are usually characterized by their large scale, which calls for dedicated classes of algorithms. They also usually present some specific structure, that can be leveraged.

The main objectives are to:

1. develop a taxonomy of optimization problems and which algorithms can be used in which settings
2. understand convergence properties and functions properties that govern them
3. master theoretical tools and techniques to derive convergence proofs
4. get hands-on practice to challenge theoretical results, in Python

I thank Cesare Molinari, Saverio Salzo and Silvia Villa for their graduate course on convex analysis.

General notational conventions

We almost always work in finite dimension and the optimization is performed over the set \mathbb{R}^d . In a Machine Learning context, the integer d will thus denote the dimension of the parameter space, while n will denote the number of samples (observed data points). We adopt the optimization convention, i.e. we write a linear system as $Ax = b$ (while statistics has $y = X\beta$, the inverse problems field writes $f = Au$, signal processing writes $y = \Phi x \dots$): the unknown is x ; the iterates are x_k or x_t . As much as possible we write x^* for the minimizer of given problems, but due to the various meanings of the star symbol (Fenchel conjugation, adjoint), it happens that we also write \hat{x} as in the statistical literature.

These conventions are meant to ease reading, but there can always be exceptions. It is anyway a good exercise for your mathematical dexterity to adapt to different notation.

0 Reminders

0.1 Calculus

Definition 0.1 (Gradient vector). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable. Then for any $x \in \mathbb{R}^d$, there exists a vector of \mathbb{R}^d , denoted $\nabla f(x)$ and called the gradient of f at x , such that for all $h \in \mathbb{R}^d$,*

$$f(x+h) = f(x) + \langle \nabla f(x), h \rangle + o(\|h\|) \quad (0.1)$$

This vector is unique.

This is a linear approximation at f , generalizing the well-known Taylor formula in 1D.

Definition 0.2 (Hessian matrix). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be twice differentiable. Then for any $x \in \mathbb{R}^d$, there exists a matrix in $\mathbb{R}^{d \times d}$, denoted $\nabla^2 f(x)$ and called the Hessian of f at x , such that for all $h \in \mathbb{R}^d$,*

$$f(x+h) = f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} h^\top \nabla^2 f(x) h + o(\|h\|^2) \quad (0.2)$$

As the gradient, if the Hessian exists it is uniquely defined.

Much like Equation (0.1), Equation (0.2) is a local approximation of f , but this time *quadratic* (Taylor of order 2). The uniqueness property is very useful to compute gradients and Hessians: if $f(x+h) = f(x) + \langle v, h \rangle + o(\|h\|)$, then v must be equal to $\nabla f(x)$. Hence a technique to compute $\nabla f(x)$ is to write $f(x+h)$, try to isolate $f(x)$ and a linear term, and then identify the slope of the linear term with $\nabla f(x)$.

When f is vector-valued, there is a generalization of the gradient, called the Jacobian.

Definition 0.3 (Jacobian). *The Jacobian of $f : x \in \mathbb{R}^n \mapsto \begin{pmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{pmatrix}$ is the matrix of partial derivatives:*

$$\left(\frac{\partial f_i}{\partial x_j}(x) \right)_{ij} \in \mathbb{R}^{m \times n} . \quad (0.3)$$

For real-valued functions, the Jacobian is the *transposed* gradient. Else, the Jacobian consists of stacked transposed gradients:

$$\mathcal{J}_f(x) = \begin{pmatrix} -\nabla f_1(x)^\top \\ \vdots \\ -\nabla f_n(x)^\top \end{pmatrix} \quad (0.4)$$

The Jacobian of the gradient is the Hessian (transposed, equal if f twice differentiable by Clairaut/Schwarz theorem; holds if second partial derivatives are continuous or more weakly if partial derivatives are continuous).

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(x) \\ \vdots \\ \frac{\partial f}{\partial x_n}(x) \end{pmatrix}$$

$$\mathcal{J}_{\nabla f}(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1}(x) & \dots & \frac{\partial^2 f}{\partial x_n \partial x_1}(x) \\ \vdots & & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) & \dots & \frac{\partial^2 f}{\partial x_n \partial x_n}(x) \end{pmatrix}$$

Proposition 0.4 (Chain rule). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^n$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^p$ be differentiable functions. Then their composition $g \circ f$ is differentiable and*

$$\mathcal{J}_{g \circ f}(x) = \mathcal{J}_g(f(x)) \mathcal{J}_f(x) . \quad (0.5)$$

0.2 Linear algebra

Theorem 0.5 (Spectral theorem). *Let $M \in \mathbb{R}^{d \times d}$ be a symmetric real matrix. Then it can be diagonalized in an orthonormal basis: there exists $U \in \mathbb{R}^{d \times d}$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ such that $UU^\top = U^\top U = \text{Id}_d$ and*

$$M = U^\top \Lambda U \quad (0.6)$$

This is the most important theorem in this section. It is very convenient to diagonalize a matrix in order to raise it to some power $M^k = U^\top \Lambda^k U$, and to simplify computations: up to a change of basis, a symmetric real matrix can be considered diagonal.

Definition 0.6. A **symmetric** matrix M is said to be **positive semidefinite** (PSD, written $M \succeq 0$) if for all $x \in \mathbb{R}^d$, $x^\top M x \geq 0$. If the inequality holds strictly for all nonzero x , M is said to be **positive definite** (PD, written $M \succ 0$).

Note that by definition, a PSD matrix must be symmetric; also notice that for a generic A , the antisymmetric part of A has no influence of $x^\top A x$, so it's not really a restriction.

Example 0.7. The typical way to create a PSD matrix is to pick $M = A^\top A$ for A of any shape, because then $x^\top M x = x^\top A^\top A x = \|Ax\|^2 \geq 0$. By this formulation, we deduce that $A^\top A$ is PD if and only if A is injective.

Proposition 0.8 (Characterization of PSD matrices). *A matrix $M \in \mathbb{R}^{d \times d}$ is PSD if and only if all its eigenvalues are positive.*

Definition 0.9 (Loewner order). *The cone of PSD matrices induces a partial order on symmetric matrices: for $A, B \in \mathbb{R}^{d \times d}$ symmetric, we write $A \succeq B$ if and only if $A - B \succeq 0$. This order is not total (find an example of two matrices which cannot be compared).*

Apart from SPD matrices, another class of very useful matrices are rank one matrices.

Proposition 0.10. *The matrix $M \in \mathbb{R}^{d \times d}$ is rank one if and only if it writes $M = xy^\top$ for $x, y \in \mathbb{R}^d$, both nonzero. If M is both symmetric and rank one, then it writes xx^\top or $-xx^\top$ for $x \in \mathbb{R}^d$.*

Remark 0.11 (Writing a matrix as sum of rank one matrices). *Let $M \in \mathbb{R}^{n \times p}$, $N \in \mathbb{R}^{p \times q}$. Then $MN = \sum_{i=1}^p M_{:i} N_{i:}$.*

The eigenvalue decomposition is a very powerful tool, but not all matrices can be diagonalized. The singular value decomposition (SVD) is a more versatile tool, generalizing the eigenvalue decomposition, that applies to any matrix (even non-square ones).

Proposition 0.12 (Singular value decomposition). *Let $A \in \mathbb{R}^{n \times m}$, let $r = \text{rank}(A)$. Then there exist an $n \times n$ orthogonal matrix U , an $m \times m$ orthogonal matrix V and a $m \times n$ diagonal matrix Σ , with r nonzero decreasing diagonal entries $\sigma_1 \geq \dots \geq \sigma_r > 0$, such that:*

$$A = U\Sigma V^\top. \quad (0.7)$$

The σ_i 's are called the singular values of A . The columns of U (resp. V) are denoted u_i (resp. v_i) and are called the left (resp. right) singular vectors of A . The SVD is often written under the more convenient form

$$A = \sum_{i=1}^r \sigma_i u_i v_i^\top. \quad (0.8)$$

Since Σ only has r nonzero entries, one often writes the SVD as $A = U\Sigma V^\top$ with $\Sigma \in \mathbb{R}^{r \times r}$, $U \in \mathbb{R}^{n \times r}$, $V \in \mathbb{R}^{m \times r}$. Note that in that case U and V are no longer orthogonal: $U^\top U = \text{Id}_r$ but $UU^\top \neq \text{Id}_n$. See [Figure 0.1](#) for an illustration of the various ways to write the SVD.

See [Appendix A](#) for the proof of [Proposition 0.12](#) and more details about the singular value decomposition.

The SVD allows generalizing the notion of inverse to noninvertible and non square matrices, through the Moore-Penrose pseudoinverse.

Definition 0.13. *Let $A \in \mathbb{R}^{n \times d}$, admitting as SVD $A = \sum_{i=1}^r \sigma_i u_i v_i^\top$. Its Moore-Penrose pseudoinverse is, by definition,*

$$A^\dagger = \sum_{i=1}^r \frac{1}{\sigma_i} v_i u_i^\top \in \mathbb{R}^{d \times n}. \quad (0.9)$$

If A is invertible, then its inverse and its Moore-Penrose pseudoinverse coincide.

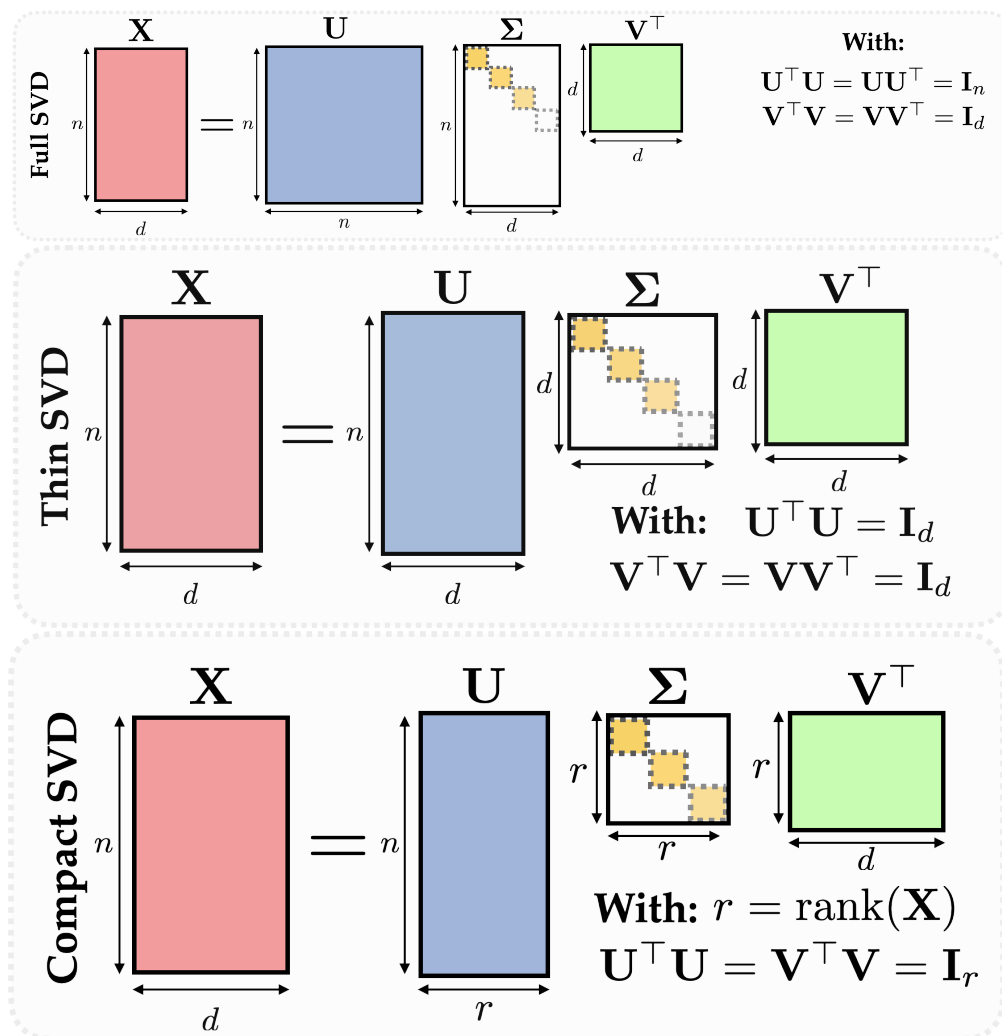


Figure 0.1: 3 possible ways to write the SVD of $X \in \mathbb{R}^{n \times d}$, in the $n \geq d$ setting. Top: with full matrices U and V . Middle: with square Σ of size $d \times d$ (with potentially vanishing entries). Bottom: with square Σ of size $r \times r$ (only nonzero entries). Picture courtesy of Titouan Vayer.

0.3 Exercises

Exercise 0.1. ☹ What are the eigenvalues of a rank one matrix? Of a symmetric rank one matrix?

Exercise 0.2 (Power method). ☹☹ Let $M \in \mathbb{R}^{n \times d}$. Consider the following iterations:

$$\begin{aligned} a_{k+1} &= Mb_k / \|Mb_k\| \\ b_{k+1} &= M^\top a_{k+1} / \|M^\top a_{k+1}\| \end{aligned}$$

where v_0 is initialized as a random Gaussian vector. What is the cost of one iteration of this method?

Show that (a_k) converges to “the” leading left singular vector of M .

What do (b_k) and (Ma_k) converge to?

Exercise 0.3. ☹☹ Let $A \in \mathbb{R}^{n \times d}$, let $y \in \mathbb{R}^n$. Show that $(A^\top A + \lambda \text{Id}_d)^{-1} \rightarrow_{\lambda \rightarrow 0+} A^\dagger y$.

Exercise 0.4. ☹☹ Provide an example of a setting where the gradient is not equal to the vector of partial derivatives.

Exercise 0.5. ☹ Show that the gradient of a function is orthogonal to the level lines of that function.

Exercise 0.6. ☹ For a differentiable function f , compute the gradient and Hessian of f^2 .

Exercise 0.7. ☹ Compute the gradients and Hessians of $x \mapsto \|x\|^2$, $x \mapsto \|x\|$, $x \mapsto a^\top x$. Is it true that the gradient of $x \mapsto \frac{1}{2}x^\top Ax$ is equal to Ax ?

Exercise 0.8. ☹ Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ be differentiable. Show that $\nabla(g \circ f)(x) = g'(f(x))\nabla f(x)$.

Exercise 0.9. ☹ Let $A \in \mathbb{R}^{n \times d}$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ defined as $g(x) = f(Ax)$ for all $x \in \mathbb{R}^d$. Show that

$$\begin{aligned} \nabla g(x) &= A^\top \nabla f(Ax) , \\ \nabla^2 g(x) &= A^\top \nabla^2 f(Ax) A . \end{aligned}$$

Exercise 0.10. ☹☹ Compute the gradient of the logdet function, $M \mapsto \log \det(M)$, defined over \mathbb{S}_{++}^n .

1 Motivation: gradient descent on ordinary least squares

Most optimization problems studied in this class are of the following form:

$$\inf_{x \in \mathcal{C}} f(x) \quad , \quad (1.1)$$

where \mathcal{C} is a subset of \mathbb{R}^d and f is a real-valued function on a subset of \mathbb{R}^d .

Formally, solving [Problem \(1.1\)](#) means to find the largest lower bound on all values $f(x)$ when x describes \mathcal{C} . This largest lower bound, called the *infimum*, may not exist, and if it exists it may not necessarily be attained, even for nicely-behaved f . The reader should mind that most of the time, we write $\min_{\mathcal{C}} f$ instead of $\inf_{\mathcal{C}} f$; this is an abuse of notation and does not mean that the infimum is attained (not even that it is finite).

In Machine Learning, rather than the smallest value taken by f , we are more interested in finding a minimizer, that is

$$x^* \in \operatorname{argmin}_{x \in \mathcal{C}} f(x) \quad , \quad (1.2)$$

meaning in finding $x^* \in \mathcal{C}$ such that for all $x \in \mathcal{C}$, $f(x^*) \leq f(x)$. In Machine Learning, x^* usually represents the parameters of a model, that once computed are used to achieve some statistical task.

1.1 Why does statistical learning need optimization?

Let $n \in \mathbb{N}$. Suppose that we have collected a dataset $(a_i, b_i)_{i \in [n]}$ of observation-target pairs, and we postulate that they are roughly linearly related, i.e. there exists $x \in \mathbb{R}^d$ such that $b_i \approx a_i^\top x$. This “ \approx ” is very vague. To make it more precise and rigorous, it is frequent to model the link as

$$\forall i \in [n], b_i = x_{\text{true}}^\top a_i + \varepsilon_i \quad (1.3)$$

with $(\varepsilon_i)_{i \in [n]}$ n independent realizations of a centered Gaussian variable of variance σ^2 .

Given the collected observations, how can we infer x_{true} ? From a statistical point of view, an answer is to maximize the conditional likelihood. Under model [\(1.3\)](#), the conditional likelihood is:

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\|a_i^\top x - b_i\|^2}{2\sigma^2}\right) \quad . \quad (1.4)$$

Minimization of the negative log-likelihood in x gives

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^d} \sum_{i=1}^n (b_i - a_i^\top x)^2 = \operatorname{argmin} \frac{1}{2} \|Ax - b\|^2 \quad , \quad (1.5)$$

where the **design matrix** $A \in \mathbb{R}^{n \times d}$ has i -th row a_i , and the **target vector** $b \in \mathbb{R}^n$ has i -th entry b_i . [Problem \(1.5\)](#) is called Ordinary Least Squares (OLS). Solving it gives an estimator for x_{true} and allows, for example, predicting a value b_{n+1} if one observes a new input point a_{n+1} . This is also an illustration that the value of the infimum is not always very interesting (if $d \geq n$ and A has rank n it is 0), contrary to the minimizer.

How to compute the solution of [Problem \(1.5\)](#)? This is a classical problem in linear algebra and there exist direct ways to do it ([Exercise 1.2](#)). But in other settings there may

not exist a closed-form solution: it is perfectly possible to postulate a different generative model than (1.3), or simply not impose a generative model, and find x by minimizing another cost/loss function. This is the setting of empirical risk minimization (ERM): one finds the parameters of a model by minimizing a cost function, often –but not always– arising from likelihood maximization.

Example 1.1 (ℓ_1 and Huber regression). *Ordinary least squares are strongly influenced by outliers: the squared loss gives a lot of importance to large differences. In the objective function (1.5) a mismatch of 2 between the true observation b_i and the model fit $a_i^\top x$ costs $2^2/2 = 2$, but a mismatch of 10 costs $10^2/2 = 50$. A more robust solution, that is less dominated by large mismatches (or “outliers”), can be found by minimizing the sum of absolute errors:*

$$\operatorname{argmin}_{x \in \mathbb{R}^d} \sum_{i=1}^n |a_i^\top x - b_i| = \operatorname{argmin}_{x \in \mathbb{R}^d} \|Ax - b\|_1 . \quad (1.6)$$

Problem (1.6) is, unfortunately, more complicated to solve, as we shall see. It can be reformulated as a Linear Program, but the cost to solve it scales very badly with d . A hybrid of Problems (1.5) and (1.6) involves the Huber loss (Figure 1.1):

$$\operatorname{argmin}_{x \in \mathbb{R}^d} \sum_{i=1}^n \operatorname{huber}_\rho(a_i^\top x - b_i) , \quad (1.7)$$

with the Huber loss defined as:

$$\operatorname{huber}_\rho(x) = \begin{cases} \frac{x^2}{2\rho} + \frac{\rho}{2} , & \text{if } |x| \leq \rho , \\ |x| , & \text{otherwise.} \end{cases} \quad (1.8)$$

The Huber loss is less sensitive to outliers than the squared ℓ_2 loss, and more optimization friendly than the ℓ_1 loss, but there is no free lunch: it relies on a parameter $\rho \in \mathbb{R}_+$ that needs to be tuned, and there is no obvious and cheap way to choose the best ρ .

For both problems, there is no closed-form solution and we need optimization techniques to approximate the minimizer.

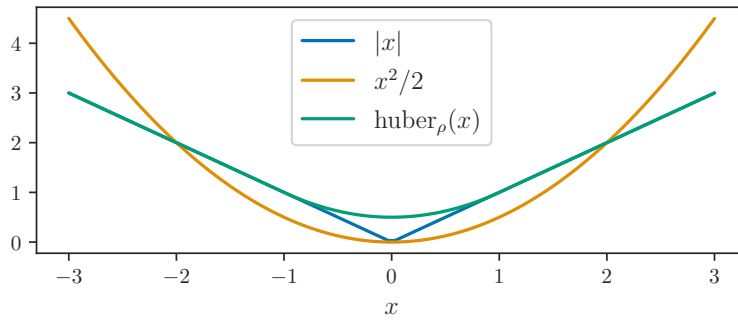


Figure 1.1: ℓ_1 , squared ℓ_2 and Huber ($\rho = 1$) losses

Finally, instead of a linear model, one may model the dependency between observation and target by a function belonging to some parametric class $\{f_x : \mathbb{R}^d \mapsto \mathbb{R} : x \in \mathbb{R}^p\}$, described by p parameters stored in x :

$$b_i \approx f_x(a_i) , \quad (1.9)$$

and, depending on what is meant by \approx , minimize a loss $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, yielding the following Empirical Risk Minimization (ERM) problem:

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^p} \sum_{i=1}^n L(f_x(a_i), b_i) . \quad (1.10)$$

As we have seen, the typical example is $f_x = x^\top \cdot$; but this framework also fits the Deep Learning framework (with f_x a composition of matrix-vector multiplications followed by pointwise application of a non-linear function).

Example 1.2 (Logistic regression). *When working with classification data ($b_i \in \{0, 1\}$), a very popular model is logistic regression¹, which has the following statistical interpretation. Suppose that the law of each b_i is a Bernoulli of parameter $\sigma(x^\top a_i)$, with σ the sigmoid function: $\sigma(x) = \frac{1}{1 + \exp(-x)}$. This model thus assumes that $\mathbb{P}(b_i = 1) = \sigma(x^\top a_i)$. Using $\sigma(-u) = 1 - \sigma(u)$, one immediately has $\mathbb{P}(b_i = 0) = 1 - \mathbb{P}(b_i = 1) = 1 - \sigma(x^\top a_i) = \sigma(-x^\top a_i)$.*

The likelihood of the observations is:

$$\prod_{i:b_i=1} \mathbb{P}(b_i = 1) \prod_{i:b_i=0} \mathbb{P}(b_i = 0) = \prod_{i=1}^n \mathbb{P}(b_i = 1)^{b_i} \mathbb{P}(b_i = 0)^{1-b_i} , \quad (1.11)$$

hence the negative log-likelihood, also known as the binary cross-entropy, is

$$- \sum_{i=1}^n b_i \log \left(\frac{1}{1 + \exp(-a_i^\top x)} \right) + (1 - b_i) \log \left(\frac{1}{1 + \exp(a_i^\top x)} \right) \quad (1.12)$$

$$= - \sum_{i=1}^n b_i \log \left(\frac{1 + \exp(a_i^\top x)}{1 + \exp(-a_i^\top x)} \right) + \log \left(\frac{1}{1 + \exp(a_i^\top x)} \right) \quad (1.13)$$

$$= \sum_{i=1}^n -b_i a_i^\top x + \log(1 + \exp(a_i^\top x)) . \quad (1.14)$$

Minimizing the above function defines the logistic regression estimator.

Remark 1.3. *In the case where the labels are in $\{-1, 1\}$, the negative log likelihood is given by:*

$$\sum_{i:b_i=1} \log(1 + \exp(-x^\top a_i)) + \sum_{i:b_i=-1} \log(1 + \exp(x^\top a_i)) \quad (1.15)$$

$$= \sum_{i=1}^n \log(1 + \exp(-b_i a_i^\top x)) . \quad (1.16)$$

¹despite the name, this is a model for classification and not for regression

1.2 Solving least squares with gradient descent

Problem (1.5) can be solved in closed-form with dedicated algorithms such as Cholesky, LU or QR decompositions (see [Exercise 1.2](#)); but ℓ_1 , Huber and generic ERM models cannot. *We need general-purpose algorithms to solve classes of similar problems.*

One of the arguably simplest minimization algorithms is gradient descent (GD). If it exists, the gradient of a function points towards the direction of maximal increase, since the directional derivative $f'(x; v) = \lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon v) - f(x)}{\epsilon}$ is given by

$$f'(x; v) = \langle \nabla f(x), v \rangle . \quad (1.17)$$

To decrease the value of f as fast as possible, it thus makes sense to move in the opposite direction of the gradient. Taking a stepsize $\eta > 0$, this gives the gradient descent iterations:

$$\begin{aligned} x_0 &\in \mathbb{R}^d , \\ x_{t+1} &= x_t - \eta \nabla f(x_t) . \end{aligned} \quad (1.18)$$

Proposition 1.4 (Convergence rate of GD on OLS). *Consider the OLS problem, and assume that there exist strictly positive scalars μ and L such that $\mu \text{Id} \preceq A^\top A \preceq L \text{Id}$. Let $\kappa = L/\mu$ denote the condition number. Then there exists a unique minimizer of least squares, x^* , and the iterates of gradient descent on ordinary least squares with stepsize $\eta = 1/L$ satisfy:*

$$\|x_t - x^*\| \leq \exp(-t/\kappa) \|x_0 - x^*\| . \quad (1.19)$$

Proof. The gradient of $x \mapsto \frac{1}{2} \|Ax - b\|^2$ is $A^\top (Ax - b)$, hence gradient descent on ordinary least squares read

$$x_{t+1} = x_t - \eta A^\top (Ax_t - b) . \quad (1.20)$$

There exists a unique minimizer x^* , and it satisfies $A^\top A x^* = A^\top b$ ([Exercises 1.1](#) and [1.2](#)), so:

$$x_{t+1} - x^* = (\text{Id} - \eta A^\top A)(x_t - x^*) . \quad (1.21)$$

Observe that

$$(1 - \eta L) \text{Id} \preceq \text{Id} - \eta A^\top A \preceq (1 - \eta \mu) \text{Id} , \quad (1.22)$$

hence

$$\|\text{Id} - \eta A^\top A\|_2 \leq q(\eta) \triangleq \max(|1 - \eta L|, |1 - \eta \mu|) . \quad (1.23)$$

Setting $\eta = 1/L$ yields $q(\eta) = 1 - \mu/L$, and using $(1 - u)^t \leq \exp(-tu)$ concludes. \square

Remark 1.5. *The convergence rate (1.19), is a very good rate in the sense that it decays fast towards 0. It is called a linear rate². Suppose that we want to get ε -close to x^* , then we need at most $\kappa \log(\|x_0 - x^*\|/\varepsilon)$ iterations of GD, which grows quite slowly as ε goes to 0.*

²this is a terrible name: it is linear in log scale

We can also obtain a rate in objective value:

$$f(x_t) - f(x^*) = \langle \nabla f(x^*), x_t - x^* \rangle + \frac{1}{2}(x_t - x^*)^\top A^\top A(x_t - x^*) \quad (1.24)$$

$$= \frac{1}{2}(x_t - x^*)^\top A^\top A(x_t - x^*) \quad (1.25)$$

$$\leq \frac{L}{2} \|x_t - x^*\|^2, \quad (1.26)$$

but as we've seen, such rates are usually less desirable than rates on x_t .

Remark 1.6 (Other choices of η). We have actually proven the stronger result:

$$\|x_t - x^*\| \leq q(\eta)^t \|x_0 - x^*\|, \quad (1.27)$$

with $q(\eta)$ defined in [Equation \(1.23\)](#). This shows that any choice of η in $]0, 2/L[$ leads to convergence, because it results in $q(\eta) < 1$. We can tune η to minimize $q(\eta)$: doing so yields $\eta = \frac{2}{L+\mu}$ and $q(\eta) = \frac{L-\mu}{L+\mu} = \frac{\kappa-1}{\kappa+1}$.

Questions: Does there always exist an L as in the assumptions of [Proposition 1.4](#)? And a μ ? To what do they correspond geometrically?

1.3 Numerical puzzles

On [Figure 1.2](#) is the numerically observed convergence rate of GD on one instance of OLS. The matrix A is 500 by 500 with i.i.d. normal entries so it is full rank, $\mu > 0$ and [Proposition 1.4](#) applies. Yet the rate we see is clearly not linear. Why?

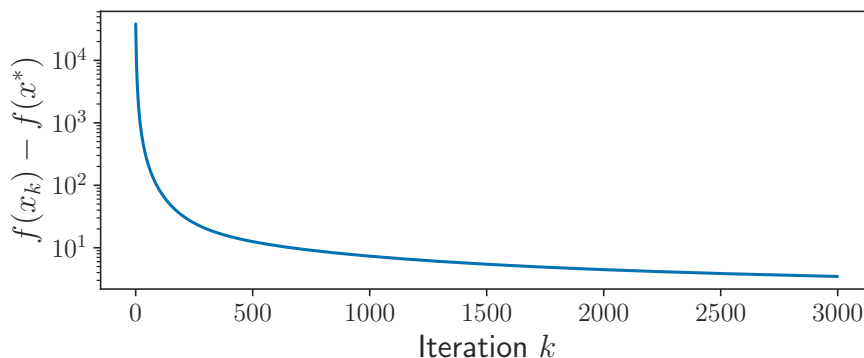


Figure 1.2: Objective convergence rate of GD on OLS on a random i.i.d. $A \in \mathbb{R}^{500 \times 500}$.

On [Figure 1.3](#), A is still random but its shape is 200×300 , so $A^\top A$ is not full rank, $\mu = 0$. Yet we see a clear linear convergence rate (numerical floating point errors kick in at iteration 1000). Why?

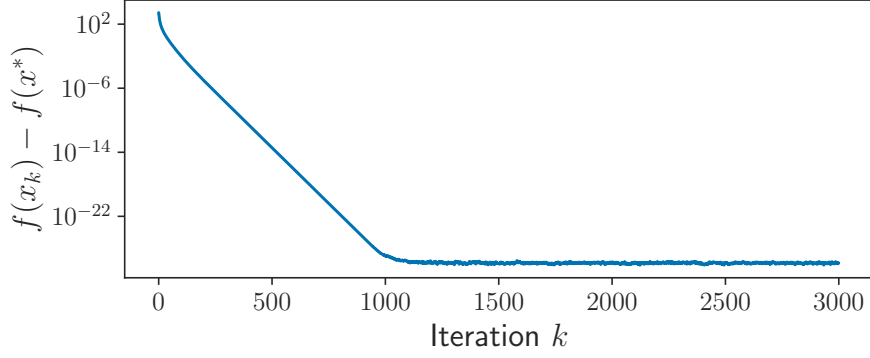


Figure 1.3: Objective convergence rate of GD on OLS on a random i.i.d. $A \in \mathbb{R}^{200 \times 300}$.

1.4 Exercises

Exercise 1.1. ☹️ Let $A \in \mathbb{R}^{n \times d}$. Show that $\text{Ker } A = \text{Ker } A^*A$.

Show that for any $b \in \mathbb{R}^n$ there exist a solution to $A^*Ax = A^*b$.

☹️☹️ Show that there does not always exist a solution to $A^*Ax = A^*b$ in the infinite dimensional case (when A is a bounded linear operator between infinite dimensional Hilbert spaces.)

Exercise 1.2. ☹️ Let $A \in \mathbb{R}^{n \times d}$, $b \in \mathbb{R}^n$. Show that solving Ordinary Least Squares:

$$\min \frac{1}{2} \|Ax - b\|^2, \quad (1.28)$$

amounts to solving $A^*Ax = A^*b$ (aka the normal equations).

Show that the set of solutions is:

$$A^\dagger b + \text{Ker } A.$$

Exercise 1.3. ☹️ When is $x \mapsto \|Ax - b\|^2$ strictly convex? Strongly convex?

Exercise 1.4 (Least squares with intercept). ☹️☹️ An intercept x_0 is a constant scalar term in the linear prediction function, that becomes $a \mapsto a^\top x + x_0$. Fitting an intercept can be done by adding a column of 1s to A . Alternatively, show that the solution of least squares with intercept,

$$(\hat{x}, \hat{x}_0) \in \underset{x \in \mathbb{R}^d, x_0 \in \mathbb{R}}{\text{argmin}} \frac{1}{2} \|Ax - b - x_0 \mathbf{1}\|^2 \quad (1.29)$$

is given by:

$$\hat{x} = \hat{x}_c, \quad (1.30)$$

$$\hat{x}_0 = \frac{1}{n} \sum_{i=1}^n (a_i^\top \hat{x} - b_i), \quad (1.31)$$

where \hat{x}_c is the solution of least squares without intercept on centered data $A_c = A - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top A$ and $b_c = b - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top b$ (versions of A and b where the mean per column has been subtracted).

Exercise 1.5 (Gradient descent on isotropic parabola). ☕ Let $A \in \mathbb{R}^{n \times d}$ be such that the condition number³ of $A^\top A$ is equal to 1. Show that gradient descent with stepsize $1/\|A^\top A\|_{\text{op}}^2$ converges in a single iteration for the problem $\min \frac{1}{2}\|Ax - b\|^2$.

³i.e. the ratio between the largest and the smallest eigenvalues of $A^\top A$.

2 Reminder on convex analysis

[Section 1](#) has shown us basic notions in optimization. We have studied gradient descent, an iterative algorithm producing a sequence of iterates x_k that aims at solving an optimization problem. For the Ordinary Least Squares objective function, we have proved two types of convergence results⁴:

- in iterates: $\|x_k - x^*\| \rightarrow 0$
- in function values: $f(x_k) - \inf f \rightarrow 0$.

Both results were *non-asymptotic*: we had information about the speed at which convergence happened.

In the sequel, we want to minimize functions beyond least squares. For starters, we'll work with convex functions, because they're roughly the only ones we can hope to minimize globally.

2.1 Convexity and minimizers

Definition 2.1 (Global and local minimizers). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$. A global minimizer of f is a point x^* such that $f(x^*) \leq f(x)$ for all $x \in \mathbb{R}^d$. A local minimizer of f is a point x^* such that there exists a neighborhood V of x^* such that $f(x^*) \leq f(x)$ for all $x \in V$.*

Definition 2.2 (Convexity). *A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if it lies below its cords:*

$$\forall x, y \in \mathbb{R}^d, \forall \lambda \in]0, 1[, \quad f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) . \quad (2.1)$$

It is strictly convex if the inequality (2.1) holds strictly.

Remark 2.3. *Equation (2.1) trivially holds for $\lambda = 0$ and $\lambda = 1$, and so we could equivalently define it for $\lambda \in [0, 1]$. However, we will later work with functions that can take the value $+\infty$, and in that case, we could end up with $\lambda f(x) = 0 \times (+\infty)$ which is not defined. So it is more rigorous to require the convexity inequality to hold only for $\lambda \in]0, 1[$.*

Definition 2.4 (Strong convexity). *Let $\mu > 0$ and $\|\cdot\|$ be a norm. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex with respect to $\|\cdot\|$ if for all $x, y \in \mathbb{R}^d$ and $\lambda \in]0, 1[$,*

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \frac{\mu}{2}\lambda(1 - \lambda)\|x - y\|^2 . \quad (2.2)$$

Note: this is the only true definition of strong convexity. Other definitions are rather characterizations in special cases. In particular, many other definitions assume that the norm is the Euclidean one, a requirement that this definition does not have. This will play an important role in Mirror Descent ([Section 12](#)).

When $\|\cdot\|$ is the Euclidean norm, then f is μ -strongly convex if and only if $f - \frac{\mu}{2}\|\cdot\|^2$ is convex ([Exercise 2.11](#)). In this sense, a strongly convex function is so convex that when you subtract a parabola with positive curvature, it remains convex.

⁴Does one imply the other? Under which condition?

Convex functions are amenable to optimization because of the nice “local to global” properties they enjoy.

Proposition 2.5 (Local minimizers are global minimizers). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function. Then all local minimizers of f are also global.*

Proposition 2.6 (Above its tangents). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex differentiable function. Then*

$$\forall x, y \in \mathbb{R}^d, \quad f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle .$$

Written as $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$, it tells that f lies above all of its tangents. This property is sometimes called “local to global”: from local information $\nabla f(x)$, we get information about the global behavior of f .

Convex differentiable functions are nice because it is easy to characterize their minimizers.

Proposition 2.7 (Optimality condition for convex differentiable functions). *Let f be a convex differentiable function. Then*

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \Leftrightarrow \nabla f(x^*) = 0 . \quad (2.3)$$

Equipped with this theoretical framework, we can move to our first convergence proofs on generic functions.

2.2 Exercises

2.2.1 Convexity

Exercise 2.1. ☛ *Show that local minimizers of convex functions are global minimizers.*

Exercise 2.2 (Pointwise supremum preserves convexity). ☛ *Let $(f_i)_I$ be a family of convex functions (not necessarily countable). Show that $x \mapsto \sup_{i \in I} f_i(x)$ is convex.*

Exercise 2.3 (Precomposition by linear operator preserves convexity). ☛ *Let $f : \mathcal{Y} \rightarrow \mathbb{R}$ be a convex function and $A : \mathcal{X} \rightarrow \mathcal{Y}$ a linear operator. Show that $f(A \cdot)$ is convex (on \mathcal{X}).*

Exercise 2.4 (Misconceptions on existence of minimizers). ☛ *Provide an example of convex function which does not admit a minimizer.*

What if the function is continuous and lower bounded?

Exercise 2.5. ☛ *Show that a strictly convex function has at most one minimizer.*

Exercise 2.6 (Jensen’s inequality). ☛ *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex. Let $n \in \mathbb{N}$, $x_1, \dots, x_n \in \mathbb{R}^d$, and let $\lambda_1, \dots, \lambda_n$ be positive scalars summing to 1. Show that $f(\sum_{i=1}^n \lambda_i x_i) \leq \sum_{i=1}^n \lambda_i f(x_i)$.*

Exercise 2.7. ☛ *Show that the sublevel sets of a convex function are convex. Find a function with convex sublevel sets which is not convex.*

Exercise 2.8. ☹☹ A function is said to be midpoint convex if for all x, y ,

$$f\left(\frac{x+y}{2}\right) \leq \frac{1}{2}f(x) + \frac{1}{2}f(y). \quad (2.4)$$

Show that a continuous midpoint convex function is convex.

What happens if you remove the continuity assumption? (this is ☹☹☹)

Exercise 2.9 (First order characterization of convex functions). ☹☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable function. Show that the following are equivalent:

1. f is convex
2. f lies above its tangents: $\forall x, y \in \mathbb{R}^d, f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle$
3. ∇f is monotone: $\forall x, y \in \mathbb{R}^d, \langle \nabla f(x) - \nabla f(y), x - y \rangle \geq 0$

Exercise 2.10 (Continuity). ☹☹☹ Show that a convex function is locally Lipschitz (hence continuous) on the interior of its domain.

Exercise 2.11 (Characterization of strongly convex functions in the Euclidean case). ☹ Show that f is μ -strongly convex with respect to the Euclidean norm if and only if $f - \frac{\mu}{2}\|\cdot\|^2$ is convex.

☹☹ Show that this does not always hold if the norm is not Euclidean.

Exercise 2.12. ☹ Show that a strongly convex function admits exactly one minimizer.

Exercise 2.13 (A convenient equality). Show that, for any $x, y \in \mathbb{R}^d$ and any λ (not necessarily in $[0, 1]$),

$$\|\lambda x + (1 - \lambda)y\|^2 = \lambda\|x\|^2 + (1 - \lambda)\|y\|^2 - \lambda(1 - \lambda)\|x - y\|^2. \quad (2.5)$$

In which sense is it a generalization of the parallelogram identity?

2.2.2 Gradient

Exercise 2.14. ☹☹ Provide an example of a setting where the gradient is not equal to the vector of partial derivatives.

Exercise 2.15. ☹ Show that the gradient of a function is orthogonal to the level lines of that function.

Exercise 2.16. ☹ For a differentiable function f , compute the gradient and Hessian of f^2 .

Exercise 2.17. ☹ Compute the gradients and Hessians of $x \mapsto \|x\|^2$, $x \mapsto \|x\|$, $x \mapsto a^\top x$. Is it true that the gradient of $x \mapsto \frac{1}{2}x^\top Ax$ is equal to Ax ?

Exercise 2.18. ☹ Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ be differentiable. Show that $\nabla(g \circ f)(x) = g'(f(x))\nabla f(x)$.

Exercise 2.19. ☹ Let $A \in \mathbb{R}^{n \times d}$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ defined as $g(x) = f(Ax)$ for all $x \in \mathbb{R}^d$. Show that

$$\begin{aligned} \nabla g(x) &= A^\top \nabla f(Ax) , \\ \nabla^2 g(x) &= A^\top \nabla^2 f(Ax) A . \end{aligned}$$

Exercise 2.20. ☹☹ Compute the gradient of the logdet function, $M \mapsto \log \det(M)$, defined over \mathbb{S}_{++}^n .

3 Gradient descent on convex functions

In [Section 1](#) we proved that for gradient descent on ordinary least squares, we can get very fast (so-called “linear”) convergence rates in terms of function values and iterates. Two numbers governed this convergence: global upper and lower bounds L and μ on the Hessian.

Our proof was very ad hoc, relying on the explicit expression of the minimizer, the gradient and the Hessian. The questions we will try to answer in the sequel are:

- Can we generalize the results of [Proposition 1.4](#) to other objective functions?
- Under which conditions?
- What are the properties that influence the convergence rate we obtain?

3.1 Convergence rates for Lipschitz convex functions

The first and weakest result we’ll prove concerns convex Lipschitz functions.

Remark 3.1. *The following proposition in fact applies to subgradient descent, a more generic algorithm which does not require the function to be differentiable. This notion will be introduced in [Section 7.3](#).*

Proposition 3.2 (Gradient descent on Lipschitz convex functions). *Let f be a convex, L -Lipschitz differentiable function admitting at least one minimizer x^* . For $t \in \mathbb{N}$, the iterates of gradient descent $x_{k+1} = x_k - \eta \nabla f(x_k)$ satisfy*

$$\min_{0 \leq k \leq t-1} f(x_k) - f(x^*) \leq \frac{L \|x_0 - x^*\|}{\sqrt{t}} , \quad (3.1)$$

and

$$f\left(\frac{1}{t} \sum_{k=0}^{t-1} x_k\right) - f(x^*) \leq \frac{L \|x_0 - x^*\|}{\sqrt{t}} , \quad (3.2)$$

if the stepsize η is taken as $\eta = \frac{\|x_0 - x^*\|}{L\sqrt{t}}$.

Let us make a few observations before the proof:

- This algorithm is a bit weird: the total number of iterations t must be known in advance to select the step size (in [Exercise 3.10](#) we get rid of this up to a slight worsening in the rate, using a decaying stepsize $\eta_k \propto 1/\sqrt{k}$).
- The objective values are not necessarily decreasing, it is not a *descent* algorithm.
- The more iterations we do, the smaller we need to take the step size.
- The stepsize depends on the unknown quantity $\|x_0 - x^*\|$. We can get rid of this dependency by bounding $\|x_0 - x^*\|$ with e.g. the diameter of the domain.
- The kind of rate is not as strong as in the OLS case: we know nothing about the last iterate x_t and only have results on the ergodic (averaged) iterates or on the best iterate.

Proof. Let x^* be a minimizer of f . For $k \in \mathbb{N}$, using convexity of f , definition of x_{k+1} and the parallelogram identity,

$$f(x_k) - f(x^*) \leq \langle \nabla f(x_k), x_k - x^* \rangle \quad (3.3)$$

$$\leq \frac{1}{\eta} \langle x_k - x_{k+1}, x_k - x^* \rangle \quad (3.4)$$

$$= \frac{1}{2\eta} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2 + \|x_k - x_{k+1}\|^2) \quad (3.5)$$

$$= \frac{1}{2\eta} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) + \frac{\eta}{2} \|\nabla f(x_k)\|^2. \quad (3.6)$$

Summing, we obtain by telescoping cancellation:

$$\sum_{k=0}^{t-1} [f(x_k) - f(x^*)] \leq \frac{1}{2\eta} (\|x_0 - x^*\|^2 - \|x_t - x^*\|^2) + \frac{\eta}{2} \sum_{k=0}^{t-1} \|\nabla f(x_k)\|^2 \quad (3.7)$$

$$\leq \frac{1}{2\eta} \|x_0 - x^*\|^2 + \frac{\eta}{2} \sum_{k=0}^{t-1} \|\nabla f(x_k)\|^2. \quad (3.8)$$

Since f is L -Lipschitz we can bound $\|\nabla f(x_k)\|$ by L , so:

$$\frac{1}{t} \sum_{k=0}^{t-1} f(x_k) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2t\eta} + \frac{\eta L^2}{2}. \quad (3.9)$$

Minimizing the RHS in η gives $\eta = \frac{R}{L\sqrt{t}}$ and the upper bound has value $\frac{RL}{\sqrt{t}}$. To obtain the rate on the iterate with minimal value of f , observe that $\min_{0 \leq k \leq t-1} f(x_k) - f(x^*) \leq \frac{1}{t} \sum_{k=0}^{t-1} f(x_k) - f(x^*)$. To obtain the rate on the ergodic iterate, apply Jensen's inequality to the LHS. \square

The $1/\sqrt{k}$ rate is quite poor: to approach the optimal value within ε , one needs $\mathcal{O}(1/\varepsilon^2)$ iterations. With more assumptions on f , it can be improved to $\mathcal{O}(1/k)$, as we shall see in the next section.

3.2 Convergence rates for convex L -smooth functions

Definition 3.3 (L -smoothness). For $L \geq 0$, a differentiable function f is said to be L -smooth if its gradient is L -Lipschitz: for all $x, y \in \text{dom } f$,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|. \quad (3.10)$$

A widely used property of L -smooth functions is that they satisfy the so-called Descent lemma.

Lemma 3.4 (Descent lemma). Let f be a L -smooth function. Then for all $x, y \in \text{dom } f$,

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2. \quad (3.11)$$

Think of this inequality for a y fixed, and involving two functions of x , f and $\phi_y = f(y) + \langle \nabla f(y), \cdot - y \rangle + \frac{L}{2} \|\cdot - y\|^2$. The function ϕ_y is a convex, isotropic parabola. Equation (3.11) says that ϕ_y globally upper bounds f ; in addition, the two functions are equal and tangent at y .

Proof. Notice that:

$$f(x) - f(y) = \int_0^1 \frac{d}{dt} f(y + t(x - y)) dt = \int_0^1 \langle \nabla f(y + t(x - y)), x - y \rangle dt . \quad (3.12)$$

Subtract $\langle \nabla f(y), x - y \rangle$, use Cauchy-Schwarz and the hypothesis on f , and conclude. \square

Lemma 3.5 (Cocoercivity of the gradient). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a L -smooth function.*

$$\forall x, y \in \mathbb{R}^d, \quad \langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 . \quad (3.13)$$

In addition to the above lemma, which is the only one we'll use here, there is a whole set of other properties and characterizations of L -smooth functions; some of them are in Exercise 3.1, and the interested reader can refer to the very complete <https://xingyuzhou.org/blog/notes/Lipschitz-gradient>.

Proposition 3.6 (Convergence rate of gradient descent on L -smooth functions). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex L -smooth function with nonempty set of minimizers. Then the iterates of gradient descent with stepsize $1/L$ converge at the rate:*

$$f(x_k) - f(x^*) \leq \frac{2L\|x_0 - x^*\|^2}{k} . \quad (3.14)$$

Remark 3.7. *An alternative proof, unifying several algorithms, is given in Section 12.2. It is more elegant and stronger, as it does not require the existence of a minimizer: it shows that for any reference point x , $f(x_k) - f(x) \leq \|x - x_0\|^2/k$. Set $g = 0$ and $D(x, y) = \frac{1}{2}\|x - y\|^2$ in Section 12.2 to obtain the framework of gradient descent on L -smooth functions.*

Proof. Let x^* be a minimizer of f . First, we show that the distance of x_k to x^* decreases:

$$\|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 + 2\langle x_k - x^*, x_{k+1} - x_k \rangle + \|x_{k+1} - x_k\|^2 \quad (3.15)$$

$$= \|x_k - x^*\|^2 - \frac{2}{L} \langle \nabla f(x_k), x_k - x^* \rangle + \frac{1}{L^2} \|\nabla f(x_k)\|^2 . \quad (3.16)$$

Lemma 3.5, combined with $\nabla f(x^*) = 0$, yields:

$$-\frac{2}{L} \langle \nabla f(x_k), x_k - x^* \rangle + \frac{1}{L^2} \|\nabla f(x_k)\|^2 \leq -\frac{1}{L^2} \|\nabla f(x_k)\|^2 \leq 0 , \quad (3.17)$$

which concludes: $\|x_{k+1} - x^*\| \leq \|x_k - x^*\|$.

We then use the convexity of f :

$$f(x_k) - f(x^*) \leq \langle \nabla f(x_k), x_k - x^* \rangle \leq \|\nabla f(x_k)\| \|x_k - x^*\| \leq \|\nabla f(x_k)\| \|x_0 - x^*\| . \quad (3.18)$$

Finally, the descent Lemma ([Lemma 3.4](#)) quantifies the decrease in objective at each iteration:

$$f(x_{k+1}) - f(x_k) \leq -\frac{1}{2L} \|\nabla f(x_k)\|^2 . \quad (3.19)$$

Introducing $\delta_k = f(x_k) - f(x^*)$, we have:

$$\delta_{k+1} - \delta_k \leq -\frac{1}{2L} \|\nabla f(x_k)\|^2 \quad (3.20)$$

$$\leq -\frac{1}{2L} \frac{\delta_k^2}{\|x_0 - x^*\|^2} . \quad (3.21)$$

The above manipulations are quite standard in optimization, but the following trick is a bit surprising the first time you see it: dividing by $\delta_k \delta_{k+1}$,

$$\frac{1}{\delta_k} - \frac{1}{\delta_{k+1}} \leq -\frac{1}{2L\|x_0 - x^*\|^2} \frac{\delta_k}{\delta_{k+1}} \quad (3.22)$$

$$\leq -\frac{1}{2L\|x_0 - x^*\|^2} , \quad (3.23)$$

since (δ_k) decreases (by (3.20)). Summing and telescoping concludes. \square

If the function is even more regular, the rate of gradient descent can further be improved.

3.3 Convergence rate for smooth and strongly convex functions

Proposition 3.8. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth and μ -strongly convex. Then it admits a unique minimizer x^* ([Exercise 2.12](#)) and the iterates of gradient descent with stepsize $1/L$ converge linearly:*

$$f(x_k) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^k (f(x_0) - f(x^*)) . \quad (3.24)$$

Note that, since $e^x \geq 1 + x$ (by convexity!), so $(1 - \frac{\mu}{L})^k \leq \exp(-\frac{\mu}{L}k)$, and this quantity is often used instead in (3.24). Since $\mu/L \leq 1$ and is usually very close to 0, one does not lose much in this majorization.

Proof. Let x^* be the minimizer of f . Being μ -strongly convex, f satisfies the Polyak-Lojasiewicz inequality ([Exercise 3.5](#)):

$$\forall x \in \mathbb{R}^d, f(x) - f(x^*) \leq \frac{1}{2\mu} \|\nabla f(x)\|^2 . \quad (3.25)$$

Combined with the descent lemma,

$$f(x_{k+1}) - f(x_k) \leq -\frac{1}{2L} \|\nabla f(x_k)\|^2 \quad (3.26)$$

$$f(x_{k+1}) - f(x^*) + f(x^*) - f(x_k) \leq -\frac{\mu}{L} (f(x_k) - f(x^*)) . \quad (3.27)$$

hence $f(x_{k+1}) - f(x^*) \leq (1 - \frac{\mu}{L})(f(x_k) - f(x^*))$. \square

Remark 3.9. Using $f(x_k) - f(x^*) \leq \langle \nabla f(x_k), x_k - x^* \rangle = \langle \nabla f(x_k) - \nabla f(x^*), x_k - x^* \rangle \leq LR^2$ where $R \triangleq \|x_0 - x^*\|$, one obtains:

$$f(x_k) - f(x^*) \leq LR^2 \left(1 - \frac{\mu}{L}\right)^k \quad (3.28)$$

3.4 Exercises

3.4.1 Convexity inequalities

Exercise 3.1. ☹☹ Let f be a convex and differentiable function. Let $L > 0$. Show that the following properties are equivalent:

1. $\forall x, y, \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$
2. $\forall x, y, f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2}\|x - y\|^2$
3. $\forall x, y, \frac{1}{L}\|\nabla f(x) - \nabla f(y)\|^2 \leq \langle x - y, \nabla f(x) - \nabla f(y) \rangle$

Exercise 3.2. ☹ Let f be a twice differentiable L -smooth function. Show that for all $x \in \mathbb{R}^d$, $\nabla^2 f(x) \preceq L\text{Id}$.

Exercise 3.3. ☹ Let f be a differentiable function. Show that the following properties are equivalent:

1. f is μ -strongly convex (with respect to the Euclidean norm)
2. $\forall x, y, f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2}\|x - y\|^2$
3. $\forall x, y, \mu\|x - y\|^2 \leq \langle x - y, \nabla f(x) - \nabla f(y) \rangle$

Exercise 3.4. ☹ Let f be a twice differentiable μ -strongly convex function. Show that for all $x \in \mathbb{R}^d$, $\mu\text{Id} \preceq \nabla^2 f(x)$.

Exercise 3.5 (Polyak-Łojasiewicz inequality). ☹☹ Let f be a μ -strongly-convex and differentiable function. Let $x^* = \text{argmin } f(x)$. Show that f satisfies the Polyak-Łojasiewicz inequality:

$$\mu(f(x) - f(x^*)) \leq \frac{1}{2}\|\nabla f(x)\|^2 .$$

Provide an example of function which is not strongly convex, but satisfies the inequality.

Exercise 3.6. ☹☹ Let f be a L -smooth μ -strongly convex function. Show that for any x, y ,

$$\frac{\mu L}{\mu + L}\|x - y\|^2 + \frac{1}{L + \mu}\|\nabla f(x) - \nabla f(y)\|^2 \leq \langle x - y, \nabla f(x) - \nabla f(y) \rangle .$$

Exercise 3.7 (3 point descent lemma). ☹ Let f be convex and L -smooth. Show that for any triplet (x, y, z) ,

$$f(x) \leq f(y) + \langle \nabla f(z), x - y \rangle + \frac{L}{2}\|x - z\|^2 .$$

3.4.2 Gradient descent

Exercise 3.8 (Exercise 13.2 in discrete time). ☹☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex, twice differentiable L -smooth function.

Show that the iterates of gradient descent on f with step size $0 < \alpha < 2/L$ have decreasing gradient norm.

Can you show it if f is not twice differentiable?

Is it still true when f is not convex?

Exercise 3.9. ☹ Provide a finite-dimensional example of convex L -smooth function f such that gradient descent with stepsize $< 2/L$ diverges.

Exercise 3.10 (Gradient descent on Lipschitz function without knowing the horizon). ☹☹ For gradient descent on a Lipschitz differentiable objective, the classical proof assumes that the total number of iterations t is known in advance, to set the fixed stepsize $\eta \propto 1/\sqrt{t}$. Show that using a decreasing stepsize $\eta_k \propto 1/\sqrt{k}$ leads to a rate of order $\log k/\sqrt{k}$ on the ergodic or best iterate.

Exercise 3.11 (Gradient descent is equivariant to rotation). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and Lipschitz differentiable (not required for the exercise, just to be in a setup where gradient descent is relevant). Let U be a rotation matrix ($U^\top U = \text{Id}$) and $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ a reparametrization of f defined as $\phi = f(U \cdot)$.

Show that gradient descent is equivariant to rotation: the iterates of GD of f and ϕ respectively, (x_k) and (y_k) , started at $x_0 = Uy_0$, satisfy $x_k = Uy_k$ for all $k \in \mathbb{N}$. How do the respective losses compare?

Show that studying the convergence of gradient descent on a quadratic can be reduced to studying gradient descent on a separable quadratic (meaning, with diagonal Hessian).

3.5 Appendix

Here is a short, alternate proof technique for the $1/k$ rate of gradient descent on convex L -smooth functions.

Starting from Equation (3.8), we use that, since f is L -smooth, $f(x^*) - f(x) \leq -\frac{1}{2L} \|\nabla f(x)\|^2$ (inequality obtained by minimizing both sides of the descent lemma), meaning that $\|\nabla f(x)\|^2 \leq \frac{L}{2}(f(x) - f(x^*))$. From Equation (3.8) we get:

$$\sum_{k=0}^{t-1} [f(x_k) - f(x^*)] \leq \frac{1}{2\eta} \|x_0 - x^*\|^2 + \eta L \sum_{k=0}^{t-1} (f(x_k) - f(x^*)) \quad (3.29)$$

$$(1 - \eta L) \sum_{k=0}^{t-1} [f(x_k) - f(x^*)] \leq \frac{1}{2\eta} \|x_0 - x^*\|^2 \quad (3.30)$$

Using the descent lemma shows that $f(x_k)$ decreases, and so the LHS is larger than $k(f(x_k) - f(x^*))$.

4 Line search

In all this section we are at iteration k and p_k is a descent direction ($\langle \nabla f(x_k), p_k \rangle < 0$); we want to select the stepsize to make $\phi(\alpha) = f(x_k + \alpha p_k)$ small.

Definition 4.1 (Armijo rule). *The Armijo rule is:*

$$f(x_k + \alpha p_k) \leq f(x_k) + \alpha c_1 \langle \nabla f(x_k), p_k \rangle, \quad (4.1)$$

for $0 < c_1 < 1$ (typically, $c_1 = 10^{-4}$).

Notice that $\phi'(0) = \langle \nabla f(x_k), p_k \rangle < 0$ and that $c_1 < 1$, so the linear part in Armijo's RHS starts above ϕ . Armijo prevents α from being too large.

Definition 4.2 (Curvature condition). *The curvature condition is:*

$$\langle \nabla f(x_k + \alpha p_k), p_k \rangle \geq c_2 \langle \nabla f(x_k), p_k \rangle, \quad (4.2)$$

for $0 < c_2 < 1$ (typically, $c_2 = 0.9$).

The curvature condition ensures that we pick an alpha so that $\phi'(\alpha)$ (the LHS) is not too negative: if it is largely negative, we should take a larger α , as this will decrease ϕ more.

If α satisfies [Equations \(4.1\) and \(4.2\)](#), we say it satisfies the Wolfe conditions.

Remark 4.3. Since $c_2 > 0$ and $\langle \nabla f(x_k), p_k \rangle < 0$, the LHS in [eq. \(4.2\)](#) is negative and enforcing it may result in a bad α . An alternative is to enforce the strong curvature condition:

$$|\langle \nabla f(x_k + \alpha p_k), p_k \rangle| \geq |c_2 \langle \nabla f(x_k), p_k \rangle|. \quad (4.3)$$

If [Equations \(4.1\) and \(4.3\)](#) are satisfied, α is said to satisfy the strong Wolfe conditions.

Together, Armijo and Wolfe rule thus avoid taking steps that are too small, or too large. But can we find an α satisfying both? The following lemma answers positively, under mild assumptions.

Lemma 4.4 (Wolfe's lemma). *Let f be C^1 and lower bounded. Assumes that $0 < c_1 < c_2$. Then there exists an open interval of \mathbb{R}_+ on which the Wolfe conditions are satisfied.*

Proof. Let $g(\alpha) = \phi(\alpha) - f(x_k) - c_1 \alpha \langle \nabla f(x_k), p_k \rangle$. We have $g(0) = 0$ and $g'(0) = (1 - c_1) \langle \nabla f(x_k), p_k \rangle < 0$ (p_k being a descent direction), so since f is C^1 , there exists a (right) neighborhood of 0 where g is negative. Let $\bar{\alpha}$ be the smallest strictly positive 0 of g (which must exist because f is bounded below). Notice that $g \leq 0$ on $[0, \bar{\alpha}]$, meaning that Armijo rule is satisfied on this interval; moreover we have equality at $\bar{\alpha}$: $\phi(\bar{\alpha}) = f(x_k) + c_1 \bar{\alpha} \langle \nabla f(x_k), p_k \rangle$.

Since g is C^1 , there exists $\tilde{\alpha} \in]0, \bar{\alpha}[$ such that:

$$\phi(\bar{\alpha}) - \phi(0) = \bar{\alpha} \phi'(\tilde{\alpha}) \quad (4.4)$$

which rewrites, dividing by $\bar{\alpha}$,

$$c_1 \langle \nabla f(x_k), p_k \rangle = (\phi(\bar{\alpha}) - \phi(0)) / \bar{\alpha} = \phi'(\tilde{\alpha}) = \langle \nabla f(x_k + \tilde{\alpha} p_k), p_k \rangle \quad (4.5)$$

Since $c_1 < c_2$ we have $\langle \nabla f(x_k + \tilde{\alpha} p_k), p_k \rangle > c_2 \langle \nabla f(x_k), p_k \rangle$, meaning that $\tilde{\alpha}$ strictly satisfies the Wolfe condition, and by continuity there exists a neighborhood of $\tilde{\alpha}$ on which this remains true. \square

5 Inertial acceleration: Nesterov and Heavy-ball

In [Proposition 3.6](#), we have shown a $\mathcal{O}(1/k)$ convergence rate of gradient descent on L -smooth convex functions⁵; moreover if the function is μ strongly convex, the rate improves – without modifying the algorithm – to $\mathcal{O}(\exp(-\frac{\mu}{L}k))$. But this can still be slow: consider a function for which $\mu/L = 10^{-3}$. Then, to divide the suboptimality by 10, the number k of iterations to perform must be such that $\exp(-\frac{\mu}{L}k) \leq 0.1$, i.e. $k \geq \frac{L}{\mu} \ln(10) \approx 2300$.

In this section, we introduce first order algorithms that are faster. Namely, their convergence rates are in $\mathcal{O}(1/k^2)$ and $\exp(-\sqrt{\frac{\mu}{L}}k)$. This means that, on the same example, only $k \geq \sqrt{\frac{L}{\mu}} \ln(10) \approx 72$ would be needed.

5.1 Polyak’s heavy ball method

To understand momentum, a very cool resource with animations is <https://distill.pub/2017/momentum/>.

For a L -smooth, μ -strongly convex function f , the Heavy Ball iterations ([Polyak, 1964](#)) are defined as:

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \underbrace{\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}(x_k - x_{k-1})}_{\text{momentum}} \quad (5.1)$$

In this algorithm, the iterate moves in the direction of the gradient through $\alpha \nabla f(x_k)$, but it also continues in the direction of its previous update step $x_k - x_{k-1}$. Actually, this is where heavy ball name get its name from: the momentum term can be seen as the inertia of a ball rolling down a valley.

In this class we will only apply it on quadratics; beware that it can diverge otherwise (see the negative result by [Goujaud et al. \(2023\)](#)). For results beyond quadratics, see [Ghadimi et al. \(2015\)](#) and the recent works of [Apidopoulos et al. \(2022\)](#); [Aujol et al. \(2023\)](#).

It may be enlightening to introduce $\gamma = \frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}} = \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$ and rewrite heavy ball as:

$$\begin{cases} m_{k+1} = \gamma m_k + (1 - \gamma) \nabla f(x_k) & \left(= (1 - \gamma) \sum_{i=0}^k \gamma^i \nabla f(x_{k-i}) \right) \\ x_{k+1} = x_k - \frac{\alpha}{1-\gamma} m_{k+1} \end{cases} \quad (5.2)$$

which makes it visible that Heavy Ball replaces the gradient by a geometrically weighted average of past gradients. These two forms are equivalent: indeed unrolling the second equation we get

$$x_{k+1} = x_k - \frac{\alpha}{1-\gamma} \gamma m_k - \alpha \nabla f(x_k) \quad (5.3)$$

$$= x_k - \alpha \nabla f(x_k) + \gamma(x_k - x_{k+1}) \quad (5.4)$$

which is [Equation \(5.1\)](#).

⁵very recent and puzzling work by [Grimmer \(2023\)](#) improves this rate to $\mathcal{O}(1/(k \ln k))$

Proposition 5.1. *Let f be a L -smooth μ -strongly convex quadratic function, let $\kappa = L/\mu$. Then for a proper choice of α , Heavy Ball has the convergence rate:*

$$f(x_k) - f(x^*) \leq Ck \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x_0 - x^*\|^2 \quad (5.5)$$

Proof. Proof sketch: since we are minimizing a quadratic, we can wlog consider that the Hessian is diagonal. This allows writing down the update of (x_j, m_j) for an arbitrary coordinate $j \in [d]$. Show that a 2×2 matrix appears, put it in triangular form, find value of all powers of said matrix, bound its spectral norm by its absolute norm, show that the two values on its diagonal are complex conjugates, conclude. \square

5.2 Nesterov acceleration

The heavy ball method suffers from two flaws: it may not converge on “nice” functions beyond quadratics, and it does not apply to functions that are only convex (although, for convex quadratics, there is a straightforward adaptation with the lowest non zero eigenvalue, everything remaining constant in the kernel of the Hessian).

Nesterov’s accelerated gradient, also called the Fast Gradient Method, solves these two issues. Schematically, the difference between heavy ball and Nesterov is the order in which gradient steps and momentum are applied: heavy ball does gradient then momentum, while Nesterov does momentum then gradient.

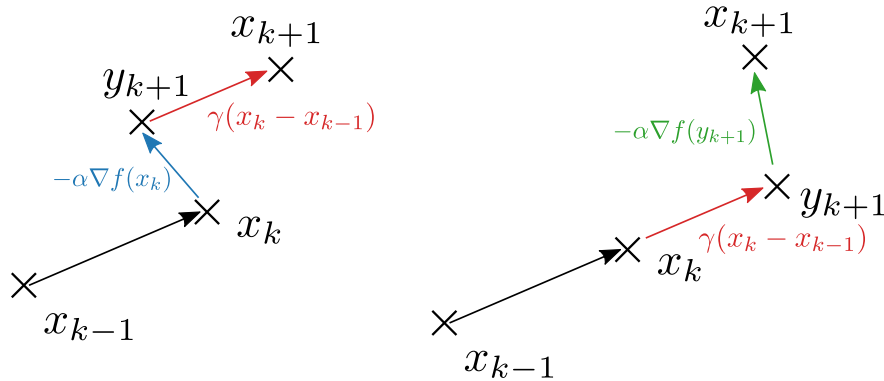


Figure 5.1: Left: heavy ball (gradient, then momentum). Right: Nesterov (momentum, then gradient).

Proposition 5.2 (Nesterov with strong convexity). *Let f be μ -strongly convex and L -smooth, with (unique) minimizer x^* . Let $\kappa = L/\mu$. Consider the following inertial algorithm: $y_0 = x_0$, and for $k \in \mathbb{N}^*$,*

$$\begin{cases} x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k) \\ y_{k+1} = x_{k+1} + \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} (x_{k+1} - x_k) \end{cases} \quad (5.6)$$

It has the following convergence rate:

$$f(x_k) - f(x^*) \leq \frac{L + \mu}{2} \left(1 - \sqrt{\frac{\mu}{L}} \right)^k \|x_0 - x^*\|^2. \quad (5.7)$$

Proof. The proof is not obvious and relies on the technique of *estimate sequences* (see [Definition 5.3](#)). It uses auxiliary functions, defined recursively:

$$\Phi_0(x) = f(x_0) + \frac{\mu}{2}\|x - x_0\|^2 \quad (5.8)$$

$$\Phi_{k+1}(x) = (1 - \sqrt{\frac{\mu}{L}})\Phi_k(x) + \sqrt{\frac{\mu}{L}}(f(y_k) + \langle \nabla f(y_k), x - y_k \rangle + \frac{\mu}{2}\|x - y_k\|^2) \quad (5.9)$$

One shows the four following results:

1. $\Phi_k(x) \leq f(x) + (1 - \sqrt{\frac{\mu}{L}})^k(\Phi_0(x) - f(x))$ (easy recursion)
2. $\nabla^2 \Phi_k(x) = m \text{ Id}$ (easy recursion), hence it is a quadratic and writes $\Phi_k(x) = \min \Phi_k + \frac{\mu}{2}\|x - \nu_k\|^2$
3. $f(x_k) \leq \min \Phi_k$ (this is the hard part)
4. $f(x_k) - f(x^*) \leq \Phi_k(x^*) - f(x^*)$

The main pain point is item 3, for which, by descent lemma and convexity successively:

$$f(x_{k+1}) \leq f(y_k) - \frac{1}{2L}\|\nabla f(y_k)\|^2 \quad (5.10)$$

$$= (1 - \kappa^{-1/2})f(y_k) + \kappa^{-1/2}f(y_k) - \frac{1}{2L}\|\nabla f(y_k)\|^2 \quad (5.11)$$

$$= (1 - \kappa^{-1/2})(f(y_k) - f(x_k)) + (1 - \kappa^{-1/2})f(x_k) + \kappa^{-1/2}f(y_k) - \frac{1}{2L}\|\nabla f(y_k)\|^2 \quad (5.12)$$

$$= (1 - \kappa^{-1/2})(f(y_k) - f(x_k)) + (1 - \kappa^{-1/2})f(x_k) + \kappa^{-1/2}f(y_k) - \frac{1}{2L}\|\nabla f(y_k)\|^2 \quad (5.13)$$

$$\leq (1 - \kappa^{-1/2})\langle \nabla f(y_k), y_k - x_k \rangle + (1 - \kappa^{-1/2})f(x_k) + \kappa^{-1/2}f(y_k) - \frac{1}{2L}\|\nabla f(y_k)\|^2 \quad (5.14)$$

By using $\nabla \Phi_{k+1}(\nu_{k+1}) = 0$ one obtains

$$\nu_{k+1} = (1 - \kappa^{-1/2})\nu_k + \kappa^{-1/2}(y_k - \frac{1}{\mu}\nabla f(y_k)) \quad (5.15)$$

Express $\Phi_{k+1}(y_k)$ with its quadratic formula and develop the term from ν_{k+1} ; also express $\Phi_{k+1}(y_k)$ using the recursion definition; at some point use $\nu_k - y_k = \sqrt{L/\mu}(y_k - x_k)$. You end up with

$$\Phi_{k+1}^* = \text{positive} + \text{the RHS term upper bounding } f(x_{k+1}) \quad (5.16)$$

It is truly a computational proof. \square

Definition 5.3 (Estimating sequence). *A sequence of function-scalar pairs (Φ_k, λ_k) is an estimating (or estimate) sequence of f if:*

1. $\lambda_k \geq 0$ for all k
2. $\lambda_k \rightarrow 0$
3. $\Phi_k(x) \leq (1 - \lambda_k)f(x) + \lambda_k\Phi_0(x)$ for any x

Estimating sequences provide a generic way to cook up convergence rates, because if (Φ_k, λ_k) is an estimate sequence and $f(x_k) \leq \min \Phi_k(x)$ one immediately has $f(x_k) - f(x^*) \leq \lambda_k(\Phi_0(x^*))$ which goes to 0 at the same rate as λ_k .

Let us make a few comments about Nesterov's algorithm:

- as κ approaches 1 (the function becomes better-conditioned), we need less and less extrapolation/acceleration and the algorithm gets closer to gradient descent.
- beware: the convergence rates holds for (x_k) , the sequence after the gradient step, but not for (y_k) . For recent work on this question, see [Attouch and Fadili \(2022\)](#)
- in practice the algorithm often leads to oscillating behavior, where the objective goes up and down periodically. This can be mitigated by *restart* ([O’Donoghue and Candes, 2015](#)).
- using the algorithm requires knowledge of μ ; when μ is unknown, one can resort to the version of [Proposition 5.4](#), or try to estimate μ ([Fercoq and Qu, 2019](#)).

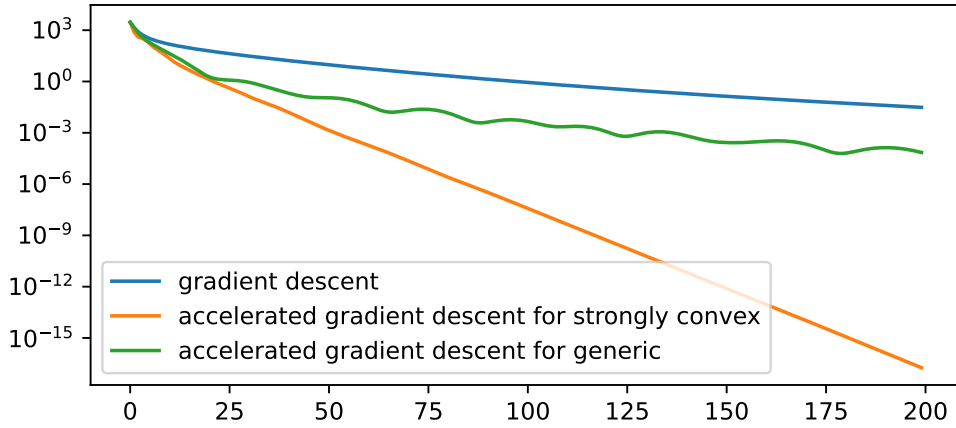


Figure 5.2: Gradient descent and Nesterov with interpolation parameter $\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$ (strongly convex) or $(t_k - 1)/t_{k+1}$ (generic). The objective is a quadratic with Hessian equal to $A^\top A$, $A \in \mathbb{R}^{100 \times 150}$ with i.i.d. Gaussian entries. All algorithms start at the same random point.

Proposition 5.4 (Nesterov without strong convexity). *Let f be convex and L -smooth, admitting at least one minimizer x^* . Let the sequence (t_k) be defined by $t_0 = 1$, and $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ (mysterious choice... let us only observe that it is equivalent to $t_k^2 = t_{k+1}^2 - t_{k+1}$, which makes the proof work⁶)*

Consider the following algorithm: $y_0 = x_0$, and for $k \in \mathbb{N}^$,*

$$\begin{cases} x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k) \\ y_{k+1} = x_{k+1} + \frac{t_k - 1}{t_{k+1}} (x_{k+1} - x_k) \end{cases} \quad (5.17)$$

Its iterates satisfy:

$$f(x_k) - f(x^*) \leq \frac{2L}{(k+1)^2} \|x_0 - x^*\|^2. \quad (5.18)$$

⁶inequality would be enough

Proof. Use the 3 point descent lemma ([Exercise 3.7](#)):

$$f(x_{k+1}) \leq f(x_k) + \langle \nabla f(y_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - y_k\|^2 \quad (5.19)$$

$$f(x_{k+1}) - f(x_k) \leq L \langle y_k - x_{k+1}, x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - y_k\|^2 \quad (5.20)$$

$$f(x_{k+1}) - f(x_k) \leq L \langle y_k - x_{k+1}, y_k - x_k \rangle - \frac{L}{2} \|x_{k+1} - y_k\|^2 \quad (5.21)$$

Use it a second time, with x^* instead of x_k (we have the right: in the derivation above, we have nothing specific about x_k):

$$f(x_{k+1}) - f(x^*) \leq L \langle y_k - x_{k+1}, y_k - x^* \rangle - \frac{L}{2} \|x_{k+1} - y_k\|^2 \quad (5.22)$$

Let $\delta_k = f(x_k) - f(x^*)$. Multiply (5.21) by $(t_k - 1)$ and add to (5.22):

$$t_k \delta_{k+1} - (t_k - 1) \delta_k \leq L \langle y_k - x_{k+1}, t_k y_k - (t_k - 1)x_k - x^* \rangle - \frac{L}{2} t_k \|y_k - x_{k+1}\|^2 \quad (5.23)$$

Multiplying by t_k , using the equality satisfied by t_k , and using the parallelogram identity:

$$t_k^2 \delta_{k+1} - t_{k-1}^2 \delta_k \leq \frac{L}{2} (2 \langle t_k(y_k - x_{k+1}), t_k y_k - (t_k - 1)x_k - x^* \rangle - \|t_k(y_k - x_{k+1})\|^2) \quad (5.24)$$

$$= \frac{L}{2} (\|t_k y_k - (t_k - 1)x_k - x^*\|^2 - \|t_k x_{k+1} - (t_k - 1)x_k - x^*\|^2) \quad (5.25)$$

Letting $u_k = t_k y_k - (t_k - 1)x_k - x^*$, we obtain, using $t_k x_{k+1} - (t_k - 1)x_k = t_{k+1} y_{k+1} - (t_{k+1} - 1)x_k$:

$$t_k^2 \delta_{k+1} - t_{k-1}^2 \delta_k \leq \frac{L}{2} (\|u_k\|^2 - \|u_{k+1}\|^2) \quad (5.26)$$

Sum up and use that, by recursion, $t_k \geq \frac{k-1}{2}$. □

6 Lower bounds for first-order methods

So far we have shown that gradient descent has rate $\mathcal{O}(1/k)$ or $\mathcal{O}(\exp(-k/\kappa))$, and with Nesterov acceleration, we improved it to $\mathcal{O}(1/k^2)$ or $\mathcal{O}(\exp(-k/\sqrt{\kappa}))$. Could we reach lower while using only the gradient information? The following shows that, in some sense, we cannot.

Definition 6.1 (First-order method). *A first-order method to minimize the function f is an algorithm that, given access to an oracle returning $g^{(i)} \in \partial f(x^{(i)})$ for each query point $x^{(i)}$, produces a sequence of iterates such that $x^{(k+1)} \in \text{Span}(g^{(0)}, \dots, g^{(k)})$. For simplicity, we also assume that $x^{(0)} = 0$.*

Proposition 6.2 (Lower bound for first-order methods, Lipschitz case). *Let $k \leq d - 1$. For any $L > 0, R > 0$, there exists a function f , L -Lipschitz on the ball of radius R , such that for any first-order method in the sense of [Definition 6.1](#),*

$$\min_{t \leq k} f(x^{(t)}) - \min_{\|x\| \leq R} f(x) \geq \frac{RL}{2(1 + \sqrt{k})}. \quad (6.1)$$

Crucial note: this is valid only for a number of iterations less than the dimension! But this is still relevant, as when the dimension d is large, one often performs less than d iterations.

Proof. Let $\gamma \geq 0, \alpha \geq 0$. Consider the function

$$f(x) = \gamma \max_{1 \leq i \leq k} x_i + \frac{\alpha}{2} \|x\|^2. \quad (6.2)$$

Its subdifferential is given by:

$$\partial f(x) = \alpha x + \gamma \text{conv}\{e_j : x_j = \max_{1 \leq i \leq k} x_i\}. \quad (6.3)$$

Hence, for $\|x\| \leq R$, subgradients have norm less than $\alpha R + \gamma$, and f is $(\alpha R + \gamma)$ -Lipschitz on the ball of radius R .

Suppose that the first-order oracle returns as subgradient $\alpha x + \gamma e_i$, where i is the smallest coordinate such that $x_i = \max_{1 \leq j \leq k} x_j$. Then, it is easy to see that $x^{(i)} \in \text{Span}(e_1, \dots, e_i)$, and thus for $i < d$, $x_d^{(i)} = 0$ and $f(x^{(i)}) \geq 0$.

Now let us compute the minimal value of f . It is equal to

$$\min \gamma v + \frac{\alpha}{2} \|x\|^2 \quad \text{s.t. } x_i - v \leq 0 \quad \forall i \in [k] \quad (6.4)$$

Hence by introducing the Lagrangian $\mathcal{L}(t, x, \lambda) = \gamma v + \frac{\alpha}{2} \|x\|^2 + \sum_{i=1}^k \lambda_i (x_i - v)$, and using the first-order optimality conditions, one gets that the minimum is reached for x^* such that $x_i^* = -\frac{\gamma}{\alpha k}$ for $i \leq k$, and $x_i^* = 0$ otherwise, yielding $f(x^*) = -\frac{\gamma^2}{2\alpha k}$. Therefore for any first-order algorithm and $i \leq k$,

$$f(x^{(i)}) - f(x^*) \geq \frac{\gamma^2}{2\alpha t}. \quad (6.5)$$

Pick $\alpha = \frac{L}{R(1+\sqrt{k})}$ and $\gamma = \frac{L\sqrt{k}}{1+\sqrt{k}}$ to prove the result (one needs to check that $\|x^*\| \leq R$, which is the case). \square

Note: the same function with $\gamma = L/2$ and $R = L/(2\alpha)$ can also be used to prove lower bounds for L -Lipschitz, α -strongly convex functions (see [Bubeck et al. \(2015, Thm 3.13\)](#)).

Proposition 6.3 (Lower bounds for first-order methods, smooth case). *Let k be an integer such that $k \leq (d-1)/2$. For any $L > 0$, there exists a L -smooth function with minimizer x^* such that for any first-order method in the sense of [Definition 6.1](#),*

$$\min_{t \leq k} f(x^{(t)}) - \min f \geq \frac{3L}{32} \frac{\|x_0 - x^*\|^2}{(k+1)^2}. \quad (6.6)$$

Again, be careful: this result is valid only for a number of iterations that does not exceed half the dimension.

Proof. For $i \leq d$, let H_i be the tridiagonal matrix with 2 on the diagonal and -1 on the over- and underdiagonal:

$$H_i = \begin{pmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & -1 \\ 0 & \dots & 0 & -1 & 2 \end{pmatrix} \in \mathbb{R}^{i \times i} \quad (6.7)$$

Let $A_i \in \mathbb{R}^{d \times d}$ be the d by d matrix such that its upper left i by i block is equal to H_i , and the rest is 0. One has $0 \preceq A_i \preceq 4\text{Id}_d$ since for all $x \in \mathbb{R}^d$:

$$x^\top A_i x = 2 \sum_{j=1}^i x_j^2 - 2 \sum_{j=1}^{i-1} x_j x_{j+1} \quad (6.8)$$

$$= \sum_{j=1}^{i-1} (x_j - x_{j+1})^2 + x_1^2 + x_i^2 \geq 0 \quad (6.9)$$

$$\leq 2 \sum_{j=1}^{i-1} (x_j^2 + x_{j+1}^2) + x_1^2 + x_i^2 \quad (6.10)$$

$$\leq 4 \sum_{j=1}^i x_j^2. \quad (6.11)$$

Let the function f to optimize be defined by:

$$f(x) = \frac{L}{8} x^\top A_{2k+1} x - \frac{L}{4} \langle x, e_1 \rangle. \quad (6.12)$$

Let $t \leq k$. Because $\nabla f(x) = \frac{L}{4}(A_{2k+1}x - e_1)$, we have that iterates of any first-order method will satisfy $x^{(t)} \in \text{Span}(e_1, \dots, e_t)$. In particular, $\langle x^{(t)}, A_{2k+1}x^{(t)} \rangle = \langle x^{(t)}, A_t x^{(t)} \rangle$. Let $f_t(x) = \frac{L}{8} x^\top A_t x - \frac{L}{4} \langle x, e_1 \rangle$, with minimal value f_t^* . It is clear that f_t^* decreases with t . We have proved that

$$f(x^{(t)}) - f^* = f_t(x^{(t)}) - f^* \geq f_t^* - f_{2k+1}^* \geq f_k^* - f_{2k+1}^*. \quad (6.13)$$

Now we compute the minimizer $x^{*,t}$ of f_t . It must satisfy $A_t x^{*,t} = e_1$. To solve this linear system, show by recursion on i that $x_{t-i}^{*,t} = (i+1)x_t^{*,t}$, then use the condition $2x_1^{*,t} - x_2^{*,t} = 1$ to conclude that:

$$x_i^{*,t} = \begin{cases} 1 - \frac{i}{t+1}, & \text{if } 1 \leq i \leq t, \\ 0 & \text{otherwise.} \end{cases} \quad (6.14)$$

Thus $f_t^* = -\frac{L}{8}x^{*,t\top}e_1 + \frac{L}{4}x^{*,t\top}e_1 = -\frac{L}{8}\frac{k}{k+1}$.

Now

$$\|x^{*,t}\|^2 = \sum_1^t \left(1 - \frac{i}{t+1}\right)^2 = \sum_1^t \left(\frac{i}{t+1}\right)^2 = \frac{t(t+1)(2t+1)}{6(t+1)^2} \leq \frac{t+1}{3} \quad (6.15)$$

Thus finally:

$$f_k^* - f_{2k+1}^* = \frac{L}{8} \left(\frac{k}{k+1} - \frac{2k-1}{2k+2} \right) \quad (6.16)$$

$$= \frac{L}{8} \frac{1}{2(k+1)} \quad (6.17)$$

$$\geq \frac{3L}{32} \frac{\|x^{*,2k+1}\|^2}{(k+1)^2} \quad (6.18)$$

where we have used $\|x^{*,2k+1}\|^2 \leq (2k+2)/3$. □

A last lower bound exists, for L -smooth, μ -strongly convex functions. It is shown in infinite dimension for commodity reasons; the reader can refer to [Nesterov et al. \(2018, Thm. 2.1.13\)](#).

7 Non-smooth convex optimization

For a differentiable objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we have seen:

1. convexity: f is above its tangents
2. Lipschitzness
3. L -smoothness, implying the descent lemma (Lemma 3.4), meaning that at any point $y \in \mathbb{R}^d$ there exists an isotropic parabola of curvature L that globally **upper** bounds f and is tangent to f at y .
4. strong convexity: at any point $y \in \mathbb{R}^d$ there exists an isotropic parabola of curvature μ that globally **lower** bounds f and is tangent to f at y .

Depending on the properties satisfied by f , we have:

- 1 gives the necessary and sufficient condition for x to be a minimizer: $\nabla f(x^*) = 0$.
- 1 + 2 gives a convergence rate of $1/\sqrt{k}$ under stepsize depending on the total number of iterations. The rate is in objective value, for the best iterate or the averaged iterates.
- 1 + 3 improves it to $1/k$ with stepsize $1/L$ not depending on the number of iterations. The rate is in objective value, on the last iterate.
- 3 + 4 sandwiches the functions between two isotropic parabolas of curvature μ and L and gives the excellent linear rate. The rate is on the distance from the last iterate to the unique minimizer.

Now we want to consider non-differentiable functions f . Why do we need to handle these?

7.1 Constrained optimization

Optimizers sometimes want the minimizer of their cost functions to belong to some set $\mathcal{C} \subset \mathbb{R}^d$.

For example, suppose that the variable x corresponds to an image, stored as a 2D array of pixels – an element of $\mathbb{R}^{d \times d}$. Then the values of the pixels are not just any real numbers: they should be in $[0, 1]$ if they correspond to gray levels for example.

If the optimization variable $x \in \mathbb{R}^d$ represents proportions (say, of various types of cells in a sane/cancerous tissue), it should have non-negative entries that sum to 1: $x \in \Delta_d$, where Δ_d is the d -dimensional *simplex*

$$\Delta_d = \{x \in \mathbb{R}^d : \forall i \in [d], x_i \geq 0, \sum_{i=1}^d x_i = 1\} . \quad (7.1)$$

Finally, suppose the practitioner is looking for the approximate solution to a linear system (using least squares) but wants a solution x^* that uses as few variables as possible. There are many ways to do this, but one popular (because convex) way to do so is to constrain the ℓ_1 norm of x :

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \|Ax - b\|^2 \quad \text{subject to} \quad \|x\|_1 \leq \tau .$$

All these examples lead us to consider constrained problems:

$$\min_{x \in \mathcal{C}} f(x) . \quad (7.2)$$

Because of the constraint, even if f is still smooth, the tools we have used so far no longer apply. Consider the basic one-dimensional problem:

$$\min_{x \in [0,1]} x . \quad (7.3)$$

This problem is friendly: the objective is convex, and the constraint set $[0, 1]$ is convex and compact. The minimizer is 0, yet $\nabla f(0) = 1 \neq 0$ and so for constrained convex optimization, the global characterization of minimizers (Proposition 2.7) does not hold.

Fortunately, we can replace some concepts seen so far with generalizations that are very well adapted. First, let us introduce a tool that allows removing the constraints.

7.2 Extended value functions

In optimization, it is often convenient to work with functions that take values in $\overline{\mathbb{R}} \triangleq \mathbb{R} \cup \{+\infty\}$. The prototypical example is the *indicator* function of a set.

Definition 7.1 (Indicator function). *In convex analysis, the indicator function $\iota_{\mathcal{C}}$ of a subset \mathcal{C} of \mathbb{R}^d is:*

$$\begin{aligned} \iota_{\mathcal{C}} : \mathbb{R}^d &\rightarrow \overline{\mathbb{R}} \\ x &\mapsto \begin{cases} 0 , & \text{if } x \in \mathcal{C} , \\ +\infty , & \text{otherwise} . \end{cases} \end{aligned} \quad (7.4)$$

When is $\iota_{\mathcal{C}}$ convex?

Note that this “optimization” indicator is different from the other, “boolean” indicator function that takes the value 1 on the set and 0 outside⁷. The indicator function conveniently allows transforming all constrained minimization problems into unconstrained ones of course at the price of working with extended value functions:

$$\operatorname{argmin}_{x \in \mathcal{C}} f(x) = \operatorname{argmin}_{x \in \mathbb{R}^d} f(x) + \iota_{\mathcal{C}}(x) . \quad (7.5)$$

Definition 7.2. *The domain of a function $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is:*

$$\operatorname{dom} f \triangleq \{x \in \mathbb{R}^d, f(x) < +\infty\} . \quad (7.6)$$

The second tool to overcome non-differentiability is a substitute for the gradient.

⁷this function is rarely convex

7.3 Subdifferential of convex functions

Definition 7.3. *The subdifferential of f at x is the set of slopes of all affine minorants of f that are exact at x :*

$$\partial f(x) = \{u \in \mathbb{R}^d : \forall y \in \mathbb{R}^d, f(y) \geq f(x) + \langle u, y - x \rangle\} . \quad (7.7)$$

An element of the subdifferential is called a subgradient.

Note that contrary to the gradient, the subdifferential is a set-valued mapping: the subdifferential at one point may contain more than one subgradient. It may also be empty at some point ([Exercise 7.1](#)); the set of points where the subdifferential is not empty is called the *domain* of the subdifferential (not to be confounded with for extended-value functions). As [Exercise 7.8](#) shows, if f is convex the domain of ∂f contains the interior of the domain of f .

The subdifferential is a generalization of the gradient in the following sense.

Proposition 7.4. *Let f be a convex differentiable function. Then the subdifferential only contains the gradient:*

$$\partial f(x) = \{\nabla f(x)\} . \quad (7.8)$$

Note: the converse is true (if the subdifferential reduces to a point, f is differentiable at this point).

The subdifferential allows an elegant characterization of the minimizers of *all* convex functions, even the nondifferentiable ones.

Proposition 7.5 (Fermat's rule). *Let f be convex. Then for all x ,*

$$x \in \operatorname{argmin} f \Leftrightarrow 0 \in \partial f(x) . \quad (7.9)$$

The proof is 1 line and left to the reader.

Example 7.6. *The subdifferential of the absolute value is:*

$$\partial |\cdot|(x) = \begin{cases} \{x/|x|\}, & \text{if } x \neq 0 , \\ [-1, 1], & \text{otherwise} . \end{cases} \quad (7.10)$$

What is the subdifferential of $\iota_{[0,1]}$?

Proposition 7.7 (Subdifferential of sum). *The subdifferential of the sum contains the sum of subdifferentials:*

$$\partial(f + g) \supset \partial f + \partial g , \quad (7.11)$$

(to be understood as the Minkowski sum). Equality does not hold in general, although the counterexamples are pathological ([Exercise 7.10](#)).

Remark 7.8 (Subgradient descent). *If you look again with nonsmooth eyes at [Proposition 3.2](#) regarding gradient descent on Lipschitz convex functions, we actually never used the fact that f was differentiable: we only started with $f(x_k) - f(x^*) \leq \langle \nabla f(x_k), x_k - x^* \rangle$. Since a subgradient of f at x_k , by definition, also satisfies this inequality, it means we can replace $\nabla f(x_k)$ by any $g_k \in \partial f(x_k)$ and the proof still holds.*

This algorithm is called subgradient descent.

The last tool we need for nonsmooth optimization is the proximal operator.

7.4 The proximal operator

So far we have brushed away the existence of minimizers, and only assumed they exist. In classical analysis, the celebrated Weierstrass theorem states that on any compact, a continuous function is bounded and reaches its minimum. We will provide a slight relaxation of this theorem that handles a slightly weaker property than continuity.

Definition 7.9 (Lower semi-continuity). *The function $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is lower semicontinuous at $x \in \text{dom } f$ if for any sequence (x_k) converging to x ,*

$$f(x) \leq \liminf_{k \rightarrow \infty} f(x_k) . \quad (7.12)$$

Any continuous function is clearly lower semicontinuous.

Why didn't we bother with lower semicontinuity before? The answer is simple but relies on a deep result: a convex function is continuous on the interior of its domain⁸. Since we were previously considering convex functions defined on \mathbb{R}^d and not taking value $+\infty$, they had to be continuous on the whole \mathbb{R}^d .

Remark 7.10. *Lower semicontinuous functions are also referred to as closed functions. This is because a function is l.s.c. if and only if its epigraph*

$$\text{epi } f = \{(x, t) : f(x) \leq t\} \subset \mathbb{R}^d \times \mathbb{R} \quad (7.13)$$

is closed.

Proposition 7.11 (Existence of minimizers, the direct method of calculus of variations). *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be a coercive lower semicontinuous function, with nonempty domain. Then f admits at least one minimizer.*

Note: this result is only valid in finite dimension, for a counterexample in infinite dimension see [Appendix B](#).

Proof. Let x_0 such that $f(x_0) < +\infty$. Since f is coercive, $f(x) \rightarrow_{\|x\| \rightarrow \infty} +\infty$. Let M such that $\|x\| \geq M \implies f(x) \geq f(x_0)$. Let \mathcal{B} be the ball of center 0, radius M . Since we are in finite dimension, \mathcal{B} is compact.

Let $f^* = \inf_{x \in \mathcal{B}} f(x)$. In general, we could have $f^* = -\infty$ (think $d = 1$, $M = 1$, $f(x) = 1/x$ except $f(0) = 0$). But we'll show that since f is l.s.c., it can't be the case.

Assume $f^* = -\infty$. We can construct a sequence (x_k) in \mathcal{B} such that $\forall k \in \mathbb{N}$, $f(x_k) \leq -k$. But \mathcal{B} is compact: we can extract a converging subsequence, that we call y_k , converging to $\tilde{x} \in \mathcal{B}$.

By lower semicontinuity of f , $f(\tilde{x}) \leq \liminf_{k \rightarrow \infty} f(y_k) = -\infty$ which is not possible since f does not take value $-\infty$.

Hence f^* is finite. Now we do the exact same reasoning to construct (x_k) such that $f(x_k) \leq f^* + 1/k$. Take a converging subsequence, then use lower semicontinuity to show that the limit point \tilde{x} satisfies $f(\tilde{x}) \leq f^*$ and thus $f(\tilde{x}) = f^*$. \square

Definition 7.12. *The function $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is said to be (or more rigorously belong to) $\Gamma_0(\mathbb{R}^d)$ if it is:*

⁸better: it is locally Lipschitz

- *convex*
- *proper (its domain is not empty)*
- *lower semicontinuous.*

To get familiar with these three notions, determine under which condition of the set $\mathcal{C} \subset \mathbb{R}^d$ does $\iota_{\mathcal{C}}$ belong to $\Gamma_0(\mathbb{R}^d)$.

Definition 7.13 (Proximal operator). *Let $f \in \Gamma_0(\mathbb{R}^d)$. The proximal operator of f evaluated at x is*

$$\text{prox}_f(x) = \underset{y \in \mathbb{R}^d}{\text{argmin}} f(y) + \frac{1}{2} \|x - y\|^2 . \quad (7.14)$$

Why is it well-defined (why does the infimum exist? Why is it attained? Why does the argmin contain exactly one point)?

The following is a useful characterization of proximal operators.

Proposition 7.14 (Characterization of proximal operators). *Let $f \in \Gamma_0(\mathbb{R}^d)$. Then $p = \text{prox}_f(x)$ if and only if $x - p \in \partial f(p)$.*

This justifies the frequent formulation $\text{prox}_f = (\text{Id} + \partial f)^{-1}$ (a.k.a. the prox is the resolvent of the subdifferential).

Proof. Apply Fermat's rule. □

Proximal operators may seem new, but without realizing it you know and have used some particular instances: if \mathcal{C} is non-empty, closed and convex, $\text{prox}_{\iota_{\mathcal{C}}}$ is the (well-defined) projection onto \mathcal{C} . For this and other reasons, proximal operators can be thought of as a generalization of projections.

Example 7.15. *The proximal operator of $\frac{\lambda}{2} \|\cdot\|^2$ is a rescaling:*

$$\text{prox}_{\frac{\lambda}{2} \|\cdot\|^2}(x) = x / (1 + \lambda) . \quad (7.15)$$

Example 7.16. *The proximal operator of the absolute value is the soft-thresholding:*

$$\text{prox}_{\gamma|\cdot|}(x) = \begin{cases} x - \lambda, & \text{if } x > \lambda, \\ 0, & \text{if } x \in [-\lambda, \lambda], \\ x + \lambda, & \text{if } x < -\lambda. \end{cases} \quad (7.16)$$

Proposition 7.17 (Firm non-expansivity of proximal operators). *Proximal operators are firmly nonexpansive:*

$$\|\text{prox}_f(x) - \text{prox}_f(y)\|^2 \leq \langle \text{prox}_f(x) - \text{prox}_f(y), x - y \rangle . \quad (7.17)$$

By Cauchy-Schwarz inequality, this immediately implies that proximal operators are nonexpansive (so firmly nonexpansive is stronger than nonexpansive, the naming makes sense):

$$\|\text{prox}_f(x) - \text{prox}_f(y)\| \leq \|x - y\| . \quad (7.18)$$

Proof. Let $p_x = \text{prox}_f(x)$, $p_y = \text{prox}_f(y)$. The prox characterization for p_x is $x - p_x \in \partial f(p_x)$, so by definition of the subdifferential:

$$f(p_y) \geq f(p_x) + \langle x - p_x, p_y - p_x \rangle . \quad (7.19)$$

Do the same but swapping the roles of x and y , sum the two inequations to cancel out the f 's, and reorder. □

7.5 The proximal point algorithm

For $f \in \Gamma_0(\mathbb{R}^d)$ and $\lambda > 0$, the proximal point algorithm consists in iterating:

$$x_{k+1} = \text{prox}_{\lambda f}(x_k) \quad (7.20)$$

One motivation for this algorithm is that, by Fermat's rule and the characterization of proxs ([Proposition 7.14](#)), x^* is a minimizer of f if and only if it is a fixed point of $\text{prox}_{\lambda f}$: $x^* = \text{prox}_{\lambda f}(x^*)$ (for any $\lambda > 0$). Since $\text{prox}_{\lambda f}$ is nonexpansive (), we just iterate applications of this operator until convergence, in a Banach-Picard spirit.

This algorithm may seem a bit weird, as, in order to minimize, we need to be able to iteratively compute $\text{prox}_{\lambda f}(x)$, which is usually as hard as minimizing f . However:

- computing $\text{prox}_{\lambda f}$ may be easier than minimizing f , for example when f is a very badly conditioned quadratic: minimizing it requires solving a badly conditioned linear system, while the objective function in the prox minimization problem is f + an isotropic quadratic, which has better conditioning and is thus easier to minimize.
- the proximal point algorithm is a very powerful analysis tool, as it turns out most algorithms can be rewritten as an instance of it (usually, “refined” versions of it: in lifted spaces, with variable metric, etc.)

TODO: connect with the fixed-point section. TODO: say that it's a weird algorithm because to minimize f it requires minimizing f + parabola. But it's a powerful analysis tool, as algorithms rewrite as the proximal point algorithm in different spaces (eg Chambolle-Pock, ADMM ([?](#)))

7.6 The proximal gradient descent algorithm

Proposition 7.18. *Let f and g be Γ_0 , such that in addition f is L -smooth. Assume that there exists a minimizer x^* to $F = f + g$. Then the iterates of proximal gradient descent with stepsize $0 < \gamma \leq 1/L$ satisfy:*

$$F(x_k) - F(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\gamma k}. \quad (7.21)$$

Remark 7.19. *As in the case of gradient descent, convergence is guaranteed also for $1/L \leq \gamma < 2/L$, see for example [Section 9](#).*

Proof.

$$x_{k+1} = \text{prox}_{\gamma g}(x_k - \gamma \nabla f(x_k)) \Leftrightarrow \frac{x_k - x_{k+1}}{\gamma} - \nabla f(x_k) \in \partial g(x_{k+1}) \quad (7.22)$$

Hence for all x ,

$$g(x_{k+1}) - g(x) \leq \left\langle \frac{x_k - x_{k+1}}{\gamma} - \nabla f(x_k), x_{k+1} - x \right\rangle \quad (7.23)$$

and by the descent lemma:

$$f(x_{k+1}) - f(x_k) \leq \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 . \quad (7.24)$$

By using $x = x_k$ and summing, we get a first nice result: the proximal gradient descent method is a descent method for $0 < \gamma < 2/L$, because

$$F(x_{k+1}) - F(x_k) \leq \left(\frac{L}{2} - \frac{1}{\gamma} \right) \|x_{k+1} - x_k\|^2 . \quad (7.25)$$

Back to our general proof, we sum [Equations \(7.23\) and \(7.24\)](#):

$$g(x_{k+1}) - g(x) + f(x_{k+1}) - f(x_k) \leq \frac{1}{\gamma} \langle x_k - x_{k+1}, x_{k+1} - x \rangle + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 \quad (7.26)$$

$$\begin{aligned} g(x_{k+1}) - g(x) + f(x_{k+1}) - f(x) &\leq \frac{1}{\gamma} \langle x_k - x_{k+1}, x_{k+1} - x \rangle + f(x_k) - f(x) + \langle \nabla f(x_k), x - x_k \rangle \\ &\quad + \frac{L}{2} \|x_{k+1} - x_k\|^2 \end{aligned} \quad (7.27)$$

Since f is convex, $f(x_k) - f(x) + \langle \nabla f(x_k), x - x_k \rangle \leq 0$. Hence

$$F(x_{k+1}) - F(x) \leq \frac{1}{\gamma} \langle x_{k+1} - x_k, x_{k+1} - x \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 \quad (7.28)$$

$$= \frac{1}{2\gamma} (\|x_k - x\|^2 - \|x_{k+1} - x\|^2) + \frac{1}{2} \left(L - \frac{1}{\gamma} \right) \|x_{k+1} - x_k\|^2 \quad (7.29)$$

$$\leq \frac{1}{2\gamma} (\|x_k - x\|^2 - \|x_{k+1} - x\|^2) \quad (7.30)$$

for $0 < \gamma \leq 1/L$. Apply in $x = x^*$, sum and use the fact that suboptimality decreases at each iteration to obtain:

$$k(F(x_k) - F(x^*)) \leq \sum_{t=1}^k F(x_t) - F(x^*) \leq \frac{1}{2\gamma} \|x_0 - x^*\|^2 . \quad (7.31)$$

□

Remark 7.20 (Inexact prox). *When the proximal operator of g cannot be computed exactly, the same convergence rates may still be obtained, provided the error in the prox computation decreases along iteration at an appropriate rate; see [Schmidt et al. \(2011\)](#).*

Due to this algorithm, nonsmooth penalization became very popular in the late 2000s, with problems such as the Lasso, sparse logistic regression, non negative least squares, low rank matrix completion⁹, etc. The reader can refer to [Bach et al. \(2012\)](#).

⁹the nuclear norm is the ℓ_1 norm of the singular values; it induces sparsity in them, i.e. low rank on the matrix. Its prox is an SVD of the target matrix, with soft thresholding applied to the eigenvalues!

7.7 Exercises

7.7.1 Subdifferential

Exercise 7.1. ☹ Provide an example of convex function which has an empty subdifferential at some point of its domain. What is the position of that point with respect to the domain?

Exercise 7.2. ☹ Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Show that for any $x, u \in \mathbb{R}^n$, $\partial(f + \langle u, \cdot \rangle)(x) = \partial f(x) + \{u\}$.

Show that $\partial(f \circ A)(x) = A^\top \partial f(Ax)$ for $A \in \mathbb{R}^{n \times d}$.

Exercise 7.3. ☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be separable: $f(x) = \sum_1^d f_i(x_i)$ where the f_i 's are functions of the real variable.

Show that $\partial f(x) = \partial f_1(x_1) \times \dots \times \partial f_d(x_d)$.

Compute the subdifferential of the ℓ_1 -norm.

Exercise 7.4. ☹ Show that the subdifferential of a function at a point writes as an intersection of half-spaces. Show that it is convex and closed.

Exercise 7.5. ☹ Show that the subdifferential of ι_C at $x \in C$ is equal to the normal cone of C at x , that is:

$$\mathcal{N}_C(x) = \{u : \forall z \in C, \langle u, z - x \rangle \leq 0\}$$

Exercise 7.6. ☹☹ Compute the subdifferential of the Euclidean norm at any point in \mathbb{R}^d .

Exercise 7.7 (Exercise 7.12 revisited). ☹☹ Show that the first order optimality condition for differentiable convex constrained optimization rewrites:

$$x^* \in \operatorname{argmin}_C f(x) \Leftrightarrow -\nabla f(x^*) \in \mathcal{N}_C(x^*) .$$

Exercise 7.8 (Non emptiness of subdifferential). ☹☹☹ Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be convex. Show that if $x \in \operatorname{int}(\operatorname{dom} f)$, $\partial f(x)$ is nonempty and compact.

Show that if x lies on the boundary of $\operatorname{dom} f$, $\partial f(x)$ is either empty or unbounded.

Exercise 7.9. ☹ Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ such that $\operatorname{dom} f$ is convex. Show that if $\partial f(x)$ is nonempty at every $x \in \operatorname{dom} f$, then f is convex.

Exercise 7.10. ☹ This exercise is in dimension 2. Let C_1 and C_2 be two closed balls of strictly positive radius, that share a single point. Show that $\partial(\iota_{C_1} + \iota_{C_2})$ is a strict superset of $\partial \iota_{C_1} + \iota_{C_2}$ at this point.

7.7.2 Constrained optimization

Exercise 7.11. ☹ Show that the indicator function ι_C is convex (resp. lower semicontinuous, resp. proper) if C is convex (resp. closer, resp. nonempty).

Exercise 7.12 (Global optimality condition for constrained convex optimisation). ☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex differentiable function, let C be a convex subset of \mathbb{R}^d . Show that $x^* \in \operatorname{argmin}_{x \in C} f(x)$ if and only if

$$\forall x \in C, \langle \nabla f(x^*), x - x^* \rangle \geq 0 .$$

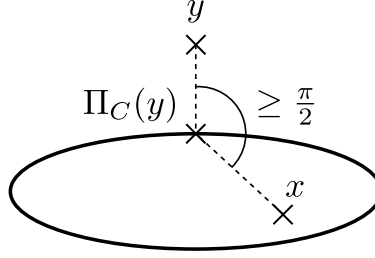


Figure 7.1: Illustration of (7.32)

Exercise 7.13. *Let C be a nonempty closed convex set. Show that for all $x \in C, y \in \mathbb{R}^d$,*

$$\langle y - \Pi_C(y), x - \Pi_C(y) \rangle \leq 0 . \quad (7.32)$$

Exercise 7.14. ☞ *Let C be a nonempty closed convex set. Show that for all $x \in C, y \in \mathbb{R}^d$,*

$$\|x - y\|^2 \geq \|x - \Pi_C(y)\|^2 + \|y - \Pi_C(y)\|^2 .$$

In particular, this shows that $\|y - x\| \geq \|x - \Pi_C(y)\|$, which says that projection can only get you closer to optimum (taking $x = x^$, if your current iterate is y).*

Exercise 7.15. ☞ *Let C be a nonempty closed convex set. Show that the overprojection, $2\Pi_C - \text{Id}$, is nonexpansive (i.e. it has Lipschitz constant 1).*

8 Duality

8.1 Motivation

Suppose we have a signal $s \in \mathbb{R}^d$ which is a noisy version of a true signal s^* , that we assume to be constant by part, and we want to recover s^* from s .

One possible way to do so would be to solve:

$$\min_x \frac{1}{2} \|x - s\|^2 + \lambda \sum_{j=1}^{d-1} |x_{j+1} - x_j| = \frac{1}{2} \|x - s\|^2 + \lambda \|Dx\|_1 \quad (8.1)$$

where D is the discrete derivation operator. As we know that penalizing the ℓ_1 norm of a vector makes this vector sparse, we will end up with a solution x such that Dx is sparse; in turn, a sparse Dx means that x is constant by parts with only a few jumps.

How do we solve [Problem \(8.1\)](#)? It is a convex problem, with a smooth + nonsmooth structure, and we know how to compute the prox of the ℓ_1 norm, so we are tempted to use the proximal gradient algorithm. But here comes a big distinction between the smooth and the nonsmooth term: if we know the gradient of f , we can easily compute the gradient of $f \circ A$. But that is no longer the case for proximal operators! In fact, for $g = \|\cdot\|_1$, computing the prox of $g \circ D$ is [Problem \(8.1\)](#) in itself, for which there is no closed form.

To try to solve [Problem \(8.1\)](#), let's do a few algebraic manipulations. We set λ to 1 for simplicity.

$$\min_x \frac{1}{2} \|x - s\|^2 + \|Dx\|_1 \quad (8.2)$$

$$\Leftrightarrow \min_{x,u} \frac{1}{2} \|x - s\|^2 + \|u\|_1 \quad \text{subject to} \quad u = Dx \quad (8.3)$$

$$\Leftrightarrow \min_{x,u} \frac{1}{2} \|x - s\|^2 + \|u\|_1 + \iota_{\{u\}}(Dx) \quad (8.4)$$

$$\Leftrightarrow \min_{x,u} \frac{1}{2} \|x - s\|^2 + \|u\|_1 + \max_y \langle y, Dx - u \rangle \quad (8.5)$$

$$\Leftrightarrow \max_y \min_{x,u} \frac{1}{2} \|x - s\|^2 + \|u\|_1 + \langle y, Dx - u \rangle \quad (8.6)$$

$$\Leftrightarrow \max_y \left(\min_x \frac{1}{2} \|x - s\|^2 + \langle D^\top y, x \rangle \right) + \left(\min_u \|u\|_1 - \langle y, u \rangle \right) \quad (8.7)$$

$$\Leftrightarrow \max_y -\frac{1}{2} \|D^\top y\|^2 + \langle D^\top y, s \rangle - \iota_{\mathcal{B}_{\|\cdot\|_\infty}}(y) \quad (8.8)$$

$$\Leftrightarrow \min_y \frac{1}{2} \|D^\top y\|^2 - \langle D^\top y, s \rangle + \iota_{\mathcal{B}_{\|\cdot\|_\infty}}(y) \quad (8.9)$$

where we have successfully used that $\iota_{\{u\}}(Dx) = \max_y \langle Dx - u, y \rangle$, switched max and min and solved the optimization problems in x and u in closed form (the first one is a quadratic minimization, the second one can be done by hand by noticing that the problem is separable in u ; one can also use [Exercise 8.2](#)).

The minimization problem obtained in Equation (8.9) is much more amenable to proximal gradient descent, because the linear operator, D^\top , is now acting on the smooth term; the nonsmooth part is $\iota_{\mathcal{B}_{\|\cdot\|_\infty}}$, whose proximal operator can be computed: it is the projection onto the unit ball of the ℓ_∞ norm.

However, we have been carefree when deriving our equivalence: it is not obvious that one can invert max and min in (8.6). In what follows, we will formalize what we have accomplished (namely, derive the dual problem from the primal problem), and derive conditions under which these problems are equivalent.

8.2 The Fenchel transform

A fundamental tool in convex analysis is the Fenchel transform (sometimes called the Fenchel-Legendre or Legendre transform), which is, in many ways, akin to the Fourier transform in signal processing.

Definition 8.1 (Fenchel conjugate). *Let $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$. The Fenchel conjugate (or transform) of f is:*

$$f^*(u) = \sup_{x \in \mathbb{R}^d} \langle u, x \rangle - f(x) . \quad (8.10)$$

Go back to the computation of Equation (8.7): you see that the Fenchel transforms of our two functions appeared!

It is always convex, even if f is not (Exercise 2.2).

Geometrical interpretation of the Fenchel transform $f^*(u) \leq \alpha$ is equivalent to

$$\forall x \in \mathbb{R}^d, \quad \langle x, u \rangle - \alpha \leq f(x) , \quad (8.11)$$

meaning that $-\alpha$ is the intercept of a global affine minorant of f of slope u . $-f^*(u)$ is thus the biggest possible intercept of such minorants – beware, there may not exist such minorants ($f^*(u) = +\infty$).

To get familiar with the Fenchel transform, compute the Fenchel transforms of $\|\cdot\|^2/2a$ ($a > 0$), $x \mapsto \langle x, b \rangle$, and $x \mapsto \iota_{\{b\}}$. What do you observe? This intuition is verified by the following result.

Proposition 8.2 (Involution on Γ_0). *If $f \in \Gamma_0(\mathbb{R}^d)$, $f^{**} = f$.*

Proof. The proof technique is (to my knowledge) not really used in other optimization proofs, but it uses a beautiful argument.

By definition, $f^{**}(x) = \sup_u \langle u, x \rangle - f^*(u)$. Let us write $\phi_u(x) = \langle u, x \rangle - f^*(u)$. What are these functions? They are affine, and they globally lower bound f by Fenchel-Young inequality: $\forall x \in \mathbb{R}^d, \phi_u(x) \leq f(x)$. So

$$f^{**}(x) = \sup_u \phi_u(x) \leq \sup_{\substack{a \text{ affine} \\ a \leq f}} a(x) . \quad (8.12)$$

In fact, these two supremums are equal: take an affine function globally lower bounding f , with slope u . It must therefore be equal to $\langle x, u \rangle - c$ for some $c \in \mathbb{R}$. For all $x \in \mathbb{R}^d$, $-c + \langle x, u \rangle \leq f(x)$, so $c \geq \sup_x \langle x, u \rangle - f(x) = f^*(u)$. Therefore $\langle \cdot, u \rangle - c \leq \phi_u$.

Hence, the supremum over all affine functions lower bounding f is equal to the supremum over the “extremal” functions ϕ_u only:

$$f^{**}(x) = \sup_{\substack{a \text{ affine} \\ a \leq f}} a(x) . \quad (8.13)$$

This shows that $f^{**} \leq f$ all the time (no requirement on convexity).

Now we show equality when f is convex. Take $y \in \mathbb{R}, x \in \mathbb{R}^d$ such that $y < f(x)$ i.e. (y, x) not in the epigraph of f . The latter is convex, so by the celebrated Hahn-Banach theorem, we can find a hyperplane separating the two. This hyperplane corresponds to an affine function (TODO what if it is vertical?) that is above y at x , at globally lower bounds f . So $f^{**}(x) > y$. This concludes the proof.

TODO this also means that f^{**} is the best pointwise: if $f^{**}(x) < g(x)$ with g convex and $g \leq f$, then we can apply Hahn-Banach too and construct an affine a that is above $(x, f^{**}(x))$ and globally below f , which is impossible. □

Proposition 8.3 (Fenchel-Young inequality). *Let $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$. Then for all $x, u \in \mathbb{R}^d$,*

$$f(x) + f^*(u) \geq \langle x, u \rangle . \quad (8.14)$$

Equality holds if and only if $u \in \partial f(x)$.

Proof. The first part is trivial by definition of f^* . Then,

$$\begin{aligned} f^*(u) = \langle x, u \rangle - f(x) &\Leftrightarrow x \in \operatorname{argmin} f - \langle \cdot, u \rangle \\ &\Leftrightarrow 0 \in \partial(f - \langle \cdot, u \rangle)(x) \\ &\Leftrightarrow u \in \partial f(x) . \end{aligned}$$

□

Proposition 8.4. *If $f \in \Gamma_0(\mathbb{R}^d)$, $\partial f^* = (\partial f)^{-1}$, in the sense that:*

$$u \in \partial f(x) \Leftrightarrow x \in \partial f^*(u) . \quad (8.15)$$

Proof. Use then Fenchel-Young equality for f and then use $f^{**} = f$. □

Proposition 8.5 (Moreau decomposition formula). *Let $f \in \Gamma_0(\mathbb{R}^d)$. Then the proximal operators of f and that of f^* are related together by:*

$$\operatorname{prox}_{\gamma f} + \gamma \operatorname{prox}_{\gamma^{-1} f^*}(\cdot / \gamma) = \operatorname{Id} \quad (8.16)$$

Proof.

$$p = \text{prox}_{\gamma f}(x) \Leftrightarrow \frac{x-p}{\gamma} \in \partial f(p) \quad (8.17)$$

$$\Leftrightarrow p \in \partial f^* \left(\frac{x-p}{\gamma} \right) \quad (8.18)$$

$$\Leftrightarrow \frac{x}{\gamma} - \frac{x-p}{\gamma} \in \partial_{\frac{1}{\gamma}} f^* \left(\frac{x-p}{\gamma} \right) \quad (8.19)$$

$$\Leftrightarrow \frac{x-p}{\gamma} = \text{prox}_{f^*/\gamma} \left(\frac{x}{\gamma} \right) . \quad (8.20)$$

□

Proposition 8.6. *For $\mu > 0$, f is μ -strongly convex iff f^* is $\frac{1}{\mu}$ -smooth.*

Definition 8.7 (Moreau envelope). *Let $f \in \Gamma_0(\mathbb{R}^d)$, let $\lambda > 0$. The Moreau envelope of f with parameter λ , denoted M_f^λ , is a function from \mathbb{R}^d to \mathbb{R} , which is the infimal convolution of f with $\frac{1}{2\lambda} \|\cdot\|^2$:*

$$M_f^\lambda(x) = (f \square \frac{1}{2\lambda} \|\cdot\|^2)(x) \quad (8.21)$$

$$= \min_{y \in \mathbb{R}^d} \frac{1}{2\lambda} \|y - x\|^2 + f(y) \quad (8.22)$$

(the inf is a min because f is Γ_0). Note: this can be defined even if f is not Γ_0 , provided the inf exists, but we will not use it in that setting.

Remark 8.8. *The Moreau envelope is closely related to the prox of f : $\text{prox}_{\lambda f}(x)$ is the value of x for which the min in (8.22) is reached. Beware however that the prox is often defined as minimizing $\frac{1}{2} \|\cdot - x\|^2 + \lambda f(\cdot)$ rather than minimizing $\frac{1}{2\lambda} \|\cdot - x\|^2 + f(\cdot)$. Of course, this does not change anything, but this allows speaking about $\text{prox}_{\lambda f}$ instead of $\text{prox}_{\lambda, f}$. However, in the Moreau problem the objective must be $\frac{1}{2\lambda} \|\cdot - x\|^2 + f(\cdot)$; if it were $\frac{1}{2} \|\cdot - x\|^2 + \lambda f(\cdot)$, it would change the value of M_f^λ (it would multiply it by λ).*

Proposition 8.9. *As the pointwise infimum of convex functions (in x), M_f^λ is convex. This would be the case even if f was not convex: we just need M_f^λ to be properly defined (i.e. the inf in (8.22) to exist).*

Example 8.10. *As usual with prox-related concepts, it is useful to understand what we obtain for $f = \iota_{\mathcal{C}}$, with $\mathcal{C} \subset \mathbb{R}^d$ nonempty closed and convex. Then,*

$$M_f^\lambda(x) = \inf_{y \in \mathbb{R}^d} \frac{1}{2\lambda} \|y - x\|^2 + \iota_{\mathcal{C}}(y) \quad (8.23)$$

$$= \inf_{y \in \mathcal{C}} \frac{1}{2\lambda} \|y - x\|^2 \quad (8.24)$$

$$= \frac{1}{2\lambda} \left(\inf_{y \in \mathcal{C}} \|y - x\| \right)^2 \quad (8.25)$$

$$= \frac{1}{2\lambda} \text{dist}(x, \mathcal{C})^2 \quad (8.26)$$

$$(8.27)$$

where $\text{dist}(x, \mathcal{C})$ is the distance from x to \mathcal{C} , well-defined because \mathcal{C} was chosen nonempty closed and convex.

Example 8.11. In dimension 1, for $\tau > 0$, $\mathcal{C} = [-\tau, \tau]$ and $f = \iota_{\mathcal{C}}$,

$$M_f^1(x) = \frac{1}{2}d_{[-\tau, \tau]}(x)^2 = \begin{cases} \frac{1}{2}(x + \tau)^2, & \text{if } x \leq -\tau \\ 0, & \text{if } -x \leq \tau \\ \frac{1}{2}(x - \tau)^2, & \text{if } \tau \leq x. \end{cases} \quad (8.28)$$

Proposition 8.12 (Gradient and smoothness of Moreau envelope). *Let $f \in \Gamma_0(\mathbb{R}^d)$, let $\lambda > 0$. The Moreau envelope M_f^λ is differentiable, with:*

$$\nabla M_f^\lambda(x) = \frac{x - \text{prox}_{\lambda f}(x)}{\lambda} \quad (8.29)$$

[Proposition 8.12](#) is an application of the following more general theorem, called Danskin's theorem or envelope theorem.

Theorem 8.13 (Danskin/envelope theorem). *Check conditions rigorously, in particular differentiability* Let \mathcal{L} be function of two variables x and y , differentiable in x , such that for every x , $\mathcal{L}(x, \cdot)$ has a unique minimizer, denoted $y^*(x)$. Denote

$$\ell^*(x) = \min_y \mathcal{L}(x, y) = \mathcal{L}(x, y^*(x)) \quad (8.30)$$

Then, $\ell^*(x)$ is differentiable, and

$$\nabla \ell^*(x) = \nabla_x \mathcal{L}(x, y^*(x)) \quad (8.31)$$

The proof of [Proposition 8.12](#) is just [Theorem 8.13](#) applied to $\mathcal{L}(x, y) = \frac{1}{2\lambda}\|x - y\|^2 + f(y)$, which is differentiable in x , for which $y^*(x) = \text{prox}_{\lambda f}(x)$ and $\nabla_x \mathcal{L}(x, y) = \frac{x - y}{\lambda}$.

Remark 8.14. [Proposition 8.12](#) sheds an interesting light on the proximal point algorithm ([Section 7.5](#)) with stepsize $\lambda > 0$ applied to $f \in \Gamma_0(\mathbb{R}^d)$. The iterations of this algorithm are given by:

$$x_{k+1} = \text{prox}_{\lambda f}(x_k) \quad (8.32)$$

In light of [Proposition 8.12](#), this rewrites:

$$x_{k+1} = x_k - \lambda \nabla M_f^\lambda(x_k), \quad (8.33)$$

which means that the proximal point algorithm on f with stepsize λ is the same as performing gradient descent on the $1/\lambda$ -smooth function M_f^λ !

Example 8.15 (Gradient of squared distance and sitance). Let $\mathcal{C} \subset \mathbb{R}^d$ be nonempty closed and convex. Let $\text{dist}(\cdot, \mathcal{C})$ denote the distance function to \mathcal{C} . Then $\frac{1}{2} \text{dist}(x, \mathcal{C})^2$ is differentiable, and

$$\nabla \frac{1}{2} \text{dist}(x, \mathcal{C})^2 = x - \text{proj}(x, \mathcal{C}) \quad (8.34)$$

Note: as a sanity check, we see that if x is in the interior of \mathcal{C} , the gradient is 0.

In addition, $\text{dist}(\cdot, \mathcal{C})$ is differentiable at points for which $\text{dist}(x, \mathcal{C}) \neq 0$, for which

$$\nabla \text{dist}(x, \mathcal{C}) = \frac{x - \text{proj}(x, \mathcal{C})}{\|x - \text{proj}(x, \mathcal{C})\|} = \frac{x - \text{proj}(x, \mathcal{C})}{\text{dist}(x, \mathcal{C})} \quad (8.35)$$

Proof. Equation (8.34) uses the fact that $\frac{1}{2} \text{dist}(\cdot, \mathcal{C})^2$ is the Moreau envelope of $\iota_{\mathcal{C}}$ (Example 8.10), together with Proposition 8.12.

The second result (Equation (8.35)) is the chain rule applied to the composition of $\sqrt{2} \cdot$ and $\frac{1}{2} \text{dist}(\cdot, \mathcal{C})^2$. \square

Example 8.16. *If we apply the gradient of the Moreau envelope formula to Example 8.11, we get as gradient:*

$$M_f^1(x) = \frac{1}{2} d_{[-\tau, \tau]}(x)^2 = \begin{cases} x + \tau, & \text{if } x \leq -\tau \\ 0, & \text{if } -\tau \leq x \leq \tau \\ x - \tau, & \text{if } x \geq \tau \end{cases} \quad (8.36)$$

Interestingly, we find that is exactly the soft-thresholding, i.e., the gradient of the Moreau envelope is a prox! **TODO** This is not a coincidence, there are deep connections.

We conclude with one very nice property of the Moreau envelope: its smoothness.

Proposition 8.17. *The Moreau envelope is $1/\lambda$ -smooth.*

Proof. • One proof relies on the direct expression of the gradient (Proposition 8.12): denoting $p_x = \text{prox}_f(x)$, $p_y = \text{prox}_f(y)$

$$\|x - p_x - y + p_y\|^2 = \|x - y\|^2 - 2\langle x - y, p_x - p_y \rangle + \|p_x - p_y\|^2 \quad (8.37)$$

$$= \|x - y\|^2 - \|p_x - p_y\|^2 \quad (8.38)$$

$$= \|x - y\|^2 \quad (8.39)$$

using $-2\langle x - y, p_x - p_y \rangle \leq -2\|p_x - p_y\|^2$ (firm nonexpansiveness of the prox, Equation (7.17)). Using $\nabla M_f^\lambda(x) - \nabla M_f^\lambda(y) = \frac{1}{\lambda} \|x - p_x - y + p_y\|$ concludes.

• The other proof is completely different and extremely beautiful in my opinion. f is convex, so $f + \frac{\lambda}{2} \|\cdot\|^2$ is λ -strongly convex, so $(f + \frac{\lambda}{2} \|\cdot\|^2)^*$ is $1/\lambda$ -smooth (Proposition 8.6). But $(f + \frac{\lambda}{2} \|\cdot\|^2)^* = f^* \square (\frac{\lambda}{2} \|\cdot\|^2)^* = f^* \square \frac{1}{2\lambda} \|\cdot\|^2 = M_f^\lambda$! \square

8.3 Fenchel-Rockafellar duality

Many problems in Machine Learning have the following form:

$$\min_{x \in \mathbb{R}^d} f(Ax) + g(x) = P(x) \quad , \quad (8.40)$$

with $A \in \mathbb{R}^{n \times d}$, $f \in \Gamma_0(\mathbb{R}^n)$ and $g \in \Gamma_0(\mathbb{R}^d)$. In this section, we will study the link between such a problem, called the primal, and a closely related one called its *dual problem*.

Definition 8.18. *The Fenchel-Rockafellar dual problem of Problem (8.40) is:*

$$\min_{y \in \mathbb{R}^n} f^*(y) + g^*(-A^\top y) = D(y) \quad . \quad (8.41)$$

Note: it is often alternatively written as a concave maximization problem

$$\max_{y \in \mathbb{R}^n} -f^*(y) - g^*(-A^\top y) \quad . \quad (8.42)$$

Where does this problem come from? It can be obtained by Lagrangian duality (Section 8.6):

$$\inf_{x \in \mathbb{R}^d} f(Ax) + g(x) = \inf_{x \in \mathbb{R}^d, z \in \mathbb{R}^n} f(z) + g(x) \quad \text{s.t. } Ax = z \quad (8.43)$$

$$= \inf_{x \in \mathbb{R}^d, z \in \mathbb{R}^n} f(z) + g(x) + \iota_{\{z\}}(Ax) \quad (8.44)$$

$$= \inf_{x \in \mathbb{R}^d, z \in \mathbb{R}^n} f(z) + g(x) + \sup_{y \in \mathbb{R}^n} \langle y, Ax - z \rangle \quad (8.45)$$

$$= \inf_{x \in \mathbb{R}^d, z \in \mathbb{R}^n} \sup_{y \in \mathbb{R}^n} f(z) + g(x) + \langle y, Ax - z \rangle \quad (8.46)$$

Let's swap inf and sup. By doing so, we lose the equivalence ($\sup \inf \leq \inf \sup$, without equality in general). We have:

$$\sup_{y \in \mathbb{R}^n} \inf_{x \in \mathbb{R}^d, z \in \mathbb{R}^n} f(z) + g(x) + \langle y, Ax - z \rangle \quad (8.47)$$

$$= \sup_{y \in \mathbb{R}^n} \inf_{x \in \mathbb{R}^d} g(x) + \langle y, Ax \rangle + \inf_{z \in \mathbb{R}^n} f(z) - \langle y, z \rangle \quad (8.48)$$

$$= \sup_{y \in \mathbb{R}^n} - \sup_{x \in \mathbb{R}^d} -g(x) + \langle -A^\top y, x \rangle - \sup_{z \in \mathbb{R}^n} -f(z) + \langle y, z \rangle \quad (8.49)$$

$$= \sup_{y \in \mathbb{R}^n} -g^*(-A^\top y) - f^*(y) , \quad (8.50)$$

that is our dual problem.

Proposition 8.19 (Weak duality). *For all $x, y \in \mathbb{R}^d \times \mathbb{R}^n$, $P(x) + D(y) \geq 0$.*

Proof. Write the sum, group f with f^* and g with g^* , then use the Fenchel-Young inequality on each of these pairs. \square

Proposition 8.20. *Let $x^*, y^* \in \mathbb{R}^d \times \mathbb{R}^n$. The following conditions are equivalent:*

1. $x^* \in \operatorname{argmin}_x P(x)$, $y^* \in \operatorname{argmin}_y D(y)$, $\inf P + \inf D = 0$.
2. $P(x^*) + D(y^*) = 0$
3. $f(Ax^*) + f^*(y^*) = \langle Ax^*, y^* \rangle$ and $g(x^*) + g^*(-A^\top y^*) = -\langle x^*, A^\top y^* \rangle$,
4. $-A^\top y^* \in \partial g(x^*)$ and $y^* \in \partial f(Ax^*)$
5. $x^* \in \partial g^*(-A^\top y^*)$ and $Ax^* \in \partial f^*(y^*)$

Note that 4 rewrites $0 \in A^ \partial f(Ax^*) + \partial g(x^*)$, and similarly 5 rewrites $0 \in \partial f^*(y^*) - A \partial g^*(-A^\top y^*)$.*

Proof. 1 same as 2 is OK.

2 same as 3: Do the same as in the proof of Proposition 8.19, and use that the sum of two positive terms is 0, if and only if both terms are 0.

3 same as 4: equality in Fenchel-Young case

4 same as 5: Proposition 8.4 \square

Definition 8.21 (Saddle point of the Lagrangian). *The pair $x^*, y^* \in \mathbb{R}^d \times \mathbb{R}^n$ is said to be a saddle-point of the Lagrangian $\mathcal{L}(x, y) = f(x) + \langle Ax, y \rangle - g^*(y)$ if:*

$$\forall x, y \in \mathbb{R}^d \times \mathbb{R}^n, \mathcal{L}(x^*, y) \leq \mathcal{L}(x^*, y^*) \leq \mathcal{L}(x, y^*) . \quad (8.51)$$

Proposition 8.22. *If $x^*, y^* \in \mathbb{R}^d \times \mathbb{R}^n$ is a saddle-point of the Lagrangian, then x^* solves the primal and y^* solves the dual.*

Proof. Write the two inequalities stating that the pair is a saddle point, make subgradients appear, and use [Proposition 8.20](#). \square

Definition 8.23 (Qualification condition). *The following inclusion is called a qualification condition:*

$$0 \in \text{int}(\text{dom } f - A \text{ dom } g) . \quad (8.52)$$

Where does it come from? It can be seen as a stronger version of “the primal objective P is proper”:

$$\begin{aligned} \text{dom } P \neq \emptyset &\Leftrightarrow \exists x \in \mathbb{R}^d, f(Ax) + g(x) < +\infty \\ &\Leftrightarrow \exists x \in \mathbb{R}^d, f(Ax) < +\infty, g(x) < +\infty \\ &\Leftrightarrow \exists x \in \mathbb{R}^d, y \in \mathbb{R}^n, y = Ax, f(y) < +\infty, g(x) < +\infty \\ &\Leftrightarrow \exists x \in \mathbb{R}^d, y \in \mathbb{R}^n, y = Ax, y \in \text{dom } f, x \in \text{dom } g \\ &\Leftrightarrow \exists x \in \text{dom } g, y \in \text{dom } f, 0 = y - Ax \\ &\Leftrightarrow 0 \in \text{dom } f - A \text{ dom } g . \end{aligned} \quad (8.53)$$

Remark 8.24. *When does QC hold? One often used property is that it holds when f is continuous at Ax for some $x \in \text{dom } g$. Indeed, by continuity, this means that there exists some δ such that with B_δ the ball of center 0 and radius δ , $Ax + B_\delta \subset \text{dom } f$. Then $B_\delta \subset \text{dom } f - Ax \subset \text{dom } f - A \text{ dom } g$ and so QC holds.*

Proposition 8.25 (Strong duality). *Suppose that QC holds. Then the dual has a solution and $\inf P + \inf D = 0$.*

Proof. The case where $\inf P = -\infty$ is easy to address. Let us now consider the case $\inf P > -\infty$. Define the function:

$$h : y \mapsto \inf_x f(Ax + y) + g(x) . \quad (8.54)$$

A few observations about this function: it does not take value $-\infty$; it is convex, its value at 0 is $h(0) = \inf P(x)$ and its domain is $\text{dom } f - A \text{ dom } g$, so QC translates into $0 \in \text{int dom } h$.

Proof: h does not take value $-\infty$: observe that by QC f is finite in a neighborhood of 0; hence $h(-\epsilon y)$ is finite for ϵ small enough. Then by convexity, since $0 = -\frac{1}{1+\epsilon}\epsilon y + (1 - \frac{1}{1+\epsilon})y$, $h(0) \leq -\frac{1}{1+\epsilon}(\epsilon y) + (1 - \frac{1}{1+\epsilon})h(y)$ and so $h(y)$ cannot be $-\infty$.

Since $0 \in \text{int dom } h$ and h is convex, the subdifferential of h at 0 is non-empty. Let $v \in \partial h(0)$. For all $y \in \mathbb{R}^d$, $h(y) \geq h(0) + \langle v, y \rangle$. Reordering, using $h(0) = \inf P$ and the definition of h as an infimum, we get that for all x, y :

$$\inf P \leq g(x) + f(Ax + y) - \langle v, y \rangle \quad (8.55)$$

$$\leq g(x) + f(Ax + y) - \langle x, -A^*v \rangle - \langle v, Ax + y \rangle . \quad (8.56)$$

Taking the infimum of the right-hand side with respect to y , then with respect to x yields:

$$\inf P \leq -f^*(v) - g^*(-A^*v) , \quad (8.57)$$

so the duality gap is 0 and v is a solution of the dual. \square

Remark 8.26. By [Proposition 8.25](#), one has that if QC holds,

$$0 \in \partial(f \circ A + g)(x) \iff 0 \in A^\top \partial f(Ax) + \partial g(x) \quad (8.58)$$

, because *TODO*.

8.4 A calculus rule for subdifferentials

So far we have only seen that $\partial f_1 + \partial f_2 \subset \partial(f_1 + f_2)$, and the inclusion can be strict ([Exercises 8.9](#) and [8.10](#)). A surprising consequence of the duality result derived in [Proposition 8.25](#) is that it allows proving that under QC, the subdifferential of the sum is the sum of subdifferentials.

Proposition 8.27. *Let QC hold. Then for all x ,*

$$\partial(f \circ A + g)(x) = A^\top \partial f(Ax) + \partial g(x) \quad (8.59)$$

Proof. We know that when QC holds, $0 \in \partial(f \circ A + g)(x) \iff 0 \in A^\top \partial f(Ax) + \partial g(x)$.

The proof uses [Exercise 7.2](#) to reduce to that case. One has

$$u \in \partial(f \circ A + g)(x) \iff 0 \in \partial(f \circ A + g)(x) - \{u\} \quad (8.60)$$

$$\iff 0 \in \partial(f \circ A + g - \langle u, \cdot \rangle)(x) \quad (8.61)$$

$$\iff 0 \in A^\top \partial f(Ax) + \partial(g - \langle u, \cdot \rangle)(x) \quad (8.62)$$

$$\iff 0 \in A^\top \partial f(Ax) + \partial(g)(x) - \{u\} \quad (8.63)$$

$$\iff u \in A^\top \partial f(Ax) + \partial(g)(x) \quad (8.64)$$

where we have used that QC holds for $(f, g - \langle u, \cdot \rangle)$ because g and $g - \langle u, \cdot \rangle$ have same domain. \square

8.5 Primal-dual algorithms

In this section, we introduce more algorithms to deal with composite problems, based on duality/splitting.

Chambolle-Pock algorithm The Chambolle-Pock algorithm solves problems of the form:

$$\min_{x \in \mathbb{R}^d} f(Ax) + g(x), \quad (8.65)$$

where $f \in \Gamma_0(\mathbb{R}^n)$, $g \in \Gamma_0(\mathbb{R}^d)$, but both need not be smooth. It requires access to their proximal operators. Let σ and τ be positive stepsizes such that $\sigma\tau\|A\|_2^2 < 1$. Let $\theta \in]0, 1]$. The Chambolle-Pock iterations write:

$$\begin{cases} x_{k+1} = \text{prox}_{\tau g}(x_k - \tau A^\top \bar{y}_k) \\ y_{k+1} = \text{prox}_{\sigma f^*}(y_k + \sigma A[(1 + \theta)x_{k+1} - \theta x_k]) \end{cases} \quad (8.66)$$

Note: The interpolation parameter θ is most often taken equal to 1. The interpolation step can be performed on y instead of x , which yields a different version of the algorithm.

Proposition 8.28. *For $\theta = 1$, the Chambolle-Pock algorithm is an instance of the preconditioned proximal point algorithm, in the space $\mathbb{R}^d \times \mathbb{R}^n$, to find a 0 of a suitable operator.*

Proof. The optimality conditions from the definitions of x_{k+1} and y_{k+1} rewrite:

$$\frac{1}{\tau}x_k - \frac{1}{\tau}x_{k+1} - A^\top y_k \in \partial g(x_{k+1}) \quad (8.67)$$

$$\frac{1}{\sigma}y_k - \frac{1}{\sigma}y_{k+1} + 2Ax_{k+1} - Ax_k \in \partial f^*(y_{k+1}) \quad (8.68)$$

In turn, this writes:

$$\begin{pmatrix} \frac{1}{\tau}\text{Id} & -A^\top \\ -A & \frac{1}{\sigma}\text{Id} \end{pmatrix} \begin{pmatrix} x_k \\ y_k \end{pmatrix} \in \left[\begin{pmatrix} \frac{1}{\tau}\text{Id} & -A^\top \\ -A & \frac{1}{\sigma}\text{Id} \end{pmatrix} + \begin{pmatrix} \partial g & A^\top \\ A^\top & \partial f^* \end{pmatrix} \right] \begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} \quad (8.69)$$

Thus for $v_k = \begin{pmatrix} x_k \\ y_k \end{pmatrix}$, $M = \begin{pmatrix} \frac{1}{\tau}\text{Id} & -A^\top \\ -A & \frac{1}{\sigma}\text{Id} \end{pmatrix}$ and $T = \begin{pmatrix} \partial g & A^\top \\ A^\top & \partial f^* \end{pmatrix}$ this writes: $v_{k+1} = (\text{Id} + M^{-1}T)^{-1}v_k$, which is the proximal point with different metric. \square

Vu-Condat algorithm [Vũ \(2013\)](#) and [Condat \(2013\)](#) independently proposed an algorithm that handles an additional smooth term h in the objective:

$$\min_{x \in \mathbb{R}^d} f(Ax) + g(x) + h(x), \quad (8.70)$$

where $f \in \Gamma_0(\mathbb{R}^n)$, $g \in \Gamma_0(\mathbb{R}^d)$, $h \in \Gamma_0(\mathbb{R}^d)$, f and g 's proximal operators are known, and h is L_h -smooth¹⁰.

For stepsizes τ and σ satisfying¹¹:

$$\frac{\tau\sigma}{\|A\|^2} + \frac{\tau L_h}{2} \leq 1, \quad (8.71)$$

the Condat-Vu algorithm reads:

$$\begin{cases} \bar{y}_k = y_k + \theta(y_k - y_{k-1}) \\ x_{k+1} = \text{prox}_{\tau g}(x_k - \tau \nabla h(x_k) - \tau A^\top \bar{y}_k), \\ y_{k+1} = \text{prox}_{\sigma f^*}(y_k + \sigma Ax_k). \end{cases} \quad (8.72)$$

Note: it can even handle a fourth smooth term, so that the saddle point formulation becomes TODO

¹⁰for completeness let us mention that it can even handle the slightly more complex case where f is replaced by $f \square F$ with F strongly convex, leading to a second smooth term, F^* , in the saddle point formulation.

¹¹this condition is in the case $\beta > 0$; otherwise inequality must be strict

ADMM The Alternating Direction Method of Multipliers (ADMM method) is designed to solve problems of the form:

$$\min_{x \in \mathbb{R}^d, z \in \mathbb{R}^n} f(z) + g(x) \quad \text{s.t.} \quad Ax + Bz = c. \quad (8.73)$$

For simplicity and to match the setup of Chambolle-Pock we take $B = -\text{Id}$ and $c = 0$.

The augmented Lagrangian of parameter ρ is:

$$\mathcal{L}^\rho(x, y, z) = f(z) + g(x) + \langle y, Ax - z \rangle + \frac{\rho}{2} \|Ax - z\|^2. \quad (8.74)$$

Note that some formulations use a *scaled* dual variable $y' = \rho y$, and rewrite the last two terms as $\frac{\rho}{2} \|Ax - z + y'\| - \frac{\rho}{2} \|y'\|^2$.

The ADMM algorithm reads:

$$\begin{cases} x_{k+1} \in \operatorname{argmin}_x \mathcal{L}^\rho(x, y_k, z_k) \\ z_{k+1} \in \operatorname{argmin}_z \mathcal{L}^\rho(x_{k+1}, y_k, z) \\ y_{k+1} = y_k + r(Ax_{k+1} - z_{k+1}) \end{cases} \quad (8.75)$$

Douglas-Rachford algorithm

$$x_{k+1} = x_k + \operatorname{prox}_g(2 \operatorname{prox}_f(x_k) - x_k) - \operatorname{prox}_f(x_k) \quad (8.76)$$

Application of the operator $\frac{1}{2}\text{Id} + \frac{1}{2}(2 \operatorname{prox}_g - \text{Id})(2 \operatorname{prox}_f - \text{Id})$.

TODO ADMM is the Douglas Rachford algorithm.

8.6 Lagrangian duality and KKT conditions

Fenchel-Rockafellar duality ([Section 8.3](#)) is a particular case of the more versatile *Lagrangian* duality. The latter is concerned with convex problems with affine constraints, that is problems of the form

$$\min_{x \in \mathbb{R}^d} f(x) \quad \text{s.t.} \quad g_i(x) \leq 0 \quad \forall i \in [m], \quad h_i(x) = 0 \quad \forall i \in [p] \quad (8.77)$$

We write the inequality constraints as ≤ 0 , because we like working with convex functions and convex sets; if g is convex its sublevel sets, and in particular $\{x : g(x) \leq 0\}$, are convex – this would not be true for $\{x : g(x) \geq 0\}$.

Definition 8.29 (Lagrangian, dual function, Lagrangian dual). *The Lagrangian associated with [Problem \(8.77\)](#) is the function*

$$\begin{aligned} \mathcal{L} : \mathbb{R}^n \times \mathbb{R}^p \times \mathbb{R}_+^m &\rightarrow \mathbb{R} \\ (x, \lambda, \mu) &\mapsto f(x) + \sum_{i=1}^n \lambda_i h_i(x) + \sum_{i=1}^p \mu_i g_i(x), \end{aligned} \quad (8.78)$$

where λ and μ are called the Lagrange multipliers. Note that the multipliers μ_i associated with the inequality constraints are positive; the rationale is that, as soon as x does not

satisfy one constraint, maximizing \mathcal{L} in λ_i or μ_i will make the Lagrangian go to $+\infty$, so that [Problem \(8.77\)](#) is equivalent to:

$$\min_{x \in \mathbb{R}^n} \max_{\lambda \in \mathbb{R}^p, \mu \in \mathbb{R}_+^m} \mathcal{L}(x, \lambda, \mu). \quad (8.79)$$

Also note that for any x feasible for [Problem \(8.77\)](#), the Lagrangian value for any $\lambda, \mu \geq 0$ upper bounds $f(x)$.

Finally, the dual function is:

$$\begin{aligned} g : \mathbb{R}^p \times \mathbb{R}_+^m &\rightarrow \mathbb{R} \\ (\lambda, \mu) &\mapsto \min_{x \in \mathbb{R}^n} \mathcal{L}(x, \lambda, \mu) \end{aligned} \quad (8.80)$$

and the Lagrange dual problem of [Problem \(8.77\)](#) is:

$$\max_{\lambda \in \mathbb{R}^p, \mu \in \mathbb{R}_+^m} g(\lambda, \mu), \quad (8.81)$$

that is, [Problem \(8.79\)](#) but with the order of the min and the max inverted.

Proposition 8.30 (Weak duality). *Without any assumption on the objective f nor on the constraints g_i and h_i , the dual optimal value is smaller than the primal optimal value (we say that weak duality holds).*

Proof. Let $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$, let $a_0, b \in \mathbb{R}^n \times \mathbb{R}^m$. Then:

$$\inf_{a \in \mathbb{R}^n} F(a, b) \leq F(a_0, b) \quad (8.82)$$

$$\sup_{b \in \mathbb{R}^m} \inf_{a \in \mathbb{R}^n} F(a, b) \leq \sup_{b \in \mathbb{R}^m} F(a_0, b) \quad (8.83)$$

The left-hand side is a constant, so the infimum of the right-hand side with respect to a_0 must be greater:

$$\sup_{b \in \mathbb{R}^m} \inf_{a \in \mathbb{R}^n} F(a, b) \leq \inf_{a \in \mathbb{R}^n} \sup_{b \in \mathbb{R}^m} F(a, b). \quad (8.84)$$

So the supremum of the infimums is smaller than the inf of the sups, which concludes. \square

The equality is not tight in general, but under some conditions, we can have strong duality, meaning that the primal and dual values are equal. There are many such conditions; one of the most popular is Slater's condition.

Proposition 8.31 (Slater's condition). *Let f be convex, let the equality constraint functions h_i be affine and let the inequality functions g_i be convex.*

If [Problem \(8.77\)](#) satisfies the so-called Slater's condition, that is,

$$\exists x \in \mathbb{R}^n, g_i(x) < 0 \ \forall i \in [m], \quad h_i(x) = 0 \ \forall i \in [p], \quad (8.85)$$

(we say the problem is strictly feasible), then strongly duality holds: the primal and dual optimal values are equal.

We now introduce a very powerful tool to solve constrained optimization problems: the Karush-Kuhn-Tucker (KKT) conditions.

Definition 8.32. The triplet (x, λ, μ) is said to satisfy the KKT conditions for [Problem \(8.77\)](#) if it is a saddle point for the Lagrangian, meaning that it satisfies:

$$\nabla_x \mathcal{L}(x, \lambda, \mu) = 0 \quad (\text{stationarity}) \quad (8.86)$$

$$g_i(x) \leq 0 \quad \forall i \in [m] \quad (\text{primal feasibility}) \quad (8.87)$$

$$h_i(x) = 0 \quad \forall i \in [p] \quad (\text{primal feasibility}) \quad (8.88)$$

$$\mu \geq 0 \quad (\text{dual feasibility}) \quad (8.89)$$

$$\mu_i g_i(x) = 0 \quad \forall i \in [m] \quad (\text{complementary slackness}) \quad (8.90)$$

For the sake of clarity, we have used a gradient for stationarity (thus implying that everything is differentiable); but the same results exist with subgradients. Actually, the complementary slackness together with dual feasibility is easier understood as $\mu_i \in \partial_{\mathbb{R}^-}(g_i(x))$

For many simple problems, it is easier to solve the above system. Under reasonable conditions, solutions of the minimization problem are amongst KKT points.

Definition 8.33 (Active set and LICQ). The active set at $x \in \mathcal{D}$, $\mathcal{A}(x)$, is the set of i 's such that inequality constraint i is active: $g_i(x) = 0$. The Linear Independence Constraint Qualification (LICQ) at x holds if $(\nabla h_i(x)) \cup (\nabla g_i(x), i \in \mathcal{A}(x))$ is an independent family (gradients of all equality constraints, gradient of active inequality constraints).

Proposition 8.34 (KKT are necessary under LICQ). If x is a local minimizer for [Problem \(8.77\)](#) and $\text{LICQ}(x)$ holds, then the KKT hold.

The proof is not easy (see Nocedal and Wright, chapter 12)

Example 8.35 (Without LICQ, KKT are not necessary). Consider $\min x \text{ s.t. } x^2 \leq 0$. The minimizer is 0. This point does not satisfy the LICQ, as for $x = 0$ the (only) constraint saturates and its gradient is 0.

And the KKT do not hold, as $\mathcal{L}(x, \mu) = x + \mu x^2$, and $\nabla_x \mathcal{L}(0, \mu) = 1 \neq 0$ for all μ .

Proposition 8.36 (KKT are sufficient in the convex + affine equality case). Let the h_i 's be affine. Let f and the g_i 's all be convex (and differentiable for easier proof, but everything works with subgradients). Then the KKT are sufficient: if they hold at x^*, λ^*, μ^* , then x^* is a global minimizer of the constrained problem.

Proof. By the requirement on f , g_i , the h_i and the fact that $\mu^* \geq 0$, the Lagrangian $\mathcal{L}(\cdot, \lambda^*, \mu^*)$ is convex (because we do not know the sign of λ^* , we must require the h_i 's to be affine for $\lambda_i h_i$ to be convex).

Because of this convexity, the KKT condition $\nabla_x \mathcal{L}(\cdot, \lambda^*, \mu^*) = 0$ implies that x^* is a minimizer of $\mathcal{L}(\cdot, \lambda^*, \mu^*)$. Therefore for any x , $\mathcal{L}(x, \lambda^*, \mu^*) \geq \mathcal{L}(x^*, \lambda^*, \mu^*) = f(x^*)$ (because of feasibility of x^* and complementary slackness).

Now if we take x feasible, the LHS is $f(x) + \sum_i \mu_i^* g_i(x)$, which is smaller than $f(x)$ ($\mu^* \geq 0$ and x feasible so $g_i(x) \leq 0$). So $f(x) \geq f(x^*)$ for any feasible x , which concludes the proof. \square

8.7 Exercises

Exercise 8.1. ☹ Show that the Fenchel transform is order-reversing: if $f \leq g$, $g^* \leq f^*$.

Exercise 8.2. ☹☹ Compute the Fenchel transform of the ℓ_p -norm for $p \in [1, +\infty]$.

Exercise 8.3 (Dual norm). ☹☹ For a norm $\|\cdot\|$, define its dual norm $\|\cdot\|_*$ (note: do not confuse it with the Fenchel transform of the norm!) as:

$$\|x\|_* = \sup_{\|y\| \leq 1} x^\top y \quad (8.91)$$

Show that the dual norm of the ℓ_p norm is the ℓ_q norm, where p and q are Hölder conjugates. Show that the Fenchel conjugate of $\frac{1}{2}\|\cdot\|^2$ is $\frac{1}{2}\|\cdot\|_*^2$.

Show that the Fenchel transform of a norm is the indicator function of the unit ball of its dual norm.

☹☹☹ Show that $\partial\|x\| = \{u \in \mathbb{R}^d : x^\top u = \|x\|, \|u\|_* \leq 1\}$. Check what you obtain for the ℓ_1 norm.

Exercise 8.4. ☹☹ Show that $x \mapsto \frac{1}{2}\|x\|^2$ is a fixed point of the Fenchel transform. Show that it is the only fixed point of the Fenchel transform.

Exercise 8.5. ☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be separable: $f(x) = \sum_1^d f_i(x_i)$ where the f_i 's are functions of the real variable.

Show that $f^*(u) = (f_i^*(u_i))_{i \in [d]}$.

Exercise 8.6 (“Lipschitzing trick”). ☹☹ Let $f \in \Gamma_0(\mathbb{R}^d)$.

Show that f is M -Lipschitz if and only if the domain of f^* is included in the Euclidean ball of center 0 and radius M .

Exercise 8.7 (Convolution smoothing). ☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex.

Show that $f \square \frac{L}{2}\|\cdot\|^2$ is L -smooth.

Exercise 8.8 (Convolution Lipschitzing). ☹ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex.

Show that $f \square L\|\cdot\|$ is L -Lipschitz.

Exercise 8.9 (A case of strict inclusion in subdifferential of sums). ☹ or ☹☹ depending on whether or not you know the result of [Exercise 7.6](#)

In \mathbb{R}^2 , let D_1 be the closed disk of center $(-1, 0)$ and radius 1, and D_2 be the closed disk of center $(0, 1)$ and radius 1. Compute $\partial(\iota_{D_1} + \iota_{D_2})(0, 0)$ and compare to $\partial\iota_{D_1}(0, 0) + \partial\iota_{D_2}(0, 0)$. Does the qualification condition hold for the two functions under consideration?

Exercise 8.10 (Another case of strict inclusion in subdifferential of sums). ☹ Let $f : \mathbb{R} \rightarrow \bar{\mathbb{R}}$ be defined as $f(x) = +\infty$ for $x < 0$, and $f(x) = -\sqrt{x}$ otherwise. Let $g : x \mapsto f(-x)$. Compare $\partial(f + g)(0)$ with $\partial f(0) + \partial g(0)$. Does the qualification condition hold for the two functions under consideration?

Exercise 8.11 (A case without strong duality). ☹️ Let

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$x \mapsto \begin{cases} x \log x - x & \text{if } x > 0 \text{ ,} \\ 0 & \text{if } x = 0 \text{ ,} \\ +\infty & \text{otherwise .} \end{cases} \quad (8.92)$$

Show that $\min_{x \in \mathbb{R}^d} f(x) + f(-x)$ has solutions. What about the dual problem?

Exercise 8.12. Let $\mathcal{D} = \mathbb{R} \times \mathbb{R}_+^*$, and $f : \mathcal{D} \rightarrow \mathbb{R}$ defined by $f(x, y) = e^{-x}$. Show that the problem

$$\min_{x, y \in \mathcal{D}} f(x, y) \quad \text{s.t.} \quad \frac{x^2}{y} \leq 0 \quad (8.93)$$

is convex. Compute its optimal value. Compute the Lagrangian dual and its optimal value. Does Slater's condition hold for this problem?

9 Convergence through fixed-point iterations

References: 40 pages “gentle and self-contained introduction and tutorial” by [Ryu and Boyd \(2016\)](#), book by [Ryu and Yin \(2022\)](#). See also the very extensive catalog of applications in [Combettes and Pesquet \(2021\)](#).

Splitting methods can be analyzed in an elegant and unified way through the lens of (monotone) operators and fixed point iterations.

9.1 Picard iterations

Let us start by citing the best-known result in fixed point theory, called the Banach (or Banach-Cacciopoli or *contraction mapping*) theorem. It applies to *contractions*.

Definition 9.1 (Contraction). *For a normed space¹² X , a contraction or contractive mapping is an operator $T : X \rightarrow X$, which has Lipschitz constant strictly less than 1, i.e.*

$$\exists q \in [0, 1[, \forall (x, y) \in X^2, \|T(x) - T(y)\| \leq q\|x - y\|. \quad (9.1)$$

Theorem 9.2 (Banach fixed-point). *Let X be a complete space and T a contractive mapping. Then T admits exactly one fixed point x^* , i.e. such that $T(x^*) = x^*$. For any $x_0 \in X$, the sequence defined by $x_{k+1} = T(x_k)$ converges to x^* , and the convergence is linear.*

Proof. Show that $T(x_k)$ is Cauchy, thus converges (X is complete), and the limit must be a fixed point. This shows the fixed point is unique, which one can also see using two fixed points a and b : $\|a - b\| = \|T(a) - T(b)\| \leq q\|a - b\|$. \square

Remark 9.3. *The iterations $x_{k+1} = T(x_k)$ are called Picard iterations.*

If the Lipschitz constant of T is 1 instead of being strictly less than 1 (such T is called *nonexpansive*) we cannot say anything (neither about eventual fixed points, nor about convergence of Picard iterations). Two good examples to keep in mind for possible behaviors are $-\text{Id}$ (a fixed point exists, Picard iterations do not converge) and translations $T(x) = x + b$ (no fixed point, no convergence).

It turns out that it is possible to refine this division: there is a subclass of nonexpansive operators, larger than contractions, for which Picard iterations converge¹³, as we will show in [Theorem 9.13](#). Let’s introduce these operators and give essential properties together with examples.

Definition 9.4 (α -averaged operators). *For $\alpha \in]0, 1]$, T is α -averaged if and only if there exists a nonexpansive operator R such that $T = (1 - \alpha)\text{Id} + \alpha R$.*

It is easy to see that such operators are nonexpansive.

Remark 9.5. *Definition 9.4 requires to be in a vector space to have addition and scalar multiplication, while Theorem 9.2 works in a metric space. There seems to be recent work extending the definition of α -averaged operators to more generic spaces ([Berdellima, 2020](#)).*

¹²the definition can straightforwardly be applied to the more general setting of metric spaces

¹³provided a fixed point exists

Proposition 9.6 (Characterization of α -averaged operators). *Let $\alpha \in]0, 1]$. The following are equivalent:*

- (i) T is α -averaged.
- (ii) $(1 - \alpha^{-1})\text{Id} + \alpha^{-1}T$ is nonexpansive.
- (iii) $\|T(x) - T(y)\|^2 \leq \|x - y\|^2 - (\alpha^{-1} - 1)\|(\text{Id} - T)(x) - (\text{Id} - T)(y)\|^2$.
- (iv) $\|T(x) - T(y)\|^2 + (1 - 2\alpha)\|x - y\|^2 \leq 2(1 - \alpha)\langle x - y, T(x) - T(y) \rangle$.

Proof. For (i) \Leftrightarrow (ii): just use the definition.

For (ii) \Leftrightarrow (iii), use the famous equality from [Exercise 2.13](#) ($\alpha^{-1} > 1$ here, but it still holds):

$$\|\lambda a + (1 - \lambda)b\|^2 = \lambda\|a\|^2 + (1 - \lambda)\|b\|^2 - \frac{1}{2}\lambda(1 - \lambda)\|a - b\|^2 \quad (9.2)$$

to get

$$\|x - y\|^2 \geq \|(1 - \alpha^{-1})(x - y) + \alpha^{-1}(T(x) - T(y))\|^2 \quad (9.3)$$

$$\begin{aligned} &= (1 - \alpha^{-1})\|x - y\|^2 \\ &\quad + \alpha^{-1}\|T(x) - T(y)\|^2 - \frac{1}{2}\alpha^{-1}(1 - \alpha^{-1})\|(\text{Id} - T)(x) - (\text{Id} - T)(y)\|^2 \end{aligned} \quad (9.4)$$

Simplify the $\|x - y\|^2$ on both sides and multiply by α .

For (iii) \Leftrightarrow (iv) just develop the last term in (iii). \square

An interesting property of the α -averaged operators is that they are stable by composition – though to be precise, one ends up with a different α for the composition.

Proposition 9.7. *The composition of an α_1 -averaged and an α_2 -averaged operator is $\frac{\alpha_1 + \alpha_2 - 2\alpha_1\alpha_2}{1 - \alpha_1\alpha_2}$ -averaged.*

The proof is computational and can be found in [Ogura and Yamada \(2002, Thm 3\(b\)\)](#).

A particularly interesting case is $\frac{1}{2}$ -averaged operators, that have their own name.

Definition 9.8. *A $\frac{1}{2}$ -averaged operator is called firmly nonexpansive.*

Proposition 9.9 (Characterization of firmly nonexpansive operators). *The following are equivalent:*

- (i) T is firmly nonexpansive
- (ii) $\|T(x) - T(y)\|^2 \leq \|x - y\|^2 - \|(\text{Id} - T)x - (\text{Id} - T)y\|^2$
- (iii) $\|T(x) - T(y)\|^2 \leq \langle T(x) - T(y), x - y \rangle$

Property (iii) is called “cocoercivity” of T , and β -cocoercivity of T means that the equality is satisfied with a positive scalar β multiplying the norm on the left-hand side.

Example 9.10. 1. *Projections and proximal operators (of convex functions) are firmly nonexpansive.*

2. *If f is L -smooth, $\frac{1}{L}\nabla f$ is firmly nonexpansive by the Baillon-Haddad theorem.*

3. *If T is firmly nonexpansive, so is $\text{Id} - T$.*

Lemma 9.11 (Opial lemma). *Let F a non empty subset and (x_k) a sequence such that:*

- (i) $\|x_k - y\|$ converges for all $y \in F$.
- (ii) Every cluster point of (x_k) belongs to F .

Then there exists $\bar{x} \in F$ such that $x_k \rightarrow \bar{x}$. For completeness, note that in infinite dimension, the convergence is only weak convergence.

Proof. First, the sequence (x_k) is bounded since $F \neq \emptyset$. So the sequence has at least one cluster point. Let y_1 and y_2 be two cluster points of the sequence. Extract $x_k^1 \rightarrow y_1$, $x_k^2 \rightarrow y_2$.

Let $y \in F$. $\lim_k \|x_k^1 - y\|^2 = \lim_k \|x_k^2 - y\|^2$ by assumption, so developing and reordering

$$\lim_k \|x_k^1\|^2 - \|x_k^2\|^2 = 2 \lim_k \langle \bar{x}, x_k^1 - x_k^2 \rangle = 2 \langle \bar{x}, y_1 - y_2 \rangle. \quad (9.5)$$

Applying this equality to \bar{x} equal to y_1 and y_2 successively, one gets $\langle y_1, y_2 - y_1 \rangle = \langle y_2, y_2 - y_1 \rangle$ so $y_1 = y_2$ and there is a single cluster point, the bounded sequence must converge. \square

The Opial lemma is frequently applied to Fejér-monotone sequences (see [Combettes \(2001\)](#) for an excellent overview of the properties of such sequences).

Definition 9.12. The sequence (x_k) is Fejér-monotone with respect to the nonempty, closed and convex set S if and only if, for all $\bar{x} \in S$:

$$\|x_{k+1} - \bar{x}\| \leq \|x_k - \bar{x}\|. \quad (9.6)$$

It follows easily that every Fejér monotone sequence is bounded, and thus has weak cluster points. One also has that $\|x_k - \bar{x}\|$ converges (being decreasing and lower bounded) for all $\bar{x} \in S$.

An application of the Opial lemma is the following theorem on the convergence of α -averaged iterations.

Theorem 9.13. Consider the Picard iterations $x_{k+1} = T(x_k)$ with $X_0 \in X$, where T is α -averaged and has **non-empty set of fixed points**.

Then

1. $(d(x_k, \text{Fix } T))$ is decreasing
2. the sequence of squared iterate distance is summable
3. $x_k \rightarrow \bar{x} \in \text{Fix } T$ (as in Opial lemma, convergence in infinite dimension is only weak).

[Theorem 9.13](#) can be used to find fixed points of any nonexpansive operator, using the following fact.

Proposition 9.14. For any $\alpha \neq 0$, R and $(1 - \alpha)\text{Id} + \alpha R$ have the same fixed points.

Proof.

$$(1 - \alpha)x + \alpha R(x) = x \Leftrightarrow \alpha x = \alpha R(x). \quad (9.7)$$

\square

Therefore, to find fixed points of a nonexpansive operator, one can pick $\alpha \in]0, 1[$ and perform Picard iterations on the (by definition) α -averaged operator $T = (1 - \alpha)\text{Id} + \alpha R$. Such iterations are called *Krasnoselskij-Mann iterations* and read:

$$x_{k+1} = x_k + \alpha(R(x_k) - x_k). \quad (9.8)$$

It is possible to take a varying λ when applying Krasnoselskij-Mann iterations. The results are preserved as long as the sequence $(\lambda_n)_n \in]0, 1[$ is such that $\sum_{\mathbb{N}} \lambda_n(1 - \lambda_n) = +\infty$.

9.2 Monotone operators

This section considers set-valued operators, denoted with the double arrow symbol: $T : X \rightrightarrows X$.

Definition 9.15. *An operator $T : X \rightrightarrows X$ is monotone if:*

$$\forall(x, y), \forall(u, v) \in (Tx) \times T(y), \langle u - v, x - y \rangle \geq 0. \quad (9.9)$$

If T is a function from \mathbb{R}^d to \mathbb{R} , it simply means that T is increasing.

Definition 9.16. *A monotone operator T is maximal monotone if there does not exist another monotone operator whose graph strictly contains its graph. Equivalently, this means that for all $x, u \in T(x)$, y and v ,*

$$\langle u - v, x - y \rangle \geq 0 \implies v \in T(y). \quad (9.10)$$

The best illustration of monotone and maximal monotone operators is the subdifferential.

Proposition 9.17. *The subdifferential of any function is a monotone operator. The subdifferential of a Γ_0 function is maximal monotone.*

Proof. The first part almost follows from the definition of the subdifferential, by writing the subgradient lower bound at $u \in \partial f(x)$ at y , and conversely at x for $v \in \partial f(y)$, then summing.

For the second part: let $f \in \Gamma_0(\mathbb{R}^d)$. Let \tilde{x}, \tilde{g} such that $\tilde{x} \notin \partial f(\tilde{x})$. We will show that there exist $x \in \mathbb{R}^d$ and $g \in \partial f(x)$ such that $\langle x - \tilde{x}, g - \tilde{g} \rangle < 0$, meaning that any extension of the graph of f , the resulting operator is not monotone.

Indeed let $x = \arg\min_y f(y) + \frac{1}{2}\|y - \tilde{x} - \tilde{g}\|^2$ – for the well-definedness of x we use the hypothesis $f \in \Gamma_0(\mathbb{R}^d)$. Thus for $g = \tilde{x} + \tilde{g} - x$, we have $g \in \partial f(x)$, and

$$\langle g - \tilde{g}, x - \tilde{x} \rangle = -\|x - \tilde{x}\|^2 = -\|g - \tilde{g}\|^2. \quad (9.11)$$

At least one of these quantities is strictly negative (and so both are), otherwise $(x, g) = (\tilde{x}, \tilde{g})$ and $\tilde{x} \in \partial f(\tilde{g})$. Thus any extension of ∂f cannot be monotone. \square

Definition 9.18 (Resolvent). *For any $T : X \rightrightarrows X$, the set-valued operator $(\text{Id} + T)^{-1}$ is called the resolvent of T .*

Proposition 9.19. *If T is monotone, its resolvent is nonexpansive (and a fortiori, it is single-valued). If T is maximal monotone, the resolvent is defined on the whole space: for any y , the inclusion $y \in x + T(x)$ has a unique solution (if T is only monotone, uniqueness is guaranteed, but not existence).*

Example 9.20. *The resolvent of the subdifferential of a Γ_0 function is simply its proximal operator.*

Theorem 9.21 (Minty's theorem). *An operator is firmly nonexpansive if and only if it is the resolvent of a maximally monotone operator.*

Proposition 9.22 (Examples of maximal monotone operator). *T is maximal monotone in the following cases:*

- T is monotone and continuous
- T is cocoercive
- $\text{Id} - T$ is nonexpansive
- T is linear and positive ($\langle x, T(x) \rangle \geq 0$)
- T is linear and skew, meaning that $T^* = -T$.

9.3 Applications

Example 9.23 (Proximal point algorithm). *Let $f \in \Gamma_0(\mathbb{R})$. By Fermat's rule, to find a minimizer of f is to find a point such that $0 \in \gamma \partial f(x^*)$, meaning a fixed point of $(\text{Id} - \gamma \partial f)^{-1}$. Picard iterations of this algorithm are precisely the proximal point algorithm.*

Example 9.24 (Gradient descent). *Let f be differentiable. A minimizer of f is a fixed point of $\text{Id} - \gamma \nabla f$, which is contractive (TODO only if f is strongly convex, take this setting?) for $\gamma < 2/L$*

Example 9.25 (Forward-backward/proximal gradient algorithm).

Example 9.26 (Douglas-Rachford algorithm).

Example 9.27 (Chambolle-Pock algorithm).

9.4 Exercises

Exercise 9.1. ☹ For $R = -\text{Id}$, check what happens when performing Picard iterations with $(1 - \alpha)\text{Id} + \alpha T$. Compare with what [Theorem 9.13](#) asserts.

Exercise 9.2. ☹ Show that if T is nonexpansive, its set of fixed points is closed and convex.

Exercise 9.3. ☹☹ Show that if T_1 and T_2 are two nonexpansive operators and $\alpha \in]0, 1[$, if $\text{Fix } T_1$ and $\text{Fix } T_2$ intersect, then $\text{Fix}(1 - \alpha)T_1 + \alpha T_2 = \text{Fix } T_1 \cap \text{Fix } T_2$. Generalize to n operators.

Exercise 9.4 (☹☹). Let T_1 and T_2 be α_1 - and α_2 -averaged, such that $\text{Fix } T_1 \cap \text{Fix } T_2 \neq \emptyset$. Show that $\text{Fix } T_1 \circ T_2 = \text{Fix } T_1 \cap \text{Fix } T_2$.

10 Second order optimization: Newton method

Let f be a twice differentiable function. We have seen that gradient descent for stepsize $1/L$ on a L smooth function can be interpreted as a Majorization-Minimization approach, using the descent lemma to upper bound f globally by an isotropic parabola (majorization step), then minimizing this upper bound (minimization step).

The issue with this majorization is that it is very coarse: in dimension 2, take $f(x) = x_1^2 + 0.001x^2$. For this function, $L = 1$ and $\mu = 0.001$: we get a linear convergence rate for gradient descent, but the constant is very poor (see also Figure 1.2), because the upper bounding parabolas that we construct have Hessian $L\text{Id}$, and $\mu \ll L$. In practice, this leads to very small steps being performed after the first one.

10.1 Newton method

We can try to preserve more information about the curvature, by using a 2nd order Taylor expansion of f – hence the name of “second-order methods”:

$$f(y) \approx f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \langle y - x, \nabla^2 f(x)(y - x) \rangle . \quad (10.1)$$

As in the MM interpretation of GD, we can derive an iterative algorithm for this: set $x = x_k$, and define x_{k+1} as the minimizer of the approximation. This yields

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k) . \quad (10.2)$$

Compared to gradient descent, the scalar stepsize has been replaced with the application of a PSD matrix, that better captures the curvature of f .

Let us illustrate the possible behaviors for Newton’s method on a simple 1-dimensional convex function, $x \mapsto \sqrt{1+x^2}$. This function is \mathcal{C}^2 , with $f'(x) = \frac{x}{\sqrt{1+x^2}}$ and $f''(x) = \frac{1}{(1+x^2)^{3/2}}$, so Newton’s method iterates are given by

$$x_{k+1} = x_k - x_k(1 + x_k^2) = -x_k^3 . \quad (10.3)$$

There are therefore, depending on $|x_0|$, three possible behaviors:

- extremely fast convergence towards the minimizer if $|x_0| < 1$
- oscillations between 1 and -1 if $|x_0| = 1$
- extremely fast divergence if $|x_0| > 1$

This example illustrates a crucial fact about the vanilla Newton method: it only has *local* convergence properties – it converges if initialized close enough to the minimizer. But when it does, it enjoys a *superlinear* convergence rate. Let us formalize this and derive a proof.

Proposition 10.1. *If $x^* = \operatorname{argmin} f(x)$, f has M -Lipschitz Hessian, $\nabla^2 f(x^*) \succeq \mu \text{Id}_n$ ($\mu > 0$), and $\|x_0 - x^*\| \leq \frac{\mu}{2M}$, then:*

$$\|x_{k+1} - x^*\| \leq \frac{M}{\mu} \|x_k - x^*\|^2 . \quad (10.4)$$

It is hard to imagine a better rate than this one: the sequence $u_k = \frac{M}{\mu} \|x_k - x^*\|$ satisfies $u_{k+1} \leq u_k^2 \leq u_0^{2^k}$. If $u_0 < 1$, the number of iterations to reach $u_k < \varepsilon$ is of the order $\log(\log \varepsilon)$, that is, in practice, basically a constant!

Proof. By the fundamental theorem of calculus, for all x, h in \mathbb{R}^d :

$$\nabla f(x+h) - \nabla f(x) = \int_0^1 \nabla^2 f(x+th)h \, dt \quad (10.5)$$

Hence

$$\nabla f(x_k) = \int_0^1 \nabla^2 f(x^* + t(x_k - x^*))(x_k - x^*) \, dt \quad (10.6)$$

and thus

$$x_{k+1} - x^* = x_k - x^* - [\nabla^2 f(x_k)]^{-1} \int_0^1 \nabla^2 f(x^* + t(x_k - x^*))(x_k - x^*) \, dt \quad (10.7)$$

$$= [\nabla^2 f(x_k)]^{-1} \int_0^1 \left(\nabla^2 f(x_k) - \nabla^2 f(x^* + t(x_k - x^*)) \right) (x_k - x^*) \, dt \quad (10.8)$$

Let's control the operator norm of the matrix being integrated, using M -Lipschitzness of the Hessian

$$\int_0^1 \|\nabla^2 f(x_k) - \nabla^2 f(x^* + t(x_k - x^*))\| \, dt \leq \int_0^1 (1-t)M\|x_k - x^*\| \, dt = \frac{M}{2}\|x_k - x^*\| \quad (10.9)$$

Finally, let's bound the operator norm of $[\nabla^2 f(x_k)]^{-1}$:

$$\nabla^2 f(x_k) \succeq \nabla^2 f(x^*) - M\|x_k - x^*\|\text{Id}_n \succeq (\mu - M \underbrace{\|x_k - x^*\|}_{\leq \frac{\mu}{2M} \text{ (induction)}})\text{Id}_n \succeq \frac{\mu}{2}\text{Id}_n \quad (10.10)$$

□

Despite the excellent convergence rate, Newton's method is not the “one algorithm to rule them all”:

- it may diverge if not started close enough to optimum
- the Hessian may not be invertible
- the Hessian is costly to compute
- the Hessian is even more costly to invert¹⁴ ($\mathcal{O}(d^3)$)

The principle of Quasi-Newton methods is to replace the inverse Hessian in [Equation \(10.2\)](#) by less costly approximations, in order to retain the very fast convergence rate while alleviating the computational burden.

¹⁴To compute $A^{-1}b$ numerically, one should solve the linear system $Ax = b$ with a LU factorization (Gaussian elimination) rather than invert A and apply to b : this is more stable and less costly. But the cost of the approach is still high

10.2 Quasi Newton methods

Algorithms in this section will have the following form:

$$\begin{cases} d_k = -B_k g_k \\ x_{k+1} = x_k + \rho_k d_k \end{cases} \quad \text{or} \quad \begin{cases} d_k = -H_k^{-1} g_k \\ x_{k+1} = x_k + \rho_k d_k \end{cases} \quad (10.11)$$

where $g_k = \nabla f(x_k)$, d_k is to be thought of as a descent direction, H_k (resp. B_k) is an approximation of the Hessian (resp. the inverse Hessian) of f at x_k , and ρ_k is a step size.

How to build such approximations? It is reasonable to impose that an approximation should satisfy the *secant condition* or *Quasi-Newton relation*:

$$g_{k+1} - g_k = H_{k+1}(x_{k+1} - x_k) \quad \text{or} \quad x_{k+1} - x_k = B_{k+1}(g_{k+1} - g_k) \quad (10.12)$$

Where does this condition come from? Suppose we have finished iteration $k+1$ and, to compute x_{k+2} , we want to construct (then minimize) ϕ_{k+1} , a quadratic approximation of f having H_{k+1} as Hessian:

$$\phi_{k+1}(x) = f(x_{k+1}) + \langle g_{k+1}, x - x_{k+1} \rangle + \frac{1}{2} \langle x - x_{k+1}, H_{k+1}(x - x_{k+1}) \rangle . \quad (10.13)$$

This quadratic approximation of f is exact at x_{k+1} ($\phi(x_{k+1}) = f(x_{k+1})$), and tangent at x_{k+1} too ($\nabla \phi_{k+1}(x_{k+1}) = \nabla f(x_{k+1})$). To impose the secant condition is simply to impose that f and ϕ_{k+1} are also tangent at the previous iterate, x_k .

The approximations in this section are updated iteratively, using rank 1 or rank 2 “corrections” from B_k (resp. H_k) to B_{k+1} (resp. H_{k+1}).

Following the notation of Nocedal and Wright, we write

$$y_k = g_{k+1} - g_k = \nabla f(x_{k+1}) - \nabla f(x_k) , \quad (10.14)$$

$$s_k = x_{k+1} - x_k . \quad (10.15)$$

10.2.1 The Broyden/SR1 formula

First, we consider rank-one updates of the Hessian, i.e. updates of the form $H_{k+1} = H_k + \sigma v v^\top$. How should we set σ and v so that it satisfies the secant condition (10.12)? We want:

$$y_k = H_{k+1} s_k \quad (10.16)$$

$$= H_k s_k + \sigma v v^\top s_k \quad (10.17)$$

$$= H_k s_k + (\sigma v^\top s_k) v \quad (\text{rotation trick } ab^\top c = b^\top c a) , \quad (10.18)$$

This does not give us v , but it gives its direction: v is proportional to $y_k - H_k s_k$. So we know¹⁵ that there exists a scalar δ such that $v = \delta(y_k - H_k s_k)$, and we can plug it back in Equation (10.17) to solve for δ and obtain that

$$H_{k+1} = H_k + \frac{1}{s_k^\top (y_k - H_k s_k)} (y_k - H_k s_k)(y_k - H_k s_k)^\top . \quad (10.19)$$

¹⁵what if $v^\top s_k = 0$?

satisfies the secant condition. This is called the Broyden or SR1 formula.

This approximation removes the cost of computing the Hessian, and also its cost of storing (we just store the scalars $s_k^\top(y_k - H_k s_k)$ and vectors $(y_k - H_k s_k)$). But at first sight, we still have the $\mathcal{O}(d^3)$ cost of computing the update (linear system solving). However, since the update of the Hessian is rank-one, by the Sherman-Morrison formula (itself a special case of the Woodbury formula, see [Exercise 10.2](#)), so is the update of the inverse Hessian $B_k = H_k^{-1}$:

$$B_{k+1} = B_k - \frac{1}{(B_k y_k - s_k)^\top y_k} (B_k y_k - s_k)(B_k y_k - s_k)^\top. \quad (10.20)$$

Therefore if we start at a simple B_0 , say Id , then we only need to do, at each iteration, rank 1 updates that are cheap to store and apply ([Exercise 10.1](#))!

10.2.2 BFGS

The BFGS update is rank two:

$$H_{k+1} = H_k + \frac{y_k y_k^\top}{y_k^\top s_k} - \frac{H_k s_k s_k^\top H_k}{s_k^\top H_k s_k}. \quad (10.21)$$

Proposition 10.2 (Inverse Hessian update for BFGS). *In terms of inverse Hessian B , the BFGS update (10.21) translates into*

$$B_{k+1} = \left(\text{Id} - \frac{s_k y_k^\top}{s_k^\top y_k} \right) B_k \left(\text{Id} - \frac{y_k s_k^\top}{s_k^\top y_k} \right) + \frac{s_k s_k^\top}{s_k^\top y_k}. \quad (10.22)$$

Proof. For lighter notation we write H , y and s for H_k , y_k , s_k . We apply the Sherman-Morrison formula twice in [Equation \(10.21\)](#). First

$$\left(H + \frac{y y^\top}{y^\top s} \right)^{-1} = H^{-1} - \frac{H^{-1} y y^\top H^{-1}}{s^\top y + y^\top H^{-1} y} \quad (10.23)$$

Then:

$$H_{k+1}^{-1} = H^{-1} - \underbrace{\frac{H^{-1} y y^\top H^{-1}}{s^\top y + y^\top H^{-1} y}}_{M_1} + \underbrace{\frac{\left(\text{Id} - \frac{H^{-1} y y^\top}{s^\top y + y^\top H^{-1} y} \right) s s^\top \left(\text{Id} - \frac{y y^\top H^{-1}}{s^\top y + y^\top H^{-1} y} \right)}{s^\top H s - s^\top H \left(H^{-1} - \frac{H^{-1} y y^\top H^{-1}}{s^\top y + y^\top H^{-1} y} \right) H s}}_{M_2} \quad (10.24)$$

The last term simplifies:

$$M_2 = \frac{s^\top y + y^\top H^{-1} y}{(s^\top y)^2} \left(\text{Id} - \frac{H^{-1} y y^\top}{s^\top y + y^\top H^{-1} y} \right) s s^\top \left(\text{Id} - \frac{y y^\top H^{-1}}{s^\top y + y^\top H^{-1} y} \right) \quad (10.25)$$

$$= \frac{1}{s^\top y} s s^\top + \frac{y^\top H^{-1} y}{(s^\top y)^2} s s^\top - \frac{1}{(s^\top y)^2} (H^{-1} y y^\top s s^\top + s s^\top y y^\top H^{-1}) + M_1 \quad (10.26)$$

$$= \frac{1}{s^\top y} s s^\top + \frac{y^\top H^{-1} y}{(s^\top y)^2} s s^\top - \frac{1}{s^\top y} (H^{-1} y s^\top + s y^\top H^{-1}) + M_1 \quad (10.27)$$

And so M_1 cancels out in Equation (10.24):

$$H_{k+1}^{-1} = H^{-1} + \frac{ss^\top}{s^\top y} + \frac{y^\top H^{-1} y s s^\top}{(s^\top y)^2} - \frac{sy^\top H^{-1} + H^{-1} y s^\top}{s^\top y} \quad (10.28)$$

$$= H^{-1} + \frac{ss^\top}{s^\top y} + \frac{sy^\top H^{-1} y s^\top}{(s^\top y)^2} - \frac{sy^\top H^{-1} + H^{-1} y s^\top}{s^\top y} \quad (10.29)$$

$$= \left(\text{Id} - \frac{sy^\top}{s^\top y} \right) H^{-1} \left(\text{Id} - \frac{ys^\top}{s^\top y} \right) + \frac{ss^\top}{s^\top y} \quad (10.30)$$

□

One very beautiful way to see this update is the following one.

Proposition 10.3 (BFGS as KL minimization). *The BFGS update Equation (10.21) is the solution of the following problem:*

$$\min_H -\log \det(H) + \langle H_k^{-1}, H \rangle \quad \text{subject to} \quad H = H^\top, Hs = y, \quad (10.31)$$

where the objective function, also known as the Stein loss, is equal to the KL divergence between two Gaussians having the same means, aka the Bregman divergence induced by the negative logdet.

Proof. Introduce the Lagrangian:

$$\mathcal{L}(H, \lambda, \Theta) = -\log \det H + \langle H_k^{-1}, H \rangle + \langle \lambda, Hs - y \rangle + \langle \Theta, H - H^\top \rangle. \quad (10.32)$$

Notice that $\langle \lambda, Hs \rangle$ is a scalar, equal to its trace $\text{tr } \lambda^\top Hs = \text{tr } s \lambda^\top H = \langle \lambda s^\top, H \rangle$ where this time the scalar product is over matrices. Similarly, $\langle \Theta, H - H^\top \rangle = \langle \Theta, H \rangle - \langle \Theta, H^\top \rangle = \langle \Theta, H \rangle - \langle \Theta^\top, H \rangle = \langle \Theta - \Theta^\top, H \rangle$.

This gives us the two gradients we need to write the first-order optimality condition over H . Using that the gradient of logdet is the inverse, we obtain:

$$\nabla_H \mathcal{L} = 0 \iff H^{-1} = H_k^{-1} + \lambda s^\top + \Theta - \Theta^\top. \quad (10.33)$$

First, we get rid of Θ : H must be symmetric hence H^{-1} too; H_k is symmetric (it's the approximate Hessian at the previous iteration) and this gives:

$$\Theta - \Theta^\top = \frac{1}{2}(s \lambda^\top - \lambda s^\top), \quad (10.34)$$

so

$$H^{-1} = H_k^{-1} + \frac{1}{2}(\lambda s^\top + s \lambda^\top). \quad (10.35)$$

To find λ , we use the other condition, $Hs = y$ aka $s = H^{-1}y$. This gives

$$\frac{1}{2}s^\top y \lambda = (s - H_k^{-1}y) - \frac{1}{2}\lambda^\top y s. \quad (10.36)$$

This may seem impossible to solve as λ appears twice, but it nevertheless tells us something: there exists a scalar a such that

$$\lambda = -\frac{2}{s^\top y} H_k^{-1} y + a s . \quad (10.37)$$

Plug back this expression to obtain H^{-1} as a function of a , then use $H^{-1}y = s$ to solve in a (it's tedious, but a good exercise). Then plugging back a into λ we obtain

$$\lambda s^\top = -\frac{2}{s^\top y} H_k^{-1} y s^\top + \frac{y^\top H_k^{-1} y + s^\top y}{(s^\top y)^2} s s^\top . \quad (10.38)$$

To conclude let's be a little lazy: $\frac{1}{2}(\lambda s^\top + s \lambda^\top)$ is the symmetric part of λs^\top . As visible in [Equation \(10.38\)](#), the term in $s s^\top$ is already symmetric, so we only need to symmetrize the other one:

$$\frac{1}{2}(\lambda s^\top + s \lambda^\top) = \frac{y^\top H_k^{-1} y + s^\top y}{(s^\top y)^2} s s^\top - \frac{1}{s^\top y} (H_k^{-1} y s^\top + s y^\top H_k^{-1}) . \quad (10.39)$$

and we recognize the computation we've been through in [Equation \(10.28\)](#), which allows to conclude invoking [Proposition 10.2](#). \square

Proposition 10.4 (BFGS as weighted Frobenius norm minimization). *Let W be any positive definite matrix satisfying $y_k = W s_k$. Then the inverse Hessian B_{k+1} solves:*

$$\min_B \|W^{1/2}(B - B_k)W^{1/2}\| \quad \text{subject to} \quad B = B^\top, By = s . \quad (10.40)$$

Proof. The proof is similar to that of [Proposition 10.3](#). Introduce the Lagrangian:

$$\mathcal{L}(B, \lambda, \Theta) = \frac{1}{2} \|W^{1/2}(B - B_k)W^{1/2}\| + \langle \lambda, By - s \rangle + \langle \Theta, B - B^\top \rangle . \quad (10.41)$$

The gradient with respect to B of the first term is $W(B - B_k)W$. For the second one, $\langle \lambda, By \rangle$ is a scalar, equal to its trace $\text{tr } \lambda^\top By = \text{tr } y \lambda^\top B = \langle \lambda y^\top, B \rangle$ where this time the scalar product is over matrices. For the last one, $\langle \Theta, B - B^\top \rangle = \langle \Theta, B \rangle - \langle \Theta, B^\top \rangle = \langle \Theta, B \rangle - \langle \Theta^\top, B \rangle = \langle \Theta - \Theta^\top, B \rangle$.

Therefore:

$$\nabla_H \mathcal{L} = 0 \iff WBW = WB_kW + \lambda y^\top + \Theta - \Theta^\top . \quad (10.42)$$

Next, we get rid of Θ : B must be symmetric hence B^{-1} too; B_k is symmetric by assumption and so:

$$\Theta - \Theta^\top = \frac{1}{2}(y \lambda^\top - \lambda y^\top) , \quad (10.43)$$

so

$$WBW = WB_kW + \frac{1}{2}(\lambda y^\top + y \lambda^\top) \quad (10.44)$$

$$B = B_k + \frac{1}{2}W^{-1}(\lambda y^\top + y \lambda^\top)W^{-1} . \quad (10.45)$$

To find λ , we use the other condition, $By = BWs = s$. Applying (10.44) to s thus yields

$$WBWs = WBy = Ws = y = WB_ky + \frac{1}{2}(\lambda y^\top s + y\lambda^\top s) \quad (10.46)$$

Hence,

$$\frac{y^\top s}{2}\lambda = y - WB_ky - \lambda^\top sy, \quad (10.47)$$

so there exists a such that

$$\lambda = -\frac{2}{y^\top s}WB_ky + ay. \quad (10.48)$$

Since $By = s$,

$$s = B_ky + \frac{1}{2}W^{-1}(\lambda y^\top + y\lambda^\top)W^{-1}y \quad (10.49)$$

$$= B_ky + \frac{1}{2}W^{-1}\left(\left(-\frac{2}{y^\top s}WB_ky + ay\right)y^\top + y\left(-\frac{2}{y^\top s}y^\top B_kW + ay^\top\right)\right)s \quad (10.50)$$

$$= \frac{1}{2}W^{-1}\left(ayy^\top + y\left(-\frac{2}{y^\top s}y^\top B_kW + ay^\top\right)\right)s \quad (10.51)$$

$$= \frac{1}{2}asy^\top s - \frac{1}{y^\top s}sy^\top B_kWs + \frac{1}{2}asy^\top s \quad (10.52)$$

$$= ay^\top ss - \frac{1}{y^\top s}y^\top B_kWss \quad (10.53)$$

$$= ay^\top ss - \frac{1}{y^\top s}y^\top B_kys \quad (10.54)$$

Hence

$$a = \frac{y^\top s + y^\top B_ky}{(y^\top s)^2}. \quad (10.55)$$

So

$$\lambda y^\top = -\frac{2}{y^\top s}WB_kyy^\top + \frac{y^\top s + y^\top B_ky}{(y^\top s)^2}yy^\top \quad (10.56)$$

$$\frac{1}{2}(\lambda y^\top + y\lambda^\top) = -\frac{1}{y^\top s}(WB_kyy^\top + yy^\top B_kW) + \frac{y^\top s + y^\top B_ky}{(y^\top s)^2}yy^\top. \quad (10.57)$$

Substituting in (10.45) and using $W^{-1}y = s$,

$$B = B_k + \frac{1}{y^\top s}W^{-1}(\lambda y^\top + y\lambda^\top)W^{-1} \quad (10.58)$$

$$= B_k - \frac{1}{y^\top s}(B_kys^\top + sy^\top B_k) + \frac{y^\top s + y^\top B_ky}{(y^\top s)^2}ss^\top \quad (10.59)$$

$$= \left(\text{Id} - \frac{sy^\top}{y^\top s}\right)B_k\left(\text{Id} - \frac{sy^\top}{s^\top y}\right) + \frac{ss^\top}{y^\top s}. \quad (10.60)$$

□

10.2.3 DFP

The DFP update is the same as BFGS, swapping the roles of H, y, s with B, s, y :

$$B_{k+1} = B_k + \frac{s_k s_k^\top}{y_k^\top s_k} - \frac{B_k y_k y_k^\top B_k}{y_k^\top B_k y_k} . \quad (10.61)$$

10.3 Exercises

Exercise 10.1. ☹ Show that any $n \times n$ matrix of rank 1 writes xy^\top with $x, y \in \mathbb{R}^n$. Show that the cost of multiplying by a rank one matrix is $2n$ instead of the usual n^2 .

Exercise 10.2 (Woodbury inverse formula). ☹☹ Let $n, k \in \mathbb{N}$. Let $U \in \mathbb{R}^{n \times k}$ and $V \in \mathbb{R}^{k \times n}$. Show that

$$(\text{Id} + UV)^{-1} = \text{Id} - U(\text{Id} + VU)^{-1}V . \quad (10.62)$$

Show that for any invertible A ,

$$(A + UV)^{-1} = A^{-1} - A^{-1}U(\text{Id} + VA^{-1}U)^{-1}VA^{-1} . \quad (10.63)$$

When $k = 1$ and U and V respectively become column and row vectors u and v^\top , the formula is known as the Sherman-Morrison formula:

$$(A + uv^\top)^{-1} = A^{-1} - \frac{1}{1 + v^\top A^{-1}u} A^{-1}uv^\top A^{-1} . \quad (10.64)$$

What is the naive cost of computing $A^{-1}uv^\top A^{-1}$? And the clever cost?

11 The finite sum structure: stochastic methods

Resources: see the very complete work of [Garrigos and Gower \(2023\)](#).

11.1 Stochastic gradient descent

We consider problems with the *finite sum structure*:

$$\min_x f(x) = \sum_{i=1}^n f_i(x) . \quad (11.1)$$

Such problems notably arise from the maximum likelihood framework of [Section 1.1](#) where each f_i is the negative log-likelihood of a data point. Computing the gradient of f normally requires computing the gradient of each f_i , that scales with n . The principle of *stochastic gradient descent*¹⁶ is to replace $\nabla f(x_t)$ at iteration $t \in \mathbb{N}$ of gradient descent by $\nabla f_{i_t}(x_t)$ for some i_t sampled uniformly in $[n]$:

$$x_{t+1} = x_t - \eta \nabla f_{i_t}(x_t) . \quad (11.2)$$

The update direction is on average correct since:

$$\mathbb{E}_{i \sim \mathcal{U}([n])}[\nabla f_i(x)] = \sum_{i'=1}^n \mathbb{P}(i = i') \nabla f_{i'}(x) = \frac{1}{n} \sum_{i'=1}^n \nabla f_{i'}(x) = \nabla f(x) . \quad (11.3)$$

In addition, and this is the crux of the method, each update of SGD has a cost independent of that of n . For the price of one gradient update (n gradients, which often amounts to $\mathcal{O}(nd)$), one can perform n updates of SGD! This is called a full pass over the dataset, or one epoch.

Like (sub)gradient descent, stochastic gradient descent has rate $1/\sqrt{k}$ for convex, Lipschitz functions. Unfortunately, unlike GD, this does not improve when using other assumptions, e.g. smoothness of the objective. For example, there is a guaranteed lower bound of order $1/k$ for smooth strongly convex functions ([Nemirovski et al., 2009](#)).

See for example the following result, which shows a fast initial convergence, but then a plateau (the suboptimality does not go to 0).

Proposition 11.1. *Let f be a μ -strongly convex function with (unique) minimizer x^* . Let $0 < \eta \leq \frac{1}{\mu}$ and assume that the iterates of stochastic descent, x_t , have bounded gradients in expectation:*

$$\forall j \in [p], \forall t \in \mathbb{N}, \mathbb{E}[\|\nabla f_j(x_t)\|^2] \leq B^2 . \quad (11.4)$$

Then

$$\mathbb{E}[\|x_t - x^*\|^2] \leq (1 - \eta\mu)^t \|x_0 - x^*\|^2 + \frac{\eta}{\mu} B^2 . \quad (11.5)$$

As usual, a few comments before moving into the proof:

¹⁶this is not a descent algorithm, as the direction followed is a descent direction for f_i , but not for f

- There's randomness so our rate is in expectation. We cannot hope for “absolute” rates without assumption since we could get very unlucky and sample always the same data point. But that is unlikely to happen: the same kind of rates with high probability can also be obtained.
- The expectation in (11.4) is over the indices (i_0, \dots, i_t) that are sampled.
- We have two bounds in the rate: an exponentially fast “forgetting” of the initial conditions, and an “irreducible” term: the upper bound does not go to 0 for fixed η . To make the second term small, we need η small, but this makes the first phase slower... See an illustration of this in Figure 11.1.
- The boundedness (in expectation) of the gradients is a strong requirement. It holds if the f_i 's are Lipschitz, which is for example true for the logistic loss, but not for the least squares loss. If the f_i 's are C^2 and the iterates converge (thus are bounded), the result hold. But this is somewhat of a circular reasoning, as we cannot know if the iterates will remain bounded or not a priori.

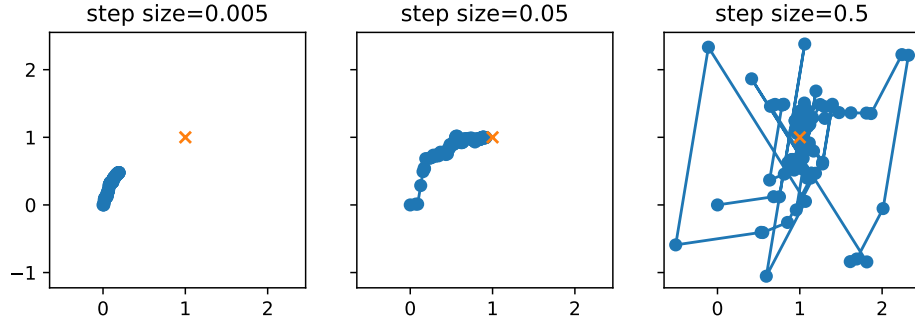


Figure 11.1: SGD iterates for three choices of step sizes. The minimizer is denoted by an orange cross. Using a larger stepsize accelerate convergence in the first steps, but slows it down closer to the minimizer.

Proof.

$$\|x_{t+1} - x^*\|^2 = \|x_t - \eta \nabla f_{i_t}(x_t) - x^*\|^2 \quad (11.6)$$

$$= \|x_t - x^*\|^2 - 2\eta \langle \nabla f_{i_t}(x_t), x_t - x^* \rangle + \eta^2 \|\nabla f_{i_t}(x_t)\|^2. \quad (11.7)$$

We can take conditional expectation with respect to x_1, \dots, x_t on both sides (denoted \mathbb{E}_t), and use that

- $\mathbb{E}_t[\|x_t - x^*\|^2] = \|x_t - x^*\|^2$ since x_t is x_t measurable.
- i_t is independent of x_1, \dots, x_t and $x_t - x^*$ is x_t measurable, so $\mathbb{E}_k[\langle \nabla f_{i_t}(x_t), x_t - x^* \rangle] = \mathbb{E}[\langle \nabla f_{i_t}(x_t), x_t - x^* \rangle] = \langle \nabla f(x_t), x_t - x^* \rangle$ (Williams, 1991, 9.10 with $r = 2$)

This gives us:

$$\mathbb{E}_k \|x_{t+1} - x^*\|^2 = \|x_t - x^*\|^2 - 2\eta \langle \nabla f(x_t), x_t - x^* \rangle + \eta^2 \mathbb{E}_t \|\nabla f_{i_t}(x_t)\|^2 \quad (11.8)$$

Then since $\nabla f(x^*) = 0$, by μ -strong convexity of f (Exercise 3.3),

$$\mu\|x_t - x^*\|^2 \leq \langle \nabla f(x_t), x_t - x^* \rangle. \quad (11.9)$$

Thus, by using Equation (11.9) and taking expectations, using $\mathbb{E}[\mathbb{E}_k[X]] = \mathbb{E}[X]$,

$$\mathbb{E}\|x_{t+1} - x^*\|^2 \leq (1 - 2\eta\mu)\mathbb{E}\|x_t - x^*\|^2 + \eta^2 B^2 \quad (11.10)$$

$$\leq (1 - 2\eta\mu)^{t+1}\mathbb{E}\|x_0 - x^*\|^2 + \eta^2 B^2 \sum_{t'=0}^t (1 - 2\eta\mu)^{t'} \quad (11.11)$$

$$\leq (1 - 2\eta\mu)^{t+1}\mathbb{E}\|x_0 - x^*\|^2 + \frac{1}{2\eta\mu}\eta^2 B^2 \quad (11.12)$$

□

It turns out that the bound on the norms of the gradient can be removed, in favor of another quantity, both more realistic and easier to check.

Definition 11.2. *The gradient noise of f is defined as:*

$$\sigma_f^* = \inf_{x^* \in \operatorname{argmin} f} \mathbb{V}[\nabla f_i(x^*)] \quad (11.13)$$

where $\mathbb{V}[X] = \mathbb{E}[\|X - \mathbb{E}[X]\|^2]$.

If the minimizer is unique, then the inf can be removed; note only depends on the function f , and not on the algorithm used to minimize it. Finally, observe that in the *interpolating* regime where x^* not only minimizes f , but also each f_i (e.g., in the case of least squares when a solution to $Ax = b$ exist), one has $\sigma_f^* = 0$ since all $\nabla f_i(x^*)$ are 0.

This quantity is super useful as it allows to bound the variance of the gradients at any point when the functions f_i are L -smooth and convex.

Proposition 11.3. *Assume that the f_i 's are convex and L -smooth. Then, for all $x \in \mathbb{R}^d$ and all $x^* \in \operatorname{argmin} f$,*

$$\mathbb{E}[\|\nabla f_i(x)\|^2] \leq 4L(f(x) - f(x^*)) + 2\sigma_f^*, \quad (11.14)$$

where the expectation is over the index i .

Proof. Since $(a + b)^2 \leq 2a^2 + 2b^2$,

$$\|\nabla f_i(x)\|^2 \leq 2\|\nabla f_i(x) - \nabla f_i(x^*)\|^2 + 2\|\nabla f_i(x^*)\|^2 \quad (11.15)$$

$$\mathbb{E}[\|\nabla f_i(x)\|^2] \leq 2\mathbb{E}[\|\nabla f_i(x) - \nabla f_i(x^*)\|^2] + 2\mathbb{E}[\|\nabla f_i(x^*)\|^2] \quad (11.16)$$

$$\leq 2\mathbb{E}[\|\nabla f_i(x) - \nabla f_i(x^*)\|^2] + 2\sigma_f^* \quad (11.17)$$

Then, use the following fact: for any $x, y \in \mathbb{R}^d$,

$$\frac{1}{2L}\|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq f_i(y) - f_i(x) + \langle \nabla f_i(x), y - x \rangle, \quad (11.18)$$

which, multiplied by $1/n$ and summed, yields:

$$\frac{1}{2L}\mathbb{E}[\|\nabla f_i(x) - \nabla f_i(y)\|^2] \leq f(y) - f(x) + \langle \nabla f(x), y - x \rangle. \quad (11.19)$$

Apply this in $x = x^*$, together with $\nabla f(x^*) = 0$, to conclude. □

Equipped with [Proposition 11.3](#), we can prove the following result, analogous to [Proposition 11.1](#).

Proposition 11.4. *Assume that all the f_i 's are convex and L -smooth, assume that f is μ -strongly convex. Then SGD with stepsize $0 < \eta < 1/(2L)$ satisfies:*

$$\mathbb{E}[\|x^{t+1} - x^*\|^2] \leq (1 - \eta\mu)^k \|x_0 - x^*\|^2 + \frac{2\eta\sigma_f^*}{\mu}. \quad (11.20)$$

In particular, in the interpolating regime ($\sigma_f^ = 0$), we get a linear rate as in the case of GD!*

Proof. As in the proof of [Proposition 11.1](#) we start from the same equation but this time use convexity of f to bound the second term is the RHS instead:

$$\mathbb{E}_k \|x_{t+1} - x^*\|^2 = \|x_t - x^*\|^2 - 2\eta \langle \nabla f(x_t), x_t - x^* \rangle + \eta^2 \mathbb{E}_k \|\nabla f_{i_t}(x_t)\|^2 \quad (11.21)$$

$$\leq \|x_t - x^*\|^2 - 2\eta (f(x_t) - f(x^*)) + \eta^2 \mathbb{E}_k \|\nabla f_{i_t}(x_t)\|^2 \quad (11.22)$$

$$\leq (1 - \eta\mu) \|x_t - x^*\|^2 + 2\eta(2\eta L - 1)(f(x_t) - f(x^*)) + 2\eta^2 \sigma_f^* \quad (11.23)$$

$$\leq (1 - \eta\mu) \|x_t - x^*\|^2 + 2\eta^2 \sigma_f^* \quad (11.24)$$

where we have used [Proposition 11.3](#) and the fact that $\eta \leq 1/2L$. The rest of the proof follows by recursively applying the result above and bounding the telescopic series by its sum. \square

11.2 Variance reduction

The main problem of SGD is the variance in the gradients. For example, if the gradient noise ([Definition 11.2](#)) is nonzero, if the algorithm is started at optimum it will move away from it. One needs to reduce the stepsize to stop moving too much, but then there's a hard tradeoff because you still want to go fast to the minimizer. Two remedies are:

- to use decreasing stepsizes, of the form C/t^α . This works, but requires tuning both C and α by trial and error. In ML and especially in DL, one can use adaptive *learning rates schedulers* that monitor some test loss decrease, and diminish the stepsize every time the test loss reaches a plateau.
- use *minibatches* to compute the gradient estimate, i.e. sample $B_t \subset [n]$ and use $\frac{1}{|B_t|} \sum_{i \in B_t} \nabla f_i(x^t)$ instead of $\nabla f_{i_t}(x^t)$. This reduces the variance, but requires computing more gradients (parallelism and memory size matter a lot in practice).

In the sequel we use an alternative to reduce the variance of the gradient estimate.

GD: Fixed learning rate, $O(1/k)$ or linear, no problem structure exploitation. SGD: decreasing stepsize, $O(1/\sqrt{k})$ or $O(1/k)$. Can we get the best of both worlds?

The idea of variance reduction is the following: random variable X , estimate $\mathbb{E}[X]$ but with lower variance. Use a random variable Y correlated with X with mean easier to estimate, introduce $\theta_\alpha = \alpha(X - Y) + \mathbb{E}[Y]$ for $\alpha \in [0, 1]$. Unbiased if $\alpha = 1$ or $\mathbb{E}X = \mathbb{E}Y$.

	GD	SGD	Variance reduced
∇f_i required per iteration	n	1	1
convex rate	$1/\sqrt{k}$	$1/\sqrt{k}$	$1/\sqrt{k}$
smooth convex rate	$1/k$	$1/\sqrt{k}$	$1/k$
smooth strongly convex rate	linear	$1/k$	linear

Variance:

$$\text{Var}[\theta_\alpha] = \mathbb{E}[(\alpha(X - Y) - \alpha\mathbb{E}[(X - Y)])^2] \quad (11.25)$$

$$= \alpha^2 \mathbb{E}[(X - \mathbb{E}X - (Y - \mathbb{E}Y))^2] \quad (11.26)$$

$$= \alpha^2 (\text{Var } X + \text{Var } Y - 2 \text{cov}(X, Y)) \quad (11.27)$$

can be lower than $\text{Var } X$, either using a low α , either a Y strongly correlated with X (so that $\text{Var } Y - 2 \text{cov}(X, Y) \leq 0$).

The two main algorithms are SAGA (itself an improvement over SAG [Schmidt et al. \(2013\)](#); [Defazio et al. \(2014\)](#)) and SVRG [Johnson and Zhang \(2013\)](#).

The SAGA algorithm at iteration t does:

$$\begin{cases} j \sim \mathcal{U}(n) \\ g_j^t = \nabla f_j(x) \\ g_i^t = g_i^{t-1} \quad \text{for } i \neq j \\ x^t = x^{t-1} - \gamma[g_j^t - g_j^{t-1} + \frac{1}{n} \sum_{i=1}^n g_i^{t-1}] \end{cases} \quad (11.28)$$

A few remarks:

- the recomputation of $\frac{1}{n} \sum_{i=1}^n g_i^{k-1}$ is a priori expensive ($\mathcal{O}(nd)$) but since only one element in the sum changes at every iteration,
- the memory cost is high: $\mathcal{O}(nd)$. But for generalized models ($f_i(x) = \phi_i(a_i^\top x)$), it can be reduced to $\mathcal{O}(n)$ by only storing the scalars $\phi'(a_i^\top x^t)$
- how to initialize the table of gradients g ? Authors suggest to first perform a full pass of SGD, updating all f_i 's sequentially
- SAGA can also handle an additional non smooth term in the objective through its prox-SAGA variant
- when the dataset is composed of sparse vectors (a frequent text with encoded text data, e.g. tf-idf transforms), individual gradients are sparse (in the GLM setting) but the mean is dense, which is costly. There exists a clever remedy (see eq 8 in “Improved Asynchronous Parallel Optimization Analysis for Stochastic Incremental Methods”, Leblond et al).
- the last A in SAGA stands for SAG Amélioré (sic); the SAG (Stochastic Average Gradient) used: $x^t = x^{t-1} - \gamma[\frac{g_j^t - g_j^{t-1}}{n} + \frac{1}{n} \sum_{i=1}^n g_i^{t-1}]$, which are biased updates.

The Stochastic Variance Reduced Gradient (SVRG, [Johnson and Zhang 2013](#)) uses a snapshot point \tilde{x} which is periodically updated:

$$x^t = x^{t-1} - \gamma[\nabla f_j(x^k) - \nabla f_j^{\tilde{x}} + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})] \quad (11.29)$$

$$\tilde{x} = \begin{cases} x^t & \text{if } t = 0 \pmod{m} \\ \tilde{x} & \text{otherwise.} \end{cases} \quad (11.30)$$

When one updates \tilde{x} , one computes the full gradient in order to have only two stochastic gradients to compute at each update of x^t . The parameter to tune is the update frequency of \tilde{x} (which requires a full gradient computation).

12 Mirror descent

The mirror descent algorithm is a generalization of gradient descent, with similar properties, that allows handling “different geometries”. See [Section 13.2](#) for an analysis of its continuous-time version.

Key to the Mirror descent algorithm is the concept of Bregman divergence.

Definition 12.1 (Legendre function). *A function $J : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ is Legendre if and only if:*

- *it is in $\Gamma_0(\mathbb{R}^d)$*
- *it is strictly convex on (every convex subset of) the interior of its domain*
- *it is essentially smooth, meaning that it is continuously differentiable on the interior of its domain, and $\|\nabla J(x_k)\| \rightarrow \infty$ for any sequence x_k converging to a point on the boundary of $\text{dom } J$.*

Legendre functions are nice because they enjoy the following properties regarding their gradient and subdifferential.

Proposition 12.2. *Let J be Legendre. Then*

1. $\text{dom } \partial J = \int \text{dom } J$
2. $\nabla J : \int \text{dom } J \rightarrow \int \text{dom } J^*$ is continuous and bijective, with inverse $\nabla J^* : \int \text{dom } J^* \rightarrow \int \text{dom } J$, which is also continuous.

Those are Theorems 26.1 and 26.4 in [Rockafellar \(1970\)](#).

The half-squared ℓ_2 norm is Legendre, but both $\|\cdot\|_1$ and $\|\cdot\|_1 + \frac{1}{2}\|\cdot\|_2^2$ are not: the first is not strictly convex, the second is not continuously differentiable at 0.

Definition 12.3 (Bregman divergence). *Let J be a convex function. The Bregman divergence induced by J is*

$$D_J : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \quad (12.1)$$

$$(x, y) \mapsto J(x) - J(y) - \langle \nabla J(y), x - y \rangle.$$

For brevity, we may sometimes omit the J subscript and write D when clear from context. J is often referred to as the potential or the entropy of D_J .

Let us cite a few properties of D_J :

- By convexity of J , it is positive.
- It is usually not symmetric (see [Exercise 12.1](#)).
- It is point-separating if J is strictly convex.
- It is the difference at x between J and its linear lower bound at y .
- It is convex in x but usually not in y .
- L -smoothness of J amounts to upper bounding $D_J(x, y)$ by $\frac{L}{2}\|x - y\|^2$. Similarly, μ -strong convexity J amounts to lower bounding $D_J(x, y)$ by $\frac{\mu}{2}\|x - y\|^2$.

Also note that a Bregman divergence is not like a distance: it is more like a *square* distance.

Proposition 12.4 (Three-point equality). *Let $J : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$ be differentiable. For any $x, y, z \in \text{dom } f$,*

$$D_J(x, y) + D_J(z, x) - D_J(z, y) = \langle \nabla J(x) - \nabla J(y), x - z \rangle. \quad (12.2)$$

Proof. For any $x, y, z \in \text{dom } f$,

$$\begin{aligned} D_J(x, y) + D_J(z, x) - D_J(z, y) &= J(x) - J(y) - \langle \nabla J(y), x - y \rangle \\ &\quad + J(z) - J(x) - \langle \nabla J(x), z - x \rangle \\ &\quad - J(z) + J(y) + \langle \nabla J(y), z - y \rangle \end{aligned} \quad (12.3)$$

$$= \langle \nabla J(y), z - y - x + y \rangle - \langle \nabla J(x), z - x \rangle \quad (12.4)$$

$$= \langle \nabla J(x) - \nabla J(y), x - z \rangle \quad (12.5)$$

□

Example 12.5. *The Bregman divergence induced by $\frac{1}{2}\|\cdot\|^2$ is also $\frac{1}{2}\|\cdot\|^2$.*

Example 12.6 (Kullback-Leibler divergence). *The following potential, called negative entropy, is particularly adapted to optimization on the simplex:*

$$\begin{aligned} H : (\mathbb{R}_+^*)^d &\rightarrow \mathbb{R} \\ x &\mapsto \sum_{i=1}^d x_i \log x_i \end{aligned} \quad (12.6)$$

It is convex, it is even 1/2-strongly convex with respect to the ℓ_1 norm, meaning that for all x, y , for any $\lambda \in]0, 1[$

$$H(\lambda x + (1 - \lambda)y) \leq \lambda H(x) + (1 - \lambda)H(y) - \frac{1}{2}\lambda(1 - \lambda)\|x - y\|_1^2. \quad (12.7)$$

This result is known as Pinsker's inequality. The Bregman divergence associated to H is:

$$D(x, y) = \sum_{i=1}^d x_i \log x_i - y_i \log y_i - (1 + \log y_i)(x_i - y_i) \quad (12.8)$$

$$= \sum_{i=1}^d x_i \log \frac{x_i}{y_i} - \sum_{i=1}^d x_i + \sum_{i=1}^d y_i \quad (12.9)$$

which is called the Kullback-Leibler divergence. It is very often considered for vectors x and y on the simplex, in which case the last two sums cancel. This Bregman divergence induces a nicely behaved projection onto the simplex. For $y \in \mathbb{R}_{++}^d$ (the case where y has vanishing entries is easy to handle similarly), let us compute

$$\operatorname{argmin}_{x \in \Delta_d} \sum_{i=1}^d x_i \log \frac{x_i}{y_i} - \sum_{i=1}^d x_i + \sum_{i=1}^d y_i \quad (12.10)$$

Since the problem is convex, with one affine constraint $\sum_i x_i = 1$ and inequality conditions $-x_i \leq 0$ for which Slater's condition holds, this problem can be solved with the Lagrangian. For $\lambda \in \mathbb{R}$ and $\mu \in \mathbb{R}_+^d$, the Lagrangian is:

$$\mathcal{L}(x, \lambda, \mu) = \sum_{i=1}^d x_i \log \frac{x_i}{y_i} - \sum_{i=1}^d x_i + \lambda \left(\sum_{i=1}^d x_i - 1 \right) - \sum_{i=1}^d \mu_i x_i \quad (12.11)$$

The KKT conditions on are:

$$\begin{cases} \mu_i x_i = 0 \\ \log \frac{x_i}{y_i} + 1 - 1 + \lambda - \mu_i = 0, \text{ aka } x_i = y_i \exp(\lambda - \mu_i) \\ \sum_{i=1}^d x_i = 1 \end{cases} \quad (12.12)$$

The second equation shows that x_i cannot have value 0, hence $1 = \sum_{i=1}^d x_i = \exp(\lambda) \sum_{i=1}^d y_i$, hence the solution is given by:

$$\operatorname{argmin}_{x \in \Delta_d} D(x, y) = \left(\frac{y_i}{\sum_{i=1}^d y_i} \right)_{i \in [d]} \quad (12.13)$$

which is just a normalization. Hence, in an algorithm, we can interpret a normalization as a projection “in the Bregman sense”.

12.1 The mirror descent algorithm

Remember that for gradient descent, we have that

$$x_{k+1} = x_k - \eta \nabla f(x_k) = \operatorname{argmin}_x f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2\eta} \|x - x_k\|^2; \quad (12.14)$$

one possible interpretation is that GD tries to minimize the linear approximation of f at x_k , while not moving too far away from x_k . The penalty for moving away from x_k is $\frac{1}{2\eta} \|x - x_k\|^2$.

But we may be in a case where the penalty is not isotropic: moving away in some directions should be penalized more. This is the generalization of Mirror Descent: replace $\frac{1}{2\eta} \|x - x_k\|^2$ by $\frac{1}{\eta} D_J(x, x_k)$, using some potential J .

Definition 12.7. The mirror descent algorithm with stepsize $\eta > 0$ is given by the following implicit equation:

$$\nabla J(x_{k+1}) = \nabla J(x_k) - \eta \nabla f(x_k). \quad (12.15)$$

It satisfies

$$x_{k+1} = \operatorname{argmin}_x f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{\eta} D_J(x, x_k). \quad (12.16)$$

If $J \in \Gamma_0(\mathbb{R}^d)$, we can invert its gradient by [Proposition 8.4](#): $(\nabla J)^{-1} = \nabla J^*$. The mirror descent algorithm then reads:

$$x_{k+1} = \nabla J^*(\nabla J(x_k) - \eta \nabla f(x_k)). \quad (12.17)$$

The convergence rates of mirror descent are similar to those of gradient descent.

Proposition 12.8. *Let f be a convex function, let $x^* \in \operatorname{argmin} f$. Let f be L -Lispchitz in some norm $\|\cdot\|$ (not necessarily the ℓ_2 one): all its gradients are bounded by L in the dual norm $\|\cdot\|^*$ (Exercise 8.3). Let J be ρ -strongly convex, still with respect to $\|\cdot\|$. Let $R = D(x^*, x_0)$. Then for $\eta = \sqrt{\frac{2\rho R}{L^2 t}}$ the iterates of mirror descent satisfy:*

$$\frac{1}{t} \sum_{k=0}^{t-1} f(x_k) - f(x^*) \leq \sqrt{\frac{2RL^2}{\rho t}}. \quad (12.18)$$

Note that choosing the right η requires knowing $D(x^*, x_0)$; in the case of projected mirror descent, $D(x^*, x_0)$ can be replaced by the diameter of the feasible set (measured in D).

Proof. The proof is similar to that of subgradient descent (Proposition 3.2). By convexity of f ,

$$f(x_k) - f(x^*) \leq \langle \nabla f(x_k), x_k - x^* \rangle \quad (12.19)$$

$$= \frac{1}{\eta} \langle J(x_k) - J(x_{k+1}), x_k - x^* \rangle \quad (12.20)$$

$$= \frac{1}{\eta} [D(x_k, x_{k+1}) + D(x^*, x_k) - D(x^*, x_{k+1})] \quad (12.21)$$

Sum to telescope:

$$\sum_{k=0}^{t-1} f(x_k) - f(x^*) \leq \frac{1}{\eta} \left[\sum_{k=0}^{t-1} D(x_k, x_{k+1}) \right] + D(x^*, x_0) - D(x^*, x_t) \quad (12.22)$$

$$\leq \frac{1}{\eta} \left[\sum_{k=0}^{t-1} D(x_k, x_{k+1}) \right] + D(x^*, x_0) \quad (12.23)$$

It remains to bound the first term, for which we use strong convexity. Strong convexity means that $D(x_k, x_{k+1})$ is lower bounded, which is not very useful here. So we make the opposite of a Bregman divergence appear:

$$D(x_k, x_{k+1}) = J(x_k) - J(x_{k+1}) - \langle \nabla J(x_{k+1}), x_k - x_{k+1} \rangle \quad (12.24)$$

$$= \langle \nabla J(x_{k+1}) - \nabla J(x_k), x_{k+1} - x_k \rangle - D(x_{k+1}, x_k) \quad (12.25)$$

$$= \langle -\eta \nabla f(x_k), x_{k+1} - x_k \rangle - D(x_{k+1}, x_k) \quad (12.26)$$

$$\leq \eta L \|x_{k+1} - x_k\| - \frac{\rho}{2} \|x_{k+1} - x_k\|^2 \quad (12.27)$$

$$\leq \frac{\eta^2 L^2}{2\rho} \quad (12.28)$$

where we have used Cauchy-Schwarz in the $\|\cdot\|$ and maximized the concave quadratic $\eta La - \frac{\rho}{2} a^2$.

Plugging back in Equation (12.23) and divide by t :

$$\frac{1}{t} \sum_{k=0}^{t-1} f(x_k) - f(x^*) \leq \eta \frac{L^2}{2\rho} + \frac{D(x^*, x_0)}{\eta t}$$

Minimize RHS in η gives $\eta = \sqrt{\frac{2\rho R}{L^2 t}}$ and an upper bound value of $\sqrt{\frac{2RL^2}{\rho t}}$. \square

This $1/\sqrt{k}$ rates improves to $1/k$ when f is L -smooth. In the following section, we prove the result for mirror descent, but also for many other algorithms.

12.2 One proof to rule them all

This proof handles the proximal point algorithm, gradient descent, proximal gradient descent, mirror descent, projected mirror descent... It can be found in Tseng (2010).

Lemma 12.9 (3 points identity). *Let $z_t = \operatorname{argmin}_x \Psi(x) + D(x, z)$. Then for all x ,*

$$\Psi(x) + D(x, z) \geq \Psi(z_t) + D(z_t, z) + D(x, z_t). \quad (12.29)$$

Proof. Use the first order optimality condition $0 = \nabla \Psi(z_t) + \nabla J(z_t) - \nabla J(x)$, then write convexity of Ψ with this gradient and use Proposition 12.4. \square

Proposition 12.10. *Let $F = f + \tau g$ with f convex L -smooth and g convex. Let $\ell_F(x, y)$ be the linearization of the objective at y :*

$$\ell_F(x, y) = f(y) + \langle \nabla f(y), x - y \rangle + \tau g(x). \quad (12.30)$$

Let the sequence x_k be defined by

$$x_{k+1} = \operatorname{argmin}_x \ell_F(x, y) + LD(x, x_k). \quad (12.31)$$

Then for all x ,

$$F(x_k) - F(x) \leq \frac{D(x, x_0)}{k}. \quad (12.32)$$

If a minimizer of F exists, we obtain a convergence rate.

We recover

- the proximal point algorithm by taking $f = 0$ and $J = \frac{1}{2}\|\cdot\|^2$. In that case we can take any L , so: there is no constraint on the stepsize in the proximal point algorithm!
- gradient descent for $g = 0$ and $J = \frac{1}{2}\|\cdot\|^2$. Note that this specific proof constrains the stepsize to be smaller than $1/L$, but we have convergence results for up to $\gamma < 2/L$.
- proximal gradient descent for $J = \frac{1}{2}\|\cdot\|^2$
- mirror descent and projected mirror descent for generic J

Proof. By convexity and smoothness of f ,

$$F(x) \geq \ell_F(x, y) \geq F(x) - \frac{L}{2} \|x - y\|^2. \quad (12.33)$$

Then:

$$F(x_{k+1}) \leq \ell_F(x_{k+1}, x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2 \quad (12.34)$$

$$\leq \ell_F(x_{k+1}, x_k) + \frac{L}{\rho} D(x_{k+1}, x_k) \quad (12.35)$$

$$\leq \ell_F(x, x_k) + \frac{L}{\rho} [D(x, x_k) - D(x, x_{k+1})] \quad (12.36)$$

$$\leq F(x) + \frac{L}{\rho} [D(x, x_k) - D(x, x_{k+1})] \quad (12.37)$$

This has the form $e_{k+1} \leq \Delta_k - \Delta_{k+1}$; $e_k = F(x_k) - F(x)$ is decreasing by the descent lemma. By summing and telescoping, we obtain $(k+1)e_{k+1} \leq \Delta_0 - \Delta_{k+1} \leq \Delta_0$. \square

Example 12.11 (Mirror descent on the simplex: the exponentiated gradient algorithm). Recall that the mirror descent iterates are given by $x_{k+1} = \nabla J^*(\nabla J(x_k) - \eta \nabla f(x_k))$, and we have

$$x = \nabla J^*(u) \Leftrightarrow u = \nabla J(x) \Leftrightarrow u = 1 + \log x \Leftrightarrow x = \exp(u - 1). \quad (12.38)$$

So, the mirror descent algorithm with the KL divergence is:

$$x_{k+1} = \exp(1 + \log x_k - \eta \nabla f(x_k) - 1) = x_k \exp(-\eta \nabla f(x_k)), \quad (12.39)$$

which is known as the exponentiated gradient algorithm. When considering a Bregman projection step on the simplex, one adds a normalisation step as seen in [Example 12.6](#).

$$x_{k+1} = \operatorname{argmin}_{x \in \Delta_d} \eta \langle \nabla f(x_k), x \rangle + D(x, x_k) = \frac{x_k \exp(-\eta \nabla f(x_k))}{\sum_{i=1}^d (x_k \exp(-\eta \nabla f(x_k)))_i}. \quad (12.40)$$

12.3 Online learning and mirror descent

The online learning framework is the following: we don't minimize a function but, at time $t \in \mathbb{N}$, take a decision $x_t \in \mathbb{R}^d$, incur cost $f_t(x_t)$ (f_t possibly chosen in an adversarial fashion; we don't know it when picking x_t !), receive/observe $g_t \in \partial f_t(x_t)$ (information: how could we have done better for this stage), and can use it to take a new decision x_{t+1} .

How to measure the performance? The regret of an online algorithm

Definition 12.12 (Regret). The regret at horizon $T \in \mathbb{N}$ of an algorithm/a sequence of actions x_t is:

$$R(T) = \sum_{t=1}^T f_t(x_t) - \min_x \sum_{t=1}^T f_t(x), \quad (12.41)$$

that is, the difference between the incurred cost and the minimal cost for a fixed strategy which knows all the f_t 's in hindsight.

Basic algorithm: subgradient descent. Works exactly as in the fixed f case under the same assumptions on f .

TODO example $f_t = x, 1 - x$ even/odd

Additional references

Nick Harvey lecture notes: <https://www.cs.ubc.ca/~nickhar/F18-531/Notes20.pdf>

Super cool connection: MD is Natural Gradient descent on the dual Riemman manifold (Raskutti and Mukherjee, 2015).

Accelerated mirror descent (Krichene et al., 2015).

Exercise 12.1 (Only symmetric Bregman divergences). *Let J be three times differentiable such that D_J is symmetric. Show that $J(x) = x^\top Ax$ for some matrix A .*

13 Continuous view: gradient flows, mirror flows, and Nesterov

In this section, we take the limit of algorithms as the stepsize goes to zero and times become continuous. Studying the continuous time versions of optimization algorithms offers analogies with the discrete case, while often being less technically involved. It allows to get intuitions on the proof techniques; while the latter cannot always be straightforwardly translated, they can be mimicked.

The most popular example is the gradient flow¹⁷.

13.1 Gradient flow

The gradient flow of f is defined as the solution¹⁸ of the following Ordinary Differential Equation (ODE):

$$\begin{aligned} x(0) &= x_0, \\ \dot{x}(t) &= -\nabla f(x(t)) . \end{aligned} \tag{13.1}$$

Suppose we want to numerically approximate the solution $x(t)$ to this equation. We pick a time discretization step $\eta > 0$, and define time instants $t_k = k\eta$ for $k \in \mathbb{N}$. We will construct a sequence x_k that should approximate the continuous version $x(t_k)$. We approximate the derivative $\dot{x}(t_k)$ by finite differences, as $\frac{x_{k+1} - x_k}{\eta}$.

Thus a discrete approximation of $x(t)$ for $t = \gamma k$ is:

$$x_{k+1} - x_k = -\eta \nabla f(x_k) . \tag{13.2}$$

This is a *forward* Euler discretization scheme, that gives the iterations of gradient descent. Thus, it is not so surprising that the gradient flow and gradient descent iterations will share many properties.

Note that, on the other hand, one could chose a *backward* Euler scheme, also called *implicit* (because it does not give an explicit value of x_k):

$$x_{k+1} - x_k = -\eta \nabla f(x_{k+1}) . \tag{13.3}$$

Rewriting this as $(\text{Id} + \eta \nabla f)(x_{k+1}) = x_k$, it becomes visible that this scheme is the proximal point algorithm! As we shall see, many other discrete algorithms can be obtained as discretizations of various ODEs. Let us come back to the gradient flow.

The gradient flow shares many properties with the iterations of gradient descent.

Proposition 13.1. *The function value $f(x(t))$ decreases along the flow. If f is convex, the gradient norm $\|\nabla f(x(t))\|$ also decreases along the flow.*

Proof. A direct calculation shows:

$$\frac{d}{dt} f(x(t)) = \langle \nabla f(x(t)), \frac{d}{dt} x(t) \rangle = -\|\nabla f(x(t))\|^2 \leq 0 . \tag{13.4}$$

¹⁷nice illustrations and insights at <https://francisbach.com/gradient-flows/>

¹⁸the solution exists and is unique if ∇f is Lipschitz, by the Cauchy-Lipschitz theorem

If f is convex, its Hessian is semidefinite positive and thus:

$$\frac{d}{dt} \|\nabla f(x(t))\|^2 = -\nabla f(x(t)) \nabla^2 f(x(t)) \nabla f(x(t)) \leq 0 . \quad (13.5)$$

□

Proposition 13.2 (Convergence rates). *Let f be convex and x^* be a minimizer of f . Then the gradient flow satisfies*

$$f(x(t)) - f(x^*) \leq \frac{\|x(0) - x^*\|^2}{2t} \quad (13.6)$$

$$\|\nabla f(x(t))\|^2 \leq \frac{f(x(0))}{t} . \quad (13.7)$$

Proof. Again, differentiating a cleverly chosen functional gives:

$$\frac{d}{dt} \frac{1}{2} \|x(t) - x^*\|^2 = -\langle \nabla f(x(t)), x(t) - x^* \rangle \leq f(x^*) - f(x(t)) . \quad (13.8)$$

First, this shows that $x(t)$ gets closer to *any* minimizer as t increases since the RHS is negative. Second, since $f(x(t))$ is decreasing,

$$f(x(t)) - f(x^*) \leq \frac{1}{t} \int_0^t f(x(s)) - f(x^*) \, ds \leq \frac{1}{2t} \|x(0) - x^*\|^2 - \frac{1}{2t} \|x(t) - x^*\|^2 , \quad (13.9)$$

which proves the first result. For the second one, since the squared gradient norm decreases and is equal to $-\frac{d}{dt} f(x(t))$ by [Equation \(13.4\)](#),

$$\|\nabla f(x(t))\|^2 \leq -\frac{1}{t} \int_0^t \frac{d}{ds} f(x(s)) \, ds = \frac{1}{t} (f(x(0)) - f(x(t))) \leq \frac{1}{t} f(x(0)) . \quad (13.10)$$

□

As for gradient descent, if f satisfies the Polyak-Łojasiewicz inequality, we can improve the rate to a linear one.

Proposition 13.3. *Let f satisfy the Polyak-Łojasiewicz inequality with parameter $\mu > 0$, let x^* be a minimizer of f . Then the gradient flow satisfies:*

$$f(x(t)) - f(x^*) \leq (f(x(0)) - f(x^*)) \exp(-2\mu t) . \quad (13.11)$$

Proof. Since by P-Ł, $-\|\nabla f(x)\|^2 \leq -2\mu(f(x) - f(x^*))$, one has

$$\frac{d}{dt} [f(x(t)) - f(x^*)] = -\|\nabla f(x(t))\|^2 \leq -2\mu(f(x(t)) - f(x^*)) \quad (13.12)$$

We can conclude by applying Gronwall's lemma, or integrating directly after making the derivative of $\ln(f(x(t)) - f(x^*))$ appear on the LHS of [Equation \(13.12\)](#):

$$\frac{1}{f(x(t)) - f(x^*)} \frac{d}{dt} (f(x(t)) - f(x^*)) \leq -2\mu . \quad (13.13)$$

□

13.2 Mirror descent flow

Definition 13.4 (Mirror descent flow). *For a Legendre function J , and a convex differentiable function f , the mirror descent flow is the solution of the ODE*

$$\nabla^2 J(x(t))\dot{x}(t) = -\nabla f(x(t)). \quad (13.14)$$

An explicit discretization of this flow with time step η gives the mirror descent equation (Equation (12.15))

Proposition 13.5. *We have:*

1. *for any minimizer x^* of f , $D(x^*, x(t))$ is decreasing*
2. *$f(x(t))$ decreases with t*
3. *$f(x(t)) - f(x^*) \leq \frac{D(x^*, x_0)}{t}$ for any minimizer x^**

Proof. 1. Let $z \in \mathbb{R}^d$.

$$\frac{d}{dt}D(z, x(t)) = \langle \nabla J(x), \dot{x} \rangle - \langle \nabla^2 J(x)\dot{x}, z - x \rangle - \langle \nabla J(x), -\dot{x} \rangle \quad (13.15)$$

$$= \langle -\nabla^2 J(x)\dot{x}, z - x \rangle \quad (13.16)$$

$$= \langle \nabla f(x), z - x \rangle \quad (13.17)$$

$$\leq f(z) - f(x). \quad (13.18)$$

Now taking z to be a minimizer yields $\frac{d}{dt}D(z, x(t)) \leq 0$.

$$\dot{f}(x(t)) = \langle \nabla f(x(t)), \dot{x}(t) \rangle = -\langle \nabla^2 J(x(t))\dot{x}(t), \dot{x}(t) \rangle \leq 0 \quad (13.19)$$

because J is convex so $\nabla^2 J$ is p.s.d. everywhere.

2.

$$\frac{1}{t}(D(x^*, x) - D(x^*, x_0)) = \frac{1}{t} \int_0^t \frac{d}{ds}D(x^*, x(s)) ds \quad (13.20)$$

$$\leq \frac{1}{t} \int_0^t f(x^*) - f(x(s)) ds \quad (13.21)$$

$$\leq \frac{1}{t} \int_0^t f(x^*) - f(x(t)) ds \quad (f(x(s)) \text{ decreases}) \quad (13.22)$$

$$= f(x^*) - f(x(t)) \quad (13.23)$$

In similar settings without decrease of f , one can also use an argument using $f(x^*) - \frac{1}{t} \int_0^t f(x(s)) ds \leq f(x^*) - f(\bar{x}(t))$ with $\bar{x}(t) = \frac{1}{t} \int_0^t x(s) ds$, but here we can do better since f decreases.

□

13.3 An ODE modeling Nesterov acceleration

To study a continuous time version of Nesterov acceleration, in a groundbreaking work, [Su et al. \(2016\)](#) considered the following equation:

$$\ddot{x}(t) + \beta(t)\dot{x}(t) + \nabla(f(x(t))) = 0, \quad (13.24)$$

where $\beta(t)$ is a specific function, to be specified later. We detail how the authors came up with such β . Their goal is to construct a Lyapunov function of the following form

$$\varepsilon(t) = \frac{A(t)}{2} \|x - x^* + C(t)\dot{x}\|^2 + B(t)(f(x(t)) - f^*) \quad (13.25)$$

If A and B are positive functions, and we manage to chose A , B and C such that $\dot{\varepsilon}(t) \leq 0$, then we will have a convergence rate $f(x(t)) - f(x^*) \leq 1/B(t)$.

Heavy but not too complex computation gives (omitting the t variable for brevity):

$$\begin{aligned} \dot{\varepsilon} &= \frac{\dot{A}}{2} \|x - x^*\|^2 + (\dot{A}C + A(\dot{C} + 1 - C\beta)) \langle x - x^*, \dot{x} \rangle - AC \langle x - x^*, \nabla f(x) \rangle \\ &\quad + \left(\frac{\dot{A}C^2}{2} + AC(\dot{C} + 1 - C\beta) \right) \|\dot{x}\|^2 + (B - AC^2) \langle \dot{x}, \nabla f(x) \rangle + \dot{B}(f(x) - f^*) \\ &\leq \frac{\dot{A}}{2} \|x - x^*\|^2 + (\dot{A}C + A(\dot{C} + 1 - C\beta)) \langle x - x^*, \dot{x} \rangle \\ &\quad + \left(\frac{\dot{A}C^2}{2} + AC(\dot{C} + 1 - C\beta) \right) \|\dot{x}\|^2 + (B - AC^2) \langle \dot{x}, \nabla f(x) \rangle + (\dot{B} - AC)(f(x) - f^*) \end{aligned} \quad (13.26)$$

using $-AC \langle x - x^*, \nabla f(x) \rangle \leq -AC(f(x) - f(x^*))$ by convexity. To simplify the expression, we take many terms equal to 0 by choosing $\dot{A} = 0$, $B = AC^2$, $\dot{B} = AC$ and $\dot{C} + 1 - C\beta = 0$. Since A is a constant and $B = AC^2$, $\dot{B} = 2AC\dot{C}$, which with the requirement $\dot{B} = AC$ yields $\dot{C} = 1/2$. Thus we take $C(t) = t/2$ and $\beta(t)$ must be equal to $3/t$.

By this choice, one obtains:

$$f(x(t)) - f(x^*) = \mathcal{O}\left(\frac{1}{B(t)}\right) = \mathcal{O}\left(\frac{1}{t^2}\right). \quad (13.28)$$

One possible discretization of the ODE (13.24) is:

$$\frac{x_{k+1} - 2x_k + x_{k-1}}{h^2} + \frac{3}{kh} \frac{x_{k+1} - x_k}{h} + \nabla f(x_k) = 0 \quad (13.29)$$

$$x_{k+1} = x_k + \left(1 - \frac{3}{k}\right)(x_k - x_{k-1}) - h^2 \nabla f(x_k) \quad (13.30)$$

This mysterious 3 has a strong influence on the properties of the solution TODO one can also obtain $\frac{k-1}{k+2}$.

Remember that the choice of momentum in Nesterov's algorithm (non strongly convex version) is given by $t_1 = 1$

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2} \quad (13.31)$$

, such that $t_{k+1}^2 - t_{k+1} = t_k$. This is actually the lower choice. Remember that in the proof one only needs $t_{k+1}^2 - t_{k+1} \leq t_k$, the choice of Nesterov is the smallest sequence satisfying this property. Since the work of [Chambolle and Dossal \(2015\)](#) (which studies FISTA but whose results straightforwardly apply to the case without prox), it has become popular to use a momentum term equal to $\frac{k}{k+\alpha}$, with the value α critically controlling the convergence of the algorithm (according to three cases: $\alpha < 3$, $\alpha = 3$, $\alpha > 3$).

Actually, one can prove that $t_k \sim k/2$ and so asymptotically the momentum term with Nesterov's choice is $\frac{k/2-1}{(k+1)/2} = \frac{k-2}{k+1} = 1 - \frac{3}{k} + o(\frac{1}{k})$, so the choices are “equivalent”.

13.4 Subgradient flow

The subgradient flow of f is the nonsmooth equivalent of the gradient flow. It is defined as the solution of the following differential inclusion:

$$-\dot{x} \in \partial f(x(t)) \quad (13.32)$$

with initial condition $x(0) = x_0$. Interestingly, it follows a “laziness principle”: even though there are potentially infinitely many possible choices of subgradient, only one is taken: the one with minimal norm.

This result, stated in [Proposition 13.7](#), is based on the following chain rule for subgradients.

Proposition 13.6. *Let I an open interval, $t \in I$, $x : I \rightarrow X$, $f : X \rightarrow \bar{\mathbb{R}}$ such that x and $f \circ x$ are differentiable at t , and $x(t) \in \text{dom } \partial f$. Then:*

$$(f \circ x)'(t) = \langle g, \dot{x}(t) \rangle, \quad \forall g \in \partial f(x(t)) \quad (13.33)$$


See 2.2.11 in *G. Garrigos' PhD thesis* for the proof.

Proposition 13.7. *Let x be a solution of (13.32). Then for all t ,*

$$-\dot{x} = \underset{u \in \partial f(x(t))}{\operatorname{argmin}} \|u\|. \quad (13.34)$$

Proof. One has by definition that $-\dot{x} \in \partial f(x(t))$. By [Proposition 13.6](#), for every $g \in \partial f(x(t))$, $\langle g, \dot{x}(t) \rangle = (f \circ x)'(t) = -\|\dot{x}(t)\|^2$, where the last equality follows from taking the particular $g = -\dot{x}$. By applying the Cauchy-Schwarz inequality, we obtain $-\|\dot{x}\|^2 \geq -\|\dot{x}\|\|g\|$, hence $\|g\| \geq \|\dot{x}\|$, which concludes the proof (the argmin is reduced to a single point because the subdifferential is closed and convex). \square

13.5 Exercises

Exercise 13.1 (Preliminaries).  Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be twice differentiable. Let $\theta : \mathbb{R}_+ \rightarrow \mathcal{X}$ be differentiable. Show that

$$\begin{aligned} \frac{d}{dt} f(\theta(t)) &= \langle \nabla f(\theta(t)), \dot{\theta}(t) \rangle, \\ \frac{d}{dt} \nabla f(\theta(t)) &= \nabla^2 f(\theta(t)) \dot{\theta}(t). \end{aligned}$$

Exercise 13.2 (Everything decreases in gradient flow). ☕ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex twice differentiable function. Let $x(t)$ be the gradient flow on f , defined as a solution of $\dot{x}(t) = -\nabla f(x(t))$.

Show that $f(x(t))$ decreases.

Show that $\|\nabla f(x(t))\|$ decreases.

Show that $\|x(t) - x^*\|$ decreases for any minimizer x^* of f .

Exercise 13.3. ☕ Do the results of [Exercise 13.2](#) hold when f is not convex?

14 Bilevel optimization

Bilevel optimization is concerned with minimizing a function with respect to a variable that it itself the solution of an optimization problem.

For example, consider the very simple example of finding the best regularization strength λ for ridge regression. We may formulate it as a bilevel optimization problem:

$$\begin{aligned} \min_{\lambda \geq 0} \quad & \frac{1}{2} \|A_{\text{test}} x^*(\lambda) - b_{\text{test}}\|^2 \\ \text{s.t.} \quad & x^*(\lambda) = \operatorname{argmin}_x \frac{1}{2} \|A_{\text{train}} x^*(\lambda) - b_{\text{train}}\|^2 + \frac{\lambda}{2} \|x\|^2 \end{aligned} \quad (14.1)$$

More generally, we formulate a bilevel optimization problem as:

$$\min_{\theta} F(v^*(\theta)) \quad \text{s.t.} \quad v^*(\theta) = \operatorname{argmin}_v G(v, \theta) \quad (14.2)$$

The minimization problem in θ is called the *upper level problem* and F is the upper objective, while the minimization problem in v is called the *lower level problem* and G is the lower objective.

Remark 14.1. *Problem (14.2) is not the most general form of a bilevel optimization problem:*

- *the outer objective function F can be a function of θ too (i.e. we consider $\min_{\theta} F(v^*(\theta), \theta)$); this is not a difficulty*
- *the inner problem can be constrained, with a constraint set depending on θ*
- *the inner problem can have multiple solutions; one then needs the notion of optimistic and pessimistic bilevel optimization, that we do not handle in this class*

14.1 The implicit approach

The most standard approach for solving [Problem \(14.1\)](#) is grid search: pick a list of candidate values for λ , $(\lambda_1, \dots, \lambda_T)$, compute $x^*(\lambda_t)$ for $t \in [T]$ and pick the λ yielding the lowest test error. This is a *zeroth-order* method for the bilevel problem.

Rather than performing zeroth order optimization, which scales very badly with the dimension of θ (think of the Ridge regularization case, but with one regularization strength per feature) one could think of performing gradient descent on F , requiring the computation of the *hypergradient* $\nabla_{\theta} F(v^*(\theta))$. By the chain rule, one has:

$$\mathcal{J}F \circ v^*(\theta) = \mathcal{J}F(v^*(\theta)) \mathcal{J}v^*(\theta) \quad (14.3)$$

$$(\nabla F \circ v^*(\theta))^{\top} = (\nabla F(v^*(\theta)))^{\top} \mathcal{J}v^*(\theta) \quad (14.4)$$

$$\nabla F \circ v^*(\theta) = (\mathcal{J}v^*(\theta))^{\top} \nabla F(v^*(\theta)). \quad (14.5)$$

where we have used the fact that $\theta \mapsto F(v^*(\theta))$ and $v \mapsto F(v)$ are scalar functions, hence their Jacobians are the transpose of their gradients. For clarity, beware of the difference between $\mathcal{J}(F \circ v^*)(\theta)$ the Jacobian of $\theta \mapsto F(v^*(\theta))$, on the left hand side, and $\mathcal{J}F(v^*(\theta))$ the Jacobian of $v \mapsto F(v)$ evaluated at $v^*(\theta)$, on the right hand side.

In the right hand side, $\nabla F(v^*(\theta))$ is usually easy to compute, as we know the form of F . However, computing the Jacobian of v^* with respect to θ is hard if one does not have access to a closed form for v^* .

A first solution is to exploit the first-order optimality condition for the inner problem, and differentiate with respect to θ . To justify differentiation (under technical assumptions) we rely on the Implicit function theorem, hence this technique is called *implicit differentiation*. One has by optimality of $v^*(\theta)$ for the inner problem:

$$\nabla_1 G(v^*(\theta), \theta) = 0 \quad (14.6)$$

where ∇_1 denotes the gradient (G is scalar-valued) with respect to the first variable. Differentiating with respect to θ :

$$\nabla_1^2 G(v^*(\theta), \theta) \mathcal{J}v^*(\theta) + \nabla_{12}^2 G(v^*(\theta), \theta) = 0 \quad (14.7)$$

where ∇_1^2 is the Hessian with respect to the first variable, and ∇_{12}^2 is the Jacobian with respect to the second variable of the gradient with respect to the first variable. To get this notation straight, remember this block matrix summary:

$$\nabla^2 G(a, b) = \begin{pmatrix} \nabla_1^2 G(a, b) & \nabla_{12}^2 G(a, b) \\ \nabla_{21}^2 G(a, b) & \nabla_2^2 G(a, b) \end{pmatrix} \quad (14.8)$$

where the first diagonal block is of size $\dim v$ and the second one is of size $\dim \theta$.

Hence, we can compute the desired Jacobian $\mathcal{J}v^*(\theta)$ by solving a linear system.

$$\nabla_1^2 G(v^*(\theta), \theta) \mathcal{J}v^*(\theta) = -\nabla_{12}^2 G(v^*(\theta), \theta) \quad (14.9)$$

and then perform gradient descent on the bilevel optimization problem.

In practice, many difficulties arise:

- one does not have access to $v^*(\theta)$, but only to some approximation. Usually, this approximation is $v^{(K)}(\theta)$, obtained by performing K steps of an iterative algorithm. One can then use automatic differentiation (Section 16) to compute $\mathcal{J}v^{(T)}(\theta)$. Even if $v^{(T)} \rightarrow v^*$, guaranteeing the convergence of the Jacobian is a tricky problem...
- the solution may not be differentiable, for example when G is not smooth, which arises frequently in regularized machine learning problems. See [Bertrand et al. \(2022\)](#) for a solution in that case.

See [Pedregosa \(2016\)](#).

Generally, most methods write the problem as trying to find a triplet (v, θ, u) satisfying:

$$\begin{cases} \nabla_1 G(v^*(\theta), \theta) = 0 \\ \nabla_1^2 G(v^*(\theta), \theta) u = \nabla F(v^*(\theta)) \\ \nabla_{12}^2 G(v^*(\theta), \theta)^\top u = 0 \end{cases} \quad (14.10)$$

and perform some variants of fixed point iterations on this system. This system is derived as follows: Introduce u satisfying $\nabla_1^2 G(v^*(\theta))u = \nabla F(v^*(\theta))$. Then the optimality on the outer

problem writes:

$$Jv^*(\theta)^\top \nabla F(x^*) = 0 \quad (14.11)$$

$$\nabla_{12}^2 G(v^*(\theta), \theta) \nabla_1^2 G(v^*(\theta), \theta)^{-1} \nabla F(x^*) = 0 \quad (14.12)$$

$$\nabla_{12}^2 G(v^*(\theta), \theta) u = 0 \quad (14.13)$$

(this is informal because of the inversion which is not really needed).

Thus the first condition in the system is the optimality condition for the inner problem, the second comes from the optimality for the outer, and the third one comes from the expression of the Jacobian of v^* .

14.2 The value fonction approach

[Liu et al. \(2022\)](#)

14.3 See also

[Bolte et al. \(2024\)](#) for “one step differentiation”.

15 Implicit regularization

Resources : Claire Boyer lecture notes: "This chapter heavily relies on the articles Hastie et al. (2019), Bartlett et al. (2021) and the lecture of Matus Telgarsky"

In classical statistical learning, it is often nice to control the norm of the solution in regression, or the margin in classification (Shalev Schwartz 2014 Understanding ML book). It is typically imposed by regularizing the objective function, using an additional square L2 term for example.

However, we have witnessed the advent of deep learning successes, where highly overparameterized, unregularized models which should have been prone to overfitting met stellar success. This has sparked new interest in the phenomenon of *implicit regularisation*, in which, amongst all solutions to an optimization problem, an algorithm always converges to a particular one – in this case, a “simple”, “good” minimizer of the empirical risk that generalizes well.

15.1 Implicit bias of GD on least squares

Consider the least square problem:

$$\min_x \frac{1}{2} \|Ax - b\|^2 \quad (15.1)$$

and the iterations of gradient descent on this problem, started at 0:

$$\begin{cases} x_0 = 0 \\ x_{k+1} = x_k - \eta A^\top (Ax_k - b) \end{cases} \quad (15.2)$$

A simple recursion shows that $x_k = \sum_0^k (\text{Id} - \eta A^\top A)^t \eta A^\top b$; if $\eta < 2/\|A\|^2$, the matrix series (called von Neumann series) converges¹⁹, to $(\text{Id} - (\text{Id} - \eta A^\top A)^\dagger)$ and so the iterates converge to $A^\dagger b$.

This solution is remarkable: it is the minimum norm solution to $A^\top Ax = A^\top b$ ²⁰.

There are many ways to show it (it can even be considered the definition of the pseudo inverse operator, I think). There are geometrical proofs and a nice proof by going to the dual TODO exercise.

Therefore, amongst all possible solutions to the linear system $A^\top Ax = A^\top b$, gradient descent “unknowingly” converges to the minimal norm one, without regularization.

What if one is interested in solutions that are minimal in other sense, e.g. minimize another functional than the Euclidean norm?

15.2 “Implicit” bias of mirror descent

J smooth and α -strongly convex, implicit bias of mirror descent on least squares? Mirror descent iterations:

¹⁹in the same way $\sum_0^{+\infty} \rho^k = (1 - \rho)^{-1}$

²⁰recall that in finite dimension there always exists such a solution, even when there is no solution to $Ax = b$; see [Exercise 1.1](#)

Proposition 15.1. *Consider the mirror descent iterations on least squares with potential smooth and strongly convex potential J :*

$$\nabla J(x_{k+1}) = \nabla J(x_k) - \eta A^*(Ax_k - b)$$

Then the iterates converge to the “ J -minimal solution”:

$$x_k \rightarrow \underset{x}{\operatorname{argmin}} J(x) \quad \text{s.t.} \quad Ax = b \quad (15.3)$$

Proof. Write $J = j + \frac{\alpha}{2} \|\cdot\|^2$ with j convex and $\alpha > 0$ (coming from strong convexity of J). Consider the explicit “min- J solution” problem and its dual:

$$\text{Primal : } \hat{x} = \underset{x}{\operatorname{argmin}} J(x) = j(x) + \frac{\alpha}{2} \|x\|_2^2 \quad \text{s.t.} \quad Ax = b \quad (15.4)$$

$$\text{Dual : } \hat{\theta} = \underset{\theta}{\operatorname{argmin}} J^*(-A^*\theta) + \langle b, \theta \rangle \quad (15.5)$$

$$\text{Link : } -A^*\hat{\theta} = \alpha \hat{x} + \nabla j(\hat{x}) \quad \text{aka} \quad \hat{x} = \operatorname{prox}_{\frac{1}{\alpha}j}(-A^*\hat{\theta}/\alpha) \quad (15.6)$$

J being strongly convex makes J^* smooth²¹. Actually, J^* is the Moreau envelope of j^* :

$$\begin{aligned} J^* &= (j + \frac{\alpha}{2} \|\cdot\|^2)^* \\ &= j^* \square \frac{1}{2\alpha} \|\cdot\|^2 \quad (\text{Fenchel conjugate of sum}) \\ &= \inf_y j^*(y) + \frac{1}{2\alpha} \| \cdot - y \|^2 \quad (\text{this is a Moreau envelope}) \end{aligned}$$

Since J is strongly convex, the dual is smooth and you can perform gradient descent on it. The gradient of J^* writes as a prox in this case. The dual is smooth, one can apply gradient descent on it:

$$\theta_{k+1} = \theta_k - \eta[-A\nabla J^*(-A^*\theta_k) + b] \quad (15.7)$$

$$= \theta_k - \eta[-A \operatorname{prox}_{\frac{1}{\alpha}j}(-A^*\theta_k) + b] \quad (\text{envelope theorem + Moreau decomposition}) \quad (15.8)$$

The formula for the gradient of Moreau envelope (??) gives:

$$\begin{aligned} \nabla J^*(x) &= \frac{1}{\alpha}(x - \operatorname{prox}_{\alpha j}(x)) \\ &= \operatorname{prox}_{\frac{1}{\alpha}j}(x/\alpha) \quad (\text{Moreau decomposition formula}) \end{aligned}$$

The dual iterates θ_k solve the dual, the final step is to map them back to the primal using the primal-dual link equation. Inspired by the link equation, let

$$x_k \triangleq \operatorname{prox}_{\frac{1}{\alpha}j}(-A^*\theta_k/\alpha), \text{ meaning } -A^*\theta_k = \alpha x_k + \nabla j(x_k) = \nabla J(x_k) \quad (15.9)$$

Plugging back x_k in dual GD iterates gives mirror descent:

$$\begin{aligned} -A^*\theta_{k+1} &= -A^*\theta_k + \eta A^*[-A\nabla J^*(-A^*\theta_k) + b] \\ \nabla J(x_{k+1}) &= \nabla J(x_k) - \eta A^*(Ax_k - b) \quad (\text{mirror descent}) \end{aligned}$$

Since $\theta_k \rightarrow \hat{\theta}$, $x_k \rightarrow \hat{x}$. □

²¹In fact, I reckon we only need strict convexity of J (in which case the interpretation as a proximal step for the primal variable x does not hold)

15.3 Logistic regression on separable data

Motivation: controlling norm avoids overfitting, controls *stability* of solution. Typical example: Tikhonov, bound singular value.

$$(A^t A)^\dagger A^\top b = \sum_1^R \frac{1}{\sigma_i} v_i u_i^\top b$$

vs

$$(A^t A + \lambda \text{Id})^\dagger A^\top b = \sum_1^R \frac{\sigma_i}{(\sigma_i + \lambda)^2} v_i u_i^\top b$$

Iterative regularization Engl 1996

Implicit bias of GD on OLS. Generalization to Mirror descent

TODO: give SVD as homework

1. the geometrical view with SVD for OLS: projection onto $\text{Span } A^\top$

Ref off the convex path: http://www.offconvex.org/2020/11/27/reg_dl_not_norm/

15.4 Matrix factorization: can everything be explained in terms of norms?

Deep linear nets for matrix factorization, ref <http://www.offconvex.org/2019/07/10/trajectories-linear-nets/>

Neyshabur thesis on implicit bias in DL: <https://arxiv.org/pdf/1709.01953.pdf>

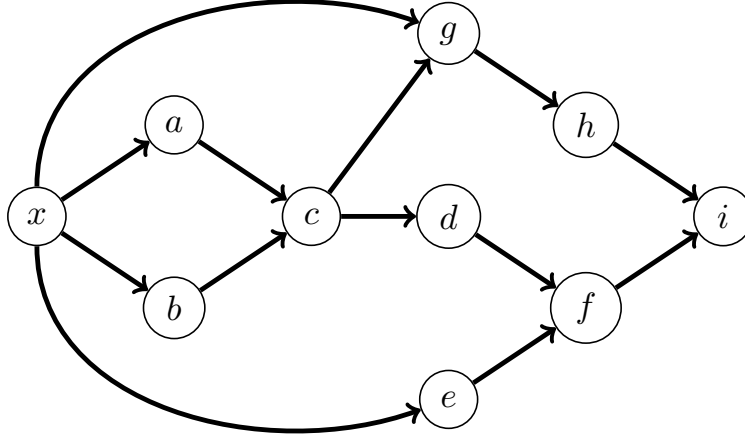


Figure 16.1: Directed acyclic graph representing the dependency graph of [Algorithm 1](#)

16 Automatic differentiation

First-order optimization methods require computing gradients of given functions. So far, all the functions were simple (least squares, logistic regression) and their gradient could be computed by hand.

We now turn to examples where the computation requires more advanced tools.

16.1 Forward and backward automatic differentiation

Note: AD acts on programs, not on functions.

Consider the following scalar function:

$$f(x) = \exp(x^2 - 3x) \sin(x) + \cos(x^2 - 3x)x \quad (16.1)$$

It can be implemented through the following program in [Algorithm 1](#), which is itself representable by a Directed Acyclic Graph (DAG) in [Figure 16.1](#).

Algorithm 1 Program implementing f

```

a = x2
b = 3x
c = a + b
d = exp(c)
e = sin(x)
f = de
g = cos(c)
h = gx
i = hf

```

In this section, we will introduce three ways to compute gradient/Jacobians of loss functions

Tuning the weights of a neural network requires computing the Jacobian of the loss with respect to various parameters. Automatic differentiation, in particular Autodiff, backpropagation

Autodiff for a one layer NN:

$$f(x) = W_2 \sigma(W_1 x) \quad (16.2)$$

Technique: vectorize everything and devectorize. Code snippet. Emphasize vector Jacobian product and Jacobian vector product (directional derivatives)

17 Variational inequalities

In all this section K is a convex bounded subset of \mathbb{R}^d . The function $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is called an operator.

Definition 17.1 (Strong and weak solutions to variational inequality problems). *We say that $x^* \in \mathbb{R}^d$ is a strong solution to the variational inequality problem, or simply to the variational inequality, associated to F (and K) if $x^* \in K$ and*

$$\forall y \in K, \langle F(x^*), x^* - y \rangle \leq 0 . \quad (17.1)$$

It is a weak solution if it belongs to K and

$$\forall y \in K, \langle F(y), x^* - y \rangle \leq 0 . \quad (17.2)$$

In the unconstrained case, a strong solution is simply a point such that $F(x^*) = 0$. Variational inequalities generalize convex (constrained) minimization problems: to solve $\min_K f(x)$ is to find a strong solution to the VI associated to $F = \nabla f$.

Under certain conditions, strong and weak solutions are the same; these conditions involve a special property on F , monotonicity.

Definition 17.2 (Monotone operator). *The operator F is monotone if*

$$\forall x, y, \langle F(x) - F(y), x - y \rangle \geq 0 . \quad (17.3)$$

The most famous class of monotone operators is probably the one composed of the convex functions' gradients (see [Exercise 2.9](#)).

Proposition 17.3 (Strong-weak equivalence under monotonicity of F). *Let F be monotone. Then any strong solution to the VI is a weak solution to the VI. The converse is also true if F is continuous and K is convex, a result known as Minty's lemma²².*

Proof. The strong \implies weak direction is easy. To prove that any weak solution is a strong solution, let x^* be a weak solution, and let $y \in K$. Let $t \in]0, 1]$. Since K is convex, $x^* + t(y - x^*) \in K$, and so:

$$\langle F(x^* + t(x^* - y)), x^* - x^* + t(x^* - y) \rangle \leq 0 \quad (17.4)$$

$$t \langle F(x^* + t(x^* - y)), x^* - y \rangle \leq 0 \quad (17.5)$$

$$\langle F(x^* + t(x^* - y)), x^* - y \rangle \leq 0 . \quad (17.6)$$

²²everything here is in finite dimension, but beware that the result may not hold in infinite dimension. Can you spot where this plays a role in the proof?

Taking the limit as $t \rightarrow 0$, by continuity of F , yields that x^* is a strong solution since y was any point. \square

Variational inequalities generalize variational optimization problems, but there is no objective anymore.

They handle more complex problems, such as min-max problems arising in game theory.

Example 17.4 (Rock paper scissors). *The rock-paper-scissors game is described by the payoff*

matrix $A = \begin{pmatrix} 0 & -1 & 1 \\ 1 & 0 & 1 \\ -1 & 1 & 0 \end{pmatrix}$. If Rock, Paper and Scissors choices are represented by 1, 2 and 3 respectively, then if player 1 plays i and player 2 plays j , player 1 receives payoff A_{ij} (and player 2 receives $-A_{ij}$).

If Player 1 (P1) chooses a random strategy represented by a vector p in the 3-d simplex Δ_3 and Player 2 (P2) chooses $q \in \Delta_3$, then the expected payoff of P1 is:

$$p^\top Aq = \sum_{i,j=1}^3 p_i A_{ij} q_j . \quad (17.7)$$

The game for P1 is to solve $\max_p \min_q p^\top Aq$, while for P2 it is to solve $\max_q \min_p -p^\top Aq$, or equivalently $\min_q \max_p p^\top Aq$.

A result, due to von Neumann and reminiscent of [Section 8](#), states that the two objectives are the same:

$$\max_{p \in \Delta_3} \min_{q \in \Delta_3} p^\top Aq = \min_{q \in \Delta_3} \max_{p \in \Delta_3} p^\top Aq . \quad (17.8)$$

More generally, variational inclusions are useful to model optimization problems of the form:

$$\min_u \max_v f(u, v) , \quad (17.9)$$

where f is differentiable, convex in u , and concave in v . Indeed, such a problem is equivalent to a VI with $F = (\nabla_u f, -\nabla_v f)$.

To solve [Problem \(17.9\)](#), one may be tempted to perform gradient descent on u and gradient ascent on v simultaneously:

$$\begin{cases} u_{k+1} = u_k - \eta \nabla_u f(u_k) \\ v_{k+1} = v_k + \eta \nabla_v f(v_k) \end{cases} \quad (17.10)$$

but unfortunately, this fails to converge even in very simple cases.

Example 17.5. *Consider the 2D problem $\min_{u \in \mathbb{R}} \max_{v \in \mathbb{R}} uv$, whose solution is $(0, 0)$. The iterates [\(17.10\)](#) on this problem read:*

$$\begin{cases} u_{k+1} = u_k - \eta v_k \\ v_{k+1} = v_k + \eta u_k \end{cases} . \quad (17.11)$$

It is easy to check that $\|(u_{k+1}, v_{k+1})\| = \sqrt{(1 + \eta^2)} \|(u_k, v_k)\|$. Hence, except in the trivial case $\eta = 0$, the iterates will always diverge (worse: they get further away from the solution at every iteration). This is visible on [Figure 17.1](#).

We must therefore come up with a better method!

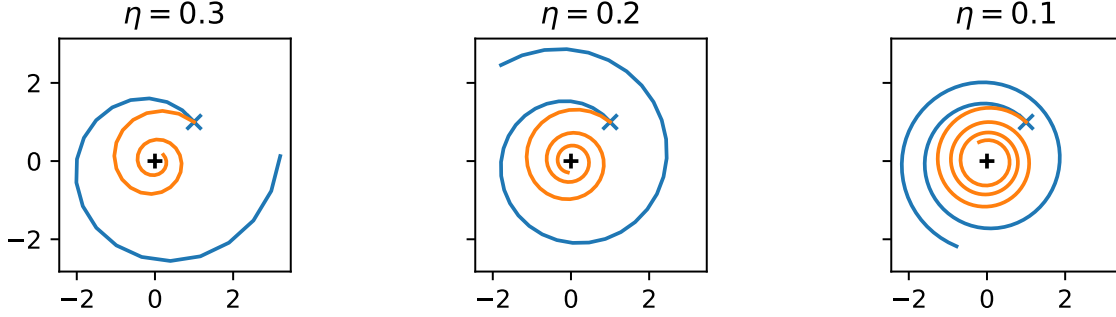


Figure 17.1: Gradient method (blue) and extragradient method (orange) for the saddle point problem $\min_{u \in \mathbb{R}} \max_{v \in \mathbb{R}} uv$, initialized at $(1, 1)$. No matter how small η , the gradient method iterates always go to infinity. On the contrary, the extragradient iterates converge to the solution, $(0, 0)$.

17.1 Solving VIs: the extragradient method

Definition 17.6 (The extragradient method). *The extragradient method²³ iterates are given by:*

$$\begin{cases} x_k = z_{k-1} - \eta F(z_{k-1}) \\ z_k = z_{k-1} - \eta F(x_k) = z_{k-1} - \eta F(z_{k-1} - \eta F(z_{k-1})) \end{cases} \quad (17.12)$$

If we take the special case $F = \nabla f$ to get back to a more familiar situation, this method can be interpreted as first taking a gradient descent step from z_k , then computing the gradient at the resulting step x_k , and using this gradient to take a gradient step on the original point z_k . See ??.

To analyze its convergence properties, we must ask ourselves, what is a relevant measure of error, since there is no objective involved. Guided by the definition of a weak solution, we will use the following value to quantify the quality of x :

$$\sup_{y \in K} \langle F(y), x - y \rangle. \quad (17.13)$$

Proposition 17.7 (Convergence rate of the extragradient method). *They have the following convergence rates on the ergodic iterates $\bar{x}_t = \frac{1}{t} \sum_1^t x_k$:*

- if F is bounded in norm by G and $\eta = \dots$, $\text{err } \bar{x}_t \leq \text{TODO}$
- if F is β Lipschitz and $\eta < 1/\beta$, rate is $/T$

Proof. In the whole proof, $y \in K$ is fixed. Let $t \in \mathbb{N}$, $k \in \llbracket 1, t \rrbracket$ By monotonicity of F ,

$$\langle F(y), x_k - y \rangle \leq \langle F(x_k), x_k - y \rangle. \quad (17.14)$$

²³we cannot talk about extragradient *descent* since there is nothing being minimized here

Now, we want to use the definitions of x_k and z_k , and so we must make $x_{k+1} - x_k$ and TODO appear. We thus write:

$$\begin{aligned}\langle F(x_k), x_k - y \rangle &= \langle F(x_k), z_k - y \rangle + \langle F(z_{k-1}), x_k - z_k \rangle + \langle F(x_k) - F(z_{k-1}), x_k - z_k \rangle \\ &= \frac{1}{\eta} \langle z_{k+1} - z_k, z_k - y \rangle + \frac{1}{\eta} \langle z_{k-1} - x_k, x_k - z_k \rangle + \langle F(x_k) - F(z_{k-1}), x_k - z_k \rangle .\end{aligned}\tag{17.15}$$

We use the parallelogram identity on the two first terms of the right-hand side:

$$\langle z_{k+1} - z_k, z_k - y \rangle = \frac{1}{2} (\|z_{k+1} - y\|^2 - \|z_{k+1} - z_k\|^2 - \|z_k - y\|^2) \tag{17.16}$$

$$\langle z_{k-1} - x_k, x_k - z_k \rangle = \frac{1}{2} (\|z_{k+1} - z_k\|^2 + \|z_{k+1} - x_k\|^2 + \|x_k - z_k\|^2) \tag{17.17}$$

Summing, the second term of the RHS of the first line and the first term of the RHS of the second line cancel out. So the RHS in (17.15) is equal to:

$$\frac{1}{2\eta} (\|z_{k+1} - y\|^2 - \|z_k - y\|^2) - \frac{1}{2\eta} (\|z_{k+1} - x_k\|^2 + \|x_k - z_k\|^2) + \langle F(x_k) - F(z_{k-1}), x_k - z_k \rangle \tag{17.18}$$

The first term will telescope when we sum over k . It remains to handle the last one, involving F .

Bounded case: by applying Cauchy-Schwarz inequality, the triangular inequality, and Young inequality, we obtain

$$\langle F(x_k) - F(z_{k-1}), x_k - z_k \rangle \leq 2G\|x_k - z_k\| \leq 2\eta G^2 + \frac{\|x_k - z_k\|^2}{2\eta} . \tag{17.19}$$

Plug back in RHS, two terms cancel out. Summing over k from 1 to t , dividing by t and choosing η accordingly concludes the proof.

Lipschitz case:

$$\langle F(x_k) - F(z_{k-1}), x_k - z_k \rangle \leq \beta \|x_k - z_{k-1}\| \|x_k - z_k\| \leq \frac{\beta}{2} \|x_k - z_{k-1}\|^2 + \frac{\beta}{2} \|x_k - z_k\|^2 . \tag{17.20}$$

(17.18) is thus upper bounded by

$$\frac{1}{2\eta} (\|z_{k+1} - y\|^2 - \|z_k - y\|^2) + \left(\frac{1}{2\eta} - \frac{\beta}{2} \right) (\|z_{k+1} - x_k\|^2 + \|x_k - z_k\|^2) . \tag{17.21}$$

The second term is negative provided $\eta \leq 1/\beta$; summing, telescoping and dividing by t concludes similarly. \square

18 Additional bibliographic resources

- [Bansal and Gupta \(2017\)](#) for a review of potential-based convergence proofs.

References

- Benjamin Grimmer. Provably faster gradient descent via long steps. *arXiv preprint arXiv:2307.06324*, 2023.
- Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- Baptiste Goujaud, Adrien Taylor, and Aymeric Dieuleveut. Provable non-accelerations of the heavy-ball method. *arXiv preprint arXiv:2307.11291*, 2023.
- Euhanna Ghadimi, Hamid Reza Feyzmahdavian, and Mikael Johansson. Global convergence of the heavy-ball method for convex optimization. In *2015 European control conference (ECC)*, pages 310–315. IEEE, 2015.
- Vassilis Apidopoulos, Nicolò Ginatta, and Silvia Villa. Convergence rates for the heavy-ball continuous dynamics for non-convex optimization, under polyak–łojasiewicz condition. *Journal of Global Optimization*, 84(3):563–589, 2022.
- J-F Aujol, Ch Dossal, and A Rondepierre. Convergence rates of the heavy-ball method under the łojasiewicz property. *Mathematical Programming*, 198(1):195–254, 2023.
- Hedy Attouch and Jalal Fadili. From the ravine method to the nesterov method and vice versa: a dynamical system perspective. *SIAM Journal on Optimization*, 32(3):2074–2101, 2022.
- Brendan O’Donoghue and Emmanuel Candes. Adaptive restart for accelerated gradient schemes. *Foundations of computational mathematics*, 15:715–732, 2015.
- Olivier Fercoq and Zheng Qu. Adaptive restart of accelerated gradient methods under local quadratic growth condition. *IMA Journal of Numerical Analysis*, 39(4):2069–2095, 2019.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Mark Schmidt, Nicolas Roux, and Francis Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. *Advances in neural information processing systems*, 24, 2011.
- Francis Bach, Rodolphe Jenatton, Julien Mairal, Guillaume Obozinski, et al. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4(1): 1–106, 2012.
- Bang Công Vũ. A splitting algorithm for dual monotone inclusions involving cocoercive operators. *Advances in Computational Mathematics*, 38:667–681, 2013.

- Laurent Condat. A primal–dual splitting method for convex optimization involving lipschitzian, proximable and linear composite terms. *Journal of optimization theory and applications*, 158(2):460–479, 2013.
- Ernest K Ryu and Stephen Boyd. Primer on monotone operator methods. *Appl. comput. math*, 15(1):3–43, 2016.
- Ernest K Ryu and Wotao Yin. *Large-scale convex optimization: algorithms & analyses via monotone operators*. Cambridge University Press, 2022.
- Patrick L Combettes and Jean-Christophe Pesquet. Fixed point strategies in data science. *IEEE Transactions on Signal Processing*, 69:3878–3905, 2021.
- Arian Berdellima. On a notion of averaged operators in cat (0) spaces. *arXiv preprint arXiv:2010.05726*, 2020.
- Nobuhiko Ogura and Isao Yamada. Non-strictly convex minimization over the fixed point set of an asymptotically shrinking nonexpansive mapping. 2002.
- Patrick L Combettes. Fejér-monotonicity in convex optimization. *Encyclopedia of optimization*, 2:106–114, 2001.
- Guillaume Garrigos and Robert M Gower. Handbook of convergence theorems for (stochastic) gradient methods. *arXiv preprint arXiv:2301.11235*, 2023.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- David Williams. *Probability with martingales*. Cambridge university press, 1991.
- Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *arXiv e-prints*, pages arXiv–1309, 2013.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in neural information processing systems*, 27, 2014.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26, 2013.
- RT Rockafellar. Convex analysis. *Princeton Math. Series*, 28, 1970.
- Paul Tseng. Approximation accuracy, gradient methods, and error bound for structured convex optimization. *Mathematical Programming*, 125(2):263–295, 2010.
- Garvesh Raskutti and Sayan Mukherjee. The information geometry of mirror descent. *IEEE Transactions on Information Theory*, 61(3):1451–1457, 2015.

- Walid Krichene, Alexandre Bayen, and Peter L Bartlett. Accelerated mirror descent in continuous and discrete time. *NeuRIPS*, 28, 2015.
- Weijie Su, Stephen Boyd, and Emmanuel J Candes. A differential equation for modeling nesterov’s accelerated gradient method: Theory and insights. *Journal of Machine Learning Research*, 17(153):1–43, 2016.
- Antonin Chambolle and Charles H Dossal. On the convergence of the iterates of” fista”. *Journal of Optimization Theory and Applications*, 166(3):25, 2015.
- Quentin Bertrand, Quentin Klopfenstein, Mathurin Massias, Mathieu Blondel, Samuel Vaiter, Alexandre Gramfort, and Joseph Salmon. Implicit differentiation for fast hyperparameter selection in non-smooth convex learning. *Journal of Machine Learning Research*, 23(149):1–43, 2022.
- Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *International conference on machine learning*, pages 737–746. PMLR, 2016.
- Bo Liu, Mao Ye, Stephen Wright, Peter Stone, and Qiang Liu. Bome! bilevel optimization made easy: A simple first-order approach. *Advances in neural information processing systems*, 35:17248–17262, 2022.
- Jérôme Bolte, Edouard Pauwels, and Samuel Vaiter. One-step differentiation of iterative algorithms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Nikhil Bansal and Anupam Gupta. Potential-function proofs for first-order methods. *arXiv preprint arXiv:1712.04581*, 2017.

A The singular value decomposition

In this section we prove [Proposition 0.12](#).

First, observe that $A^\top A$ is a symmetric semidefinite positive matrix. The symmetry is clear, and $A^\top A$ is positive semidefinite as for any $x \in \mathbb{R}^m$,

$$x^\top (A^\top A)x = (Ax)^\top Ax = \|Ax\|^2 \geq 0 . \quad (\text{A.1})$$

Therefore, by the spectral theorem, there exist $\lambda_1 \geq \dots \geq \lambda_m \geq 0$, and an orthonormal basis (v_1, \dots, v_m) of \mathbb{R}^m such that:

$$\forall i \in [m], \quad A^\top A v_i = \lambda_i v_i . \quad (\text{A.2})$$

Let $r = \text{rank } A$, which is also equal²⁴ to the rank of $A^\top A$. One thus has $\lambda_r \neq 0$ and $\lambda_{r+1} = 0$. Write $\sigma_i = \sqrt{\lambda_i}$ and let, for $j \in [r]$,

$$u_j = \frac{A v_j}{\sigma_j} . \quad (\text{A.3})$$

²⁴one short proof is to use $\text{Im } A^\top A \subset \text{Im } A^\top$, and show that both subspaces have the dimension by applying the rank theorem using $\text{Ker } A^\top A = \text{Ker } A$ (one inclusion is trivial, the other comes from [Equation \(A.1\)](#))

These vectors are orthonormal:

$$u_i^\top u_j = \frac{1}{\sigma_i \sigma_j} v_i^\top A^\top A v_j \quad (\text{A.4})$$

$$= \frac{\sigma_j}{\sigma_i} v_i^\top v_j \quad (\text{A.5})$$

$$= \delta_{ij} , \quad (\text{A.6})$$

where δ_{ij} is the Kronecker symbol (1 if $i = j$ and 0 else). We can thus complete (u_1, \dots, u_r) into a full basis of \mathbb{R}^n , (u_1, \dots, u_n) .

Define $U \in \mathbb{R}^{n \times n}$ whose rows are the u_i 's, and $V \in \mathbb{R}^{m \times n}$ whose columns are the v_i 's. Define $\Sigma \in \mathbb{R}^{n \times m}$ as a 0 matrix except its first r diagonal elements equal to $\sigma_1, \dots, \sigma_r$. Then Equation (A.3) rewrites, in matrix form:

$$U\Sigma = AV , \quad (\text{A.7})$$

that is, multiplying by V^\top on the right, and using $VV^\top = \text{Id}_m$

$$U\Sigma V^\top = A . \quad (\text{A.8})$$

This also takes the convenient form:

$$A = \sum_{i=1}^r \sigma_i u_i v_i^\top . \quad (\text{A.9})$$

Additional observations

- For $j \leq r$, one has $AA^\top u_j = AA^\top A v_j / \sigma_j = \sigma_j A v_j = \lambda_j u_j$ (where we have used $A^\top A v_j = \sigma_j^2 v_j$). This means that the u_j are eigenvectors of AA^\top (while the v_i 's are eigenvectors of $A^\top A$)²⁵.
- Notice that for $i \geq r$, u_i and v_i do not count in the SVD of A , hence U and V are sometimes reduced to their first r columns.
- The SVD of A^\top is easy to obtain from that of A .
- In Equation (A.7) we have used that $A v_j = \sigma_j u_j$ holds also for $j > r$ if by convention $\sigma_j = 0$ for $j > r$. Indeed for such a j , $v_j \in \text{Ker } A^\top A$ ($\lambda_j = 0$) and thus $v_j \in \text{Ker } A$ and both sides are 0.

B A coercive l.s.c. function without minimizer

This example has to be in infinite dimension. We consider $X = \ell^\infty$ and ϕ a Banach limit.

On X , define

$$f(x) = \sum_1^{+\infty} \frac{x_n}{n^2} + \max(0, \phi(-x)) + 4 \max(\|x\|, 0) \quad (\text{B.1})$$

²⁵one must check that for $i \geq r$ too. The associated eigenvalue is 0.

This function is continuous because by definition, a Banach limit is continuous. It is also coercive because the first two terms are lower bounded by $-\pi^2/6\|x\|$ and $-\|x\|$ respectively, so the 4 of the last term wins (actually the second term is lower bounded by 0).

Because the second term is positive, we can derive the following lower bound for f :

$$f(x) \geq -\frac{\pi^2}{6}\|x\| + 0 + 4\max(\|x\| - 1, 0). \quad (\text{B.2})$$

It turns out this is tight (it is the inf), as we can approach it arbitrarily close with sequences equal to 1 up to a certain rank, and 0 afterwards.

If $f(x) = -\pi^2/6$, we must have $\|x\| = 1$ and so $\sum_1^{+\infty} \frac{x_n}{n^2} = -\frac{\pi^2}{6}$, which implies $x_n = -1$ for all n , but then the Banach limit of $-x$ is 1, which leads to $f(x) = 1 - \frac{\pi^2}{6}$, a contradiction.