

Optimization for Data Science

Cristian Perez Jensen

July 7, 2024

Note that these are not the official lecture notes of the taught course, but only notes written by me. As such, it might contain mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

1	Risk minimization	1
1.1	Algorithms in data science	1
1.2	Empirical and expected risk	1
1.3	The map of learning	3
2	Theory of convex functions	4
2.1	Mathematical background	4
2.2	Convex sets	4
2.3	Convex functions	5
2.4	Minimizing convex functions	9
2.5	Convex programming	10
3	Gradient descent	13
3.1	Vanilla analysis	13
3.2	Lipschitz convex functions	14
3.3	Smooth functions	15
3.4	Smooth and strongly convex functions	17
4	Projected gradient descent	20
4.1	Lipschitz convex functions	20
4.2	Smooth functions	21
4.3	Smooth and strongly convex functions	23
5	Coordinate descent	24
5.1	Randomized coordinate descent	25
5.2	Importance sampling	26
5.3	Steepest coordinate descent	27
5.4	Greedy coordinate descent	29
6	Nonconvex functions	30
6.1	Trajectory analysis	32
7	The Frank-Wolfe algorithm	36
7.1	Linear minimization oracles	36
7.2	Duality gap	37
7.3	Convergence analysis	37
8	Newton's method	42
9	Quasi-Newton methods	45
10	Subgradient methods	48
10.1	Subgradient method	49
10.2	Strong convexity	51
11	Mirror descent	52
11.1	Norm and Bregman divergence	52
11.2	Mirror descent algorithm	53
12	Smoothing and proximal algorithms	56
12.1	Nesterov smoothing	56

12.2	<i>Moreau-Yosida smoothing</i>	57
12.3	<i>Proximal operators</i>	57
12.4	<i>Proximal point algorithm</i>	57
12.5	<i>Proximal gradient method</i>	58
13	<i>Stochastic optimization</i>	59
13.1	<i>Convergence analysis</i>	59
13.2	<i>Adaptive methods</i>	61
13.3	<i>Variance reduction</i>	62

List of symbols

\doteq	Equality by definition
\mathbb{R}	Set of real numbers
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbf{v} \in \mathbb{R}^n$	n -dimensional vector
$\mathbf{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
\mathbf{M}^\top	Transpose of matrix \mathbf{M}
\mathbf{M}^{-1}	Inverse of matrix \mathbf{M}
$\det(\mathbf{M})$	Determinant of \mathbf{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(x)$	Partial derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}^n$
$\nabla f(x) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$
$Jf(x) \in \mathbb{R}^{n \times m}$	Jacobian of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at point $x \in \mathbb{R}^n$
$\nabla^2 f(x) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$

1 Risk minimization

1.1 Algorithms in data science

In classical algorithm theory, an optimization problem solves a well-defined problem. For example, Kruskal's algorithm computes the minimum spanning tree of a graph. In data science, it is not as well-defined. The starting point is a learning problem, and the optimization typically happens on training data. However, even a perfect result may fail to solve the learning problem, which is a failure of the model in which the optimization algorithm was applied, rather than the optimization algorithm itself.

1.2 Empirical and expected risk

Typically, we have a data source \mathcal{X} , equipped with an unknown probability distribution. However, we do have access to independent samples $X_1, \dots, X_n \sim \mathcal{X}$. The goal is to "explain" \mathcal{X} through these samples. More specifically, we have a class \mathcal{H} of hypotheses (possible explanations of \mathcal{X}). The goal is then to select the hypothesis $H \in \mathcal{H}$ that best explains \mathcal{X} , which we measure by a *risk* (or *loss*) function $\ell : \mathcal{H} \times \mathcal{X} \rightarrow \mathbb{R}$.

The expected risk can be computed by

$$\ell(H) \doteq \mathbb{E}_{\mathcal{X}}[\ell(H, X)].$$

The best explanation is the *expected risk minimizer*,

$$H^* = \operatorname{argmin}_{H \in \mathcal{H}} \ell(H).$$

However, since we do not have access to the distribution over \mathcal{X} , we cannot compute $\ell(H)$ or H^* . Thus, we need to compute the *probably approximately correct* (PAC) hypothesis. This means that given $\delta, \epsilon > 0$, we want to produce, with probability $1 - \delta$, a hypothesis $\tilde{H} \in \mathcal{H}$ such that

$$\ell(\tilde{H}) \leq \inf_{H \in \mathcal{H}} \ell(H) + \epsilon,$$

meaning that with high probability, we approximately solve the expected risk minimization problem.

However, we can still not compute this, thus we must use our training data to compute the *empirical risk*,

$$\ell_n(H) = \frac{1}{n} \sum_{i=1}^n \ell(H, X_i).$$

This is a random variable, because it depends on the training data, which are all random variables, distributed according to the probability distribution of \mathcal{X} .

Lemma 1 (Weak law of large numbers). Let $H \in \mathcal{H}$ be a hypothesis. For any $\delta, \epsilon > 0$, there exists $n_0 \in \mathbb{N}$ such that for $n \geq n_0$,

$$|\ell_n(H) - \ell(H)| \leq \epsilon,$$

with probability at least $1 - \delta$.

Given $n \in \mathbb{N}$ and training data X_1, \dots, X_n from \mathcal{X} , we want to produce a hypothesis \tilde{H}_n such that

$$\ell_n(\tilde{H}_n) \leq \inf_{H \in \mathcal{H}} \ell_n(H) + \epsilon.$$

In an ideal world, \tilde{H}_n is also almost the best explanation for the data source \mathcal{X} ,

$$\ell(\tilde{H}_n) \leq \inf_{H \in \mathcal{H}} \ell(H) + \epsilon.$$

Remark. Note that the weak law of large numbers can only be applied to a *fixed* hypothesis, but not to the data-dependent hypothesis \tilde{H}_n . Thus, we are not always in an ideal world scenario.

A sufficient condition for an ideal world scenario is that the weak law of large numbers uniformly holds for all hypotheses with high probability. This leads us to the following theorem.

Theorem 2. Assume that for any $\delta, \epsilon > 0$, there exists $n_0 \in \mathbb{N}$ such that for $n \geq n_0$,

$$\sup_{H \in \mathcal{H}} |\ell_n(H) - \ell(H)| \leq \epsilon,$$

with probability at least $1 - \delta$. Then, for $n \geq n_0$, an approximate empirical risk minimizer \tilde{H}_n is PAC for expected risk minimization, meaning that it satisfies

$$\ell(\tilde{H}_n) \leq \inf_{H \in \mathcal{H}} \ell(H) + 3\epsilon,$$

with probability at least $1 - \delta$.

Proof.

$$\begin{aligned} \ell(\tilde{H}_n) &\leq \ell_n(\tilde{H}_n) + \epsilon \\ &\leq \inf_{H \in \mathcal{H}} \ell_n(H) + 2\epsilon \\ &\leq \inf_{H \in \mathcal{H}} \ell(H) + 3\epsilon, \end{aligned}$$

with probability at least $1 - \delta$. ■

It turns out that the assumption made by Theorem 2 holds in many relevant cases, but it is not always true.

Follows from $\sup_{H \in \mathcal{H}} |\ell_n(H) - \ell(H)| \leq \epsilon$.

\tilde{H}_n is an almost best explanation of the training data.

Follows from $\sup_{H \in \mathcal{H}} |\ell_n(H) - \ell(H)| \leq \epsilon$.

In this course, we will not learn how to pick the theory (\mathcal{H} and ℓ), but rather how to solve the optimization problems that arise in empirical risk minimization after the theory has been chosen.

1.3 The map of learning

The map of learning can be seen in Figure 1. It plots the empirical risk ($\ell_n(H_n)$, training data) against the expected risk ($\ell(H_n)$, estimated by a validation set). We can only be in the area denoted by “empirical risk minimization”, because

$$\begin{aligned}\ell_n(\tilde{H}_n) &\leq \inf_{H \in \mathcal{H}} \ell_n + \epsilon \\ &\leq \ell_n(\tilde{H}) + \epsilon \\ &\leq \ell(\tilde{H}) + 2\epsilon \\ &\leq \ell(\tilde{H}_n) + \epsilon\epsilon.\end{aligned}$$

A model is overfit when we have low empirical risk, while having high expected risk. This means that the explanation quality on the data source is much worse than on the training data. The main reason for this is that the theory (\mathcal{H} and ℓ) is so complex that it can explain almost anything.

A model is underfit when we have high empirical risk, while having high expected risk. This means that we neither explain the training data nor the data source. The main reason for this is that the theory is too simple to capture the nature of the data.

The model is learning when we have low empirical risk and low expected risk. This means that the training was successful. Generalization occurs when the expected risk is close to the empirical risk. Note that this does not mean that the explanation is good, since any “blind explanation” will generalize well due to the weak law of large numbers. Ideally, we want generalization and learning.



Figure 1. The map of learning.

2 Theory of convex functions

2.1 Mathematical background

Theorem 3 (Cauchy-Schwarz inequality). Let $u, v \in \mathbb{R}^d$, then

$$|u^\top v| \leq \|u\| \|v\|.$$

For non-zero vectors, this is equivalent to

$$-1 \leq \frac{u^\top v}{\|u\| \|v\|} \leq 1.$$

Definition 4 (Spectral norm). Let A be an $m \times d$ matrix, then

$$\|A\| \doteq \max_{v \in \mathbb{R}^d, v \neq 0} \frac{\|Av\|}{\|v\|} = \max_{\|v\|=1} \|Av\|.$$

Intuitively, this means that $\|A\|$ is the largest factor by which a vector can be stretched in length under the mapping $v \mapsto Av$.

Theorem 5 (Mean value theorem). Let $a < b$ be real numbers, and let $h : [a, b] \rightarrow \mathbb{R}$ be a continuous function that is differentiable on (a, b) . Then, there exists $c \in (a, b)$ such that

$$h'(c) = \frac{h(b) - h(a)}{b - a}.$$

Geometrically, this means that between a and b , there is a tangent to the graph of h that has the same slope; see Figure 2.

Theorem 6 (Fundamental theorem of calculus). Let $a < b$ be real numbers, and let $h : \text{dom}(h) \rightarrow \mathbb{R}$ be a differentiable function on an open domain $\text{dom}(h) \supset [a, b]$, and such that h' is continuous on $[a, b]$. Then,

$$h(b) - h(a) = \int_a^b h'(t) dt.$$



Figure 2. Illustration of the mean value theorem.

2.2 Convex sets

Definition 7 (Convex set). A set $\mathcal{C} \subseteq \mathbb{R}^d$ is convex if the line segment between any two points of \mathcal{C} lies in \mathcal{C} . I.e., if for any $x, y \in \mathcal{C}$ and any λ with $0 \leq \lambda \leq 1$, we have

$$\lambda x + (1 - \lambda)y \in \mathcal{C}.$$

Observation. Let $\mathcal{C}_i, i \in I$ be convex sets, where I is a (possibly infinite)



Figure 3. Example of a convex set in \mathbb{R}^2 .



index set. Then,

$$\mathcal{C} = \bigcap_{i \in I} \mathcal{C}_i,$$

is a convex set.

2.3 Convex functions

Definition 8 (Convexity). A function $f : \text{dom}(f) \rightarrow \mathbb{R}$ is convex if (i) $\text{dom}(f)$ is convex and (ii) for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and all $\lambda \in [0, 1]$, we have

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

Geometrically, this condition means that the line segment connecting the two points $(\mathbf{x}, f(\mathbf{x}))$ and $(\mathbf{y}, f(\mathbf{y})) \in \mathbb{R}^{d+1}$ lies pointwise above the graph of f ; see Figure 5.

Definition 9 (Epigraph). The epigraph of a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

$$\text{epi}(f) \doteq \left\{ (\mathbf{x}, \alpha) \in \mathbb{R}^{d+1} \mid \mathbf{x} \in \text{dom}(f), \alpha \geq f(\mathbf{x}) \right\}.$$

Observation. f is a convex function if and only if $\text{epi}(f)$ is a convex set.

Lemma 10 (Jensen's inequality). Let f be convex, $\mathbf{x}_1, \dots, \mathbf{x}_m \in \text{dom}(f)$, $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$ such that $\sum_{i=1}^m \lambda_i = 1$, then

$$f\left(\sum_{i=1}^m \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^m \lambda_i f(\mathbf{x}_i).$$

Proof. We will prove Jensen's inequality by induction. The base case ($m = 2$) is true by definition of convexity. Let Jensen's inequality hold for $m = k$. Consider $m = k + 1$, then

$$\begin{aligned} f\left(\sum_{i=1}^{k+1} \lambda_i \mathbf{x}_i\right) &= f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i + \lambda_{k+1} \mathbf{x}_{k+1}\right) \\ &= f\left((1 - \lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} \mathbf{x}_i + \lambda_{k+1} \mathbf{x}_{k+1}\right) \\ &\leq (1 - \lambda_{k+1}) f\left(\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} \mathbf{x}_i\right) + \lambda_{k+1} f(\mathbf{x}_{k+1}) \\ &\leq (1 - \lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} f(\mathbf{x}_i) + \lambda_{k+1} f(\mathbf{x}_{k+1}) \\ &= \sum_{i=1}^{k+1} \lambda_i f(\mathbf{x}_i). \end{aligned}$$

Since it holds for the base case and the induction step, Jensen's inequality must hold for all m . ■

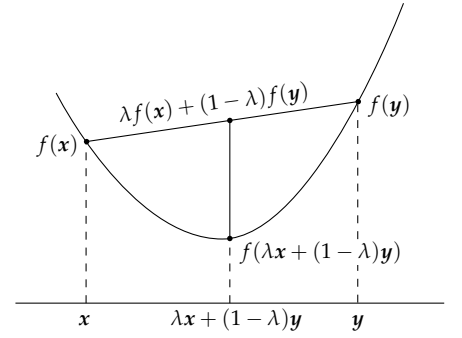


Figure 5. Illustration of the classic definition of convexity.



Figure 6. Epigraph of a convex function.

By definition of convexity.

Induction assumption, $\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} = 1$.

Lemma 11. Let f be convex and suppose that $\text{dom}(f) \subseteq \mathbb{R}^d$ is open, then f is continuous.

Definition 12 (Differentiable functions). Let $f : \text{dom}(f) \rightarrow \mathbb{R}^m$ where $\text{dom}(f) \subseteq \mathbb{R}^d$ is open. f is called differentiable at $x \in \text{dom}(f)$ if there exists an $m \times d$ matrix A and an error function $r : \mathbb{R}^d \rightarrow \mathbb{R}^m$ defined around $0 \in \mathbb{R}^d$ such that for all y in some neighborhood of x ,

$$f(y) = f(x) + A(y - x) + r(y - x),$$

where $\lim_{v \rightarrow 0} \|r(v)\|/\|v\| = 0$ (error function r is sublinear around 0). A is unique and called the Jacobian matrix of f at x .

Lemma 13 (First-order convexity). Suppose that $\text{dom}(f)$ is open and that f is differentiable. In particular, the gradient

$$\nabla f(x) \doteq \left[\frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_d} \right]$$

exists at every point $x \in \text{dom}(f)$.

Then, f is convex if and only if (i) $\text{dom}(f)$ is convex and (ii)

$$\forall x, y \in \text{dom}(f) : f(y) \geq f(x) + \nabla f(x)^\top (y - x).$$

Geometrically, this means that the graph is above all tangent hyperplanes; see Figure 8.

Proof. We will first prove that convexity implies that the first-order definition holds. Then, we will prove that the first-order definition implies convexity, making the definitions equivalent for differentiable functions f .

\Rightarrow : Suppose f is convex. Then, for all $t \in (0, 1)$,

$$\begin{aligned} f((1-t)x + ty) &\leq (1-t)f(x) + tf(y) \\ \Leftrightarrow f(x + t(y-x)) &\leq f(x) + t(f(y) - f(x)) \\ \Leftrightarrow f(y) &\geq f(x) + \frac{f(x + t(y-x)) - f(x)}{t} \\ &= f(x) + \frac{\nabla f(x)^\top t(y-x) + r(t(y-x))}{t} \\ &= f(x) + \nabla f(x)^\top (y-x) + \underbrace{\frac{r(t(y-x))}{t}}_{\rightarrow 0 \text{ for } t \rightarrow 0}. \end{aligned}$$

\Leftarrow : Define $z \doteq \lambda x + (1-\lambda)y \in \text{dom}(f)$ for $\lambda \in [0, 1]$ by convexity of

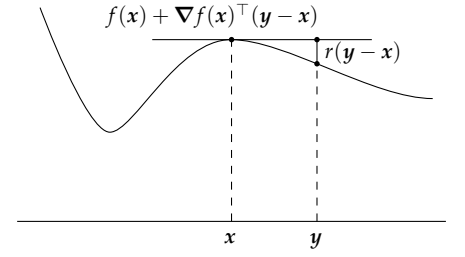


Figure 7. Graph of the affine function $f(x) + \nabla f(x)^\top (y - x)$ is a tangent hyperplane to the graph of f at $(x, f(x))$.

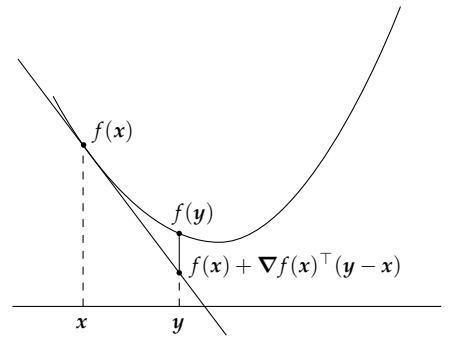


Figure 8. Illustration of the first-order characterization of convexity (Lemma 13).

Definition of convexity

$\text{dom}(f)$. Then, we have the following inequalities,

$$\begin{aligned} f(x) &\geq f(z) + \nabla f(z)^\top (x - z) \\ \lambda f(x) &\geq \lambda \left(f(z) + \nabla f(z)^\top (x - z) \right) \\ f(y) &\geq f(z) + \nabla f(z)^\top (y - z) \\ (1 - \lambda)f(y) &\geq (1 - \lambda) \left(f(z) + \nabla f(z)^\top (y - z) \right). \end{aligned}$$

This implies the following,

$$\begin{aligned} \lambda f(x) + (1 - \lambda)f(y) &\geq \lambda \left(f(z) + \nabla f(z)^\top (x - z) \right) + (1 - \lambda) \left(f(z) + \nabla f(z)^\top (y - z) \right) \\ &= f(z) + \lambda \nabla f(z)^\top (x - z) + (1 - \lambda) \nabla f(z)^\top (y - z) \\ &= f(z) + \nabla f(z)^\top (\lambda(x - z) + (1 - \lambda)(y - z)) \\ &= f(z) + \nabla f(z)^\top (\lambda x + (1 - \lambda)y - (\lambda z + (1 - \lambda)z)) \\ &= f(z) + \nabla f(z)^\top (z - z) \\ &= f(\lambda x + (1 - \lambda)y). \end{aligned}$$

■

Lemma 14 (First-order convexity alternative). Suppose that $\text{dom}(f)$ is open and that f is differentiable. Then, f is convex if and only if (i) $\text{dom}(f)$ is convex and (ii)

$$\forall x, y \in \text{dom}(f) : (\nabla f(y) - \nabla f(x))^\top (y - x) \geq 0.$$

Geometrically, this means that the gradient is monotonic.

Proof. We will first prove that it holds from left to right, and then from right to left.

\Rightarrow : If f is convex, the first-order characterization of convexity (Lemma 13) yields the following,

$$\begin{aligned} f(x) &\geq f(y) + \nabla f(y)^\top (x - y) \\ f(y) &\geq f(x) + \nabla f(x)^\top (y - x), \end{aligned}$$

for all $x, y \in \text{dom}(f)$. Adding up these two inequalities yields

$$\begin{aligned} f(x) + f(y) &\geq f(y) + \nabla f(y)^\top (x - y) + f(x) + \nabla f(x)^\top (y - x) \\ \Leftrightarrow 0 &\geq \nabla f(y)^\top (x - y) + \nabla f(x)^\top (y - x) \\ &= (\nabla f(y) - \nabla f(x))^\top (x - y) \\ \Leftrightarrow 0 &\leq (\nabla f(y) - \nabla f(x))^\top (y - x). \end{aligned}$$

\Leftarrow : Define $z \doteq (1 - t)x + ty \in \text{dom}(f)$ for $x, y \in \text{dom}(f), t \in (0, 1)$ by convexity of $\text{dom}(f)$. Observe that $z = x + t(y - x)$. Then, we have the

following inequality, according to the monotonicity of the gradient,

$$\begin{aligned} (\nabla f(z) - \nabla f(x))^\top (z - x) &\geq 0 \\ (\nabla f(x + t(\mathbf{y} - x)) - \nabla f(x))^\top (x + t(\mathbf{y} - x) - x) &\geq 0 \\ (\nabla f(x + t(\mathbf{y} - x)) - \nabla f(x))^\top (\mathbf{y} - x) &\geq 0. \end{aligned} \quad \text{Divide by } t$$

Let $h(t) = f(x + t(\mathbf{y} - x))$, then $h'(t) = \nabla f(x + t(\mathbf{y} - x))^\top (\mathbf{y} - x)$. Hence, we can rewrite the inequality as the following,

$$h'(t) \geq \nabla f(x)^\top (\mathbf{y} - x).$$

By the mean value theorem, there exists $c \in (0, 1)$ such that $h'(c) = h(1) - h(0)$. I.e.,

$$h'(c) = f(\mathbf{y}) - f(x).$$

Thus, we can rewrite the inequality to the following,

$$\begin{aligned} f(\mathbf{y}) &= f(x) + h'(c) \\ &\geq f(x) + \nabla f(x)^\top (\mathbf{y} - x), \end{aligned}$$

which is the first-order characterization of convexity (Lemma 13). ■

Lemma 15 (Second-order convexity). Suppose that $\text{dom}(f)$ is open and that f is twice differentiable. In particular, the Hessian

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_d \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_d^2} \end{bmatrix}$$

exists at every point $x \in \text{dom}(f)$ and is symmetric.

Then, f is convex if and only if (i) $\text{dom}(f)$ is convex and (ii) for all $x \in \text{dom}(f)$, we have

$$\nabla^2 f(x) \succeq 0.$$

I.e., $\nabla^2 f(x)$ is positive semidefinite (M is positive semidefinite if $x^\top M x \geq 0$ for all x and $x^\top M x > 0$ for all $x \neq \mathbf{0}$).

Geometrically, this means that f has non-negative curvature everywhere. I.e., the growth rate should be growing.

Observation (Operations that preserve convexity). Let f_1, \dots, f_m be convex functions, $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$, then

$$f \doteq \max_{i=1}^m f_i,$$

and

$$f \doteq \sum_{i=1}^m \lambda_i f_i.$$

are convex on $\text{dom}(f) \doteq \bigcap_{i=1}^m \text{dom}(f_i)$.

Let f be a convex function with $\text{dom}(f) \subseteq \mathbb{R}^d$, $g: \mathbb{R}^m \rightarrow \mathbb{R}^d$ an affine function, i.e. $g(x) = Ax + b$ for some $A \in \mathbb{R}^{d \times m}$, $b \in \mathbb{R}^d$. Then, the function $f \circ g$ is convex on $\text{dom}(f \circ g) \doteq \{x \in \mathbb{R}^m \mid g(x) \in \text{dom}(f)\}$.

2.4 Minimizing convex functions

Definition 16 (Local minimum). A local minimum of $f: \text{dom}(f) \rightarrow \mathbb{R}$ is a point x such that there exists $\epsilon > 0$ with

$$f(x) \leq f(y) \quad \forall y \in \text{dom}(f) \text{ s.t. } \|y - x\| < \epsilon.$$

Remark. This definition of local minima means that in some small neighborhood around x , x is the smallest point.

Lemma 17. Let x^* be a local minimum of a convex function $f: \text{dom}(f) \rightarrow \mathbb{R}$. Then, x^* is a global minimum, meaning that $f(x^*) \leq f(y) \quad \forall y \in \text{dom}(f)$.

Proof. Proof by contradiction. Suppose there exists $y \in \text{dom}(f)$ such that $f(y) < f(x^*)$. Define $y' \doteq \lambda x^* + (1 - \lambda)y$ for $\lambda \in (0, 1)$. From convexity, we get that $f(y') < f(x^*)$, because $f(y) < f(x^*)$. If we choose λ so close to 1 such that $\|y' - x^*\| < \epsilon$ gives the inequality $f(x^*) \leq f(y')$. This yields a contradiction, thus there cannot exist a $y \in \text{dom}(f)$ such that $f(y) < f(x^*)$, meaning that $f(x^*) \leq f(y)$ for all $y \in \text{dom}(f)$. ■

Lemma 18. Suppose that f is convex and differentiable over an open domain $\text{dom}(f)$. Let $x \in \text{dom}(f)$. If $\nabla f(x) = 0$ (critical point), then x is a global minimum.

Proof. Suppose that $\nabla f(x) = 0$. According to the first-order characterization of convexity, we have

$$f(y) \geq f(x) + \underbrace{\nabla f(x)^\top (y - x)}_{=0} = f(x)$$

for all $y \in \text{dom}(f)$, so x is a global minimum. ■

Definition 19 (Strict convexity). A function $f: \text{dom}(f) \rightarrow \mathbb{R}$ is strictly convex if (i) $\text{dom}(f)$ is convex and (ii) for all $x \neq y \in \text{dom}(f)$ and all $\lambda \in (0, 1)$, we have

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y).$$

An example of a strictly convex function can be found in Figure 5, while an example of a function that is not strictly convex can be found in Figure 10.



Figure 9. The maximum operator over m convex functions is a convex function. As can be seen, the epigraph of f is convex.

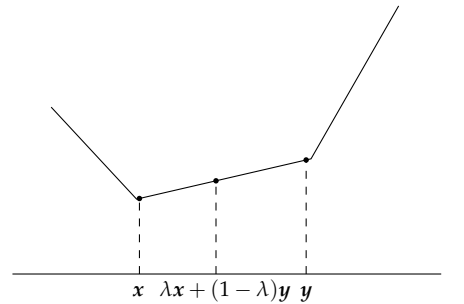


Figure 10. A non-strictly convex function.

Lemma 20. Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be strictly convex. Then, f has at most one global minimum.

But, how do we know a minimizer exists? To be able to find this out, we must use the Weierstrass theorem.

Definition 21 (Sublevel set). $f : \mathbb{R}^d \rightarrow \mathbb{R}, \alpha \in \mathbb{R}$. The set $f^{\leq \alpha} \doteq \{x \in \mathbb{R}^d \mid f(x) \leq \alpha\}$ is the α -sublevel set of f .

Theorem 22 (Weierstrass). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous function, and suppose there is a non-empty and bounded sublevel set $f^{\leq \alpha}$. Then, f has a global minimum.

Proof. We know that f attains a minimum over the closed and bounded set $f^{\leq \alpha}$ at some x^* , because f is continuous. This x^* is also a global minimum as it has value $f(x^*) \leq \alpha$, while any $x \notin f^{\leq \alpha}$ has value $f(x) > \alpha \geq f(x^*)$. ■

Since convex functions are continuous by Lemma 11, the Weierstrass theorem also applies to convex functions. Thus, to figure out whether a convex function has a minimizer, we need to find any non-empty and bounded sublevel set of this function.

2.5 Convex programming

In standard form, optimization problems look like the following,

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, \dots, m \\ & && h_j(x) = 0, \quad j = 1, \dots, p. \end{aligned}$$

The problem domain is then $\mathcal{D} = \bigcap_{i=1}^m \text{dom}(f_i) \cap \bigcap_{j=1}^p \text{dom}(h_j)$.

In a convex program, all f_i are convex functions, and all h_j are affine functions with domain \mathbb{R}^d .

Definition 23 (Lagrange dual function). Given an optimization problem in standard form, its Lagrangian is the function $L : \mathcal{D} \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ given by

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \nu_j h_j(x).$$

The λ_i, ν_j are called Lagrange multipliers. The Lagrange dual function is the function $g : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R} \cup \{-\infty\}$ defined by

$$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu).$$



Figure 11. Illustration of a sublevel set of a non-convex function.

Only the (λ, ν) pairs with $g(\lambda, \nu) > -\infty$ are interesting.

Lemma 24 (Weak Lagrange duality). Let \mathbf{x} be a feasible solution ($f_i(\mathbf{x}) \leq 0$ for $i = 1, \dots, m$ and $h_j(\mathbf{x}) = 0$ for $j = 1, \dots, p$). Let g be the Lagrange dual function and $\lambda \in \mathbb{R}^m, \nu \in \mathbb{R}^p$ such that $\lambda \geq \mathbf{0}$. Then

$$g(\lambda, \nu) \leq f_0(\mathbf{x}).$$

Proof.

$$g(\lambda, \nu) \leq L(\mathbf{x}, \lambda, \nu) = f_0(\mathbf{x}) + \underbrace{\sum_{i=1}^m \lambda_i f_i(\mathbf{x})}_{\leq 0} + \underbrace{\sum_{j=1}^p \nu_j h_j(\mathbf{x})}_{=0} \leq f_0(\mathbf{x}).$$

■

However, we want to know what the best lower bound is that we can get in this way. For this, we must choose $\lambda \geq \mathbf{0}, \nu$ such that $g(\lambda, \nu)$ is maximized. This can be phrased as another optimization problem, called the Lagrange dual,

$$\begin{aligned} & \text{maximize} && g(\lambda, \nu) \\ & \text{subject to} && \lambda \geq \mathbf{0}. \end{aligned}$$

Observation. The Lagrange dual is a convex program, even if the primal is not.

Theorem 25. Suppose that a convex program has a feasible solution $\tilde{\mathbf{x}}$ that in addition satisfies $f_i(\tilde{\mathbf{x}}) < 0, i = 1, \dots, m$ (Slater point). Then, the infimum value of the primal equals the supremum value of its Lagrange dual. Moreover, if the value is finite, it is attained by a feasible solution of the dual,

$$\inf f_0(\mathbf{x}) = \max g(\lambda, \nu).$$

A case of particular interest is that strong duality holds and the joint value is attained in both the primal and dual problem.¹ If all f_i and h_j are differentiable, then the Karush-Kuhn-Tucker (KKT) conditions provide necessary and, under convexity, also sufficient conditions for this case to occur.

¹ This would mean that $\min f_0(\mathbf{x}) = \max g(\lambda, \nu)$.

Definition 26 (Zero duality gap). Let $\tilde{\mathbf{x}}$ be feasible for the primal and $(\tilde{\lambda}, \tilde{\nu})$ feasible for the Lagrange dual. The primal and dual solutions $\tilde{\mathbf{x}}$ and $(\tilde{\lambda}, \tilde{\nu})$ are said to have zero duality gap if $f_0(\tilde{\mathbf{x}}) = g(\tilde{\lambda}, \tilde{\nu})$.

If the primal and dual have zero duality gap, then $\min f_0(\mathbf{x}) = \max g(\lambda, \nu)$.

The consequence of zero duality gap is the *master equation*,

$$\begin{aligned}
 f_0(\tilde{\mathbf{x}}) &= g(\tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{v}}) \\
 &= \inf_{\mathbf{x} \in \mathcal{D}} \left(f_0(\mathbf{x}) + \sum_{i=1}^m \tilde{\lambda}_i f_i(\mathbf{x}) + \sum_{j=1}^p \tilde{v}_j h_j(\mathbf{x}) \right) \\
 &\leq f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \underbrace{\tilde{\lambda}_i f_i(\tilde{\mathbf{x}})}_{\leq 0} + \sum_{j=1}^p \underbrace{\tilde{v}_j h_j(\tilde{\mathbf{x}})}_{=0} \\
 &\leq f_0(\tilde{\mathbf{x}}),
 \end{aligned}$$

which means that all inequalities turn into equalities. Thus,

$$\tilde{\lambda}_i f_i(\tilde{\mathbf{x}}) = 0, \quad i = 1, \dots, m,$$

which is called complementary slackness, because if $\tilde{\lambda}_i \neq 0$, then $f_i(\tilde{\mathbf{x}}) = 0$, and vice versa. Furthermore, if all f_i and h_j are differentiable, then

$$\nabla f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \tilde{\lambda}_i \nabla f_i(\tilde{\mathbf{x}}) + \sum_{j=1}^p \tilde{v}_j \nabla h_j(\tilde{\mathbf{x}}) = \mathbf{0},$$

which is called the vanishing Lagrangian gradient condition.

In summary, we get the following result.

Theorem 27 (Karush-Kuhn-Tucker necessary conditions). Let $\tilde{\mathbf{x}}$ and $(\tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{v}})$ be feasible solutions of the primal optimization problem and its Lagrange dual, respectively, with zero duality gap. If all f_i and h_j are differentiable, then

$$\begin{aligned}
 \tilde{\lambda}_i f_i(\tilde{\mathbf{x}}) &= 0, \quad i = 1, \dots, m \\
 \nabla f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \tilde{\lambda}_i \nabla f_i(\tilde{\mathbf{x}}) + \sum_{j=1}^p \tilde{v}_j \nabla h_j(\tilde{\mathbf{x}}) &= \mathbf{0}.
 \end{aligned}$$

Complementary slackness

Vanishing Lagrangian gradient

If we have a convex program, then the KKT conditions are sufficient for zero duality gap. The motivation behind this is that they may be easier to solve for than solving the primal optimization problem. However, we cannot always count on the KKT conditions to be solvable, because they are only guaranteed if there are primal and dual solutions of zero duality gap. But, if the primal has a Slater point², then the KKT conditions are equivalent to the existence of a primal minimizer.

² A point $\tilde{\mathbf{x}}$ such that $f_i(\tilde{\mathbf{x}}) < 0$ for all $i = 1, \dots, m$.

3 Gradient descent

Gradient descent is an optimization algorithm that aims to find the global minimizer. We assume that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, differentiable, and has a global minimizer \mathbf{x}^* . Then, the goal of gradient descent is, for $\epsilon > 0$, to find a $\mathbf{x} \in \mathbb{R}^d$ such that

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \epsilon.$$

It works by first choosing a $\mathbf{x}_0 \in \mathbb{R}^d$ and then iteratively updating this value by

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t),$$

for timesteps $t = 0, 1, \dots$ and stepsize $\gamma > 0$.

3.1 Vanilla analysis

We want to be able to bound $f(\mathbf{x}_t) - f(\mathbf{x}^*)$, which tells us how far off we are from the minimum at timestep t . For this, we can use the first-order characterization of convexity,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}^* - \mathbf{x}_t) \iff f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*).$$

Using the gradient descent update rule, we get

$$\nabla f(\mathbf{x}_t) = \frac{\mathbf{x}_t - \mathbf{x}_{t+1}}{\gamma},$$

which gives the following bound

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{1}{2\gamma} \left(\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right) \\ &= \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \left(\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right). \end{aligned}$$

$2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$ (cosine theorem).

$$\mathbf{x}_t - \mathbf{x}_{t+1} = \gamma \nabla f(\mathbf{x}_t).$$

Summing this up over the first T iterations, we get an upper bound on the summed error,

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \sum_{t=0}^{T-1} \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \left(\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2 \right) && \text{Telescoping sum.} \\ &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. && \|\mathbf{x}_T - \mathbf{x}^*\|^2 \geq 0. \end{aligned}$$

We can also use this to get a bound on the average error,

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{1}{T} \left(\frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \right).$$



Figure 12. Gradient descent updates.

Take a small step into the direction of the negative gradient to move toward the minimum.

3.2 Lipschitz convex functions

Theorem 28. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable with a global minimum \mathbf{x}^* . Furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ and $\|\nabla f(\mathbf{x})\| \leq B$ for all \mathbf{x} . Choosing the stepsize,

$$\gamma \doteq \frac{R}{B\sqrt{T}},$$

gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{RB}{\sqrt{T}}.$$

Proof. We can plug our bounds into the vanilla analysis,

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\leq \frac{\gamma}{2} B^2 T + \frac{1}{2\gamma} R^2 \\ &\doteq q(\gamma). \end{aligned}$$

Now, we want to choose γ such that the bound $q(\gamma)$ is minimized. We can compute the minimum by computing the derivative,

$$q'(\gamma) = \frac{1}{2} B^2 T - \frac{R^2}{2\gamma^2},$$

and solving for $q'(\gamma) = 0$, which yields

$$\gamma = \frac{R}{B\sqrt{T}}.$$

Then, we can compute the bound by

$$q\left(\frac{R}{B\sqrt{T}}\right) = RB\sqrt{T}.$$

Dividing by T yields the result. ■

We want to find out how many iterations we would need to ensure that the average error is bounded by ϵ . We can use Theorem 28 to compute a lower bound on the number of iterations,

$$\frac{RB}{\sqrt{T}} \leq \epsilon \implies T \geq \frac{R^2 B^2}{\epsilon^2}.$$

So, the amount of iterations until convergence is of the order $\mathcal{O}(1/\epsilon^2)$. This means that we need at most $10000 \cdot R^2 B^2$ iterations to achieve an error of 0.01.

A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is Lipschitz continuous if there exists an M , such that $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^m$:

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq M \|\mathbf{x} - \mathbf{y}\|.$$

This holds if and only if the gradient is bounded.

3.3 Smooth functions

Definition 29 (Smoothness). Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be differentiable, $\mathcal{X} \subseteq \text{dom}(f)$ convex, $L \in \mathbb{R}_+$. f is smooth (with parameter L) over \mathcal{X} if $\forall x, y \in \mathcal{X}$:

$$\underbrace{f(y) \leq f(x) + \nabla f(x)^\top (y - x)}_{\text{first-order convexity}} + \frac{L}{2} \|x - y\|^2.$$

Remark. This definition does not require convexity.

Geometrically, this definition of smoothness means that the graph f is below a not too steep tangent paraboloid at $(x, f(x))$; see Figure 13. Furthermore, we can easily check if f is smooth with parameter L if the following function is convex,

$$g(x) \doteq \frac{L}{2} x^\top x - f(x),$$

over $\text{dom}(g) \doteq \text{dom}(f)$.

Observation (Operations that preserve smoothness). Let f_1, \dots, f_m be smooth with parameters L_1, \dots, L_m and let $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$. Then, the function $f \doteq \sum_{i=1}^m \lambda_i f_i$ is smooth with parameter $\sum_{i=1}^m \lambda_i L_i$.

Let f be smooth with parameter L and let $g(x) = Ax + b$. Then, the function $f \circ g$ is smooth with parameter $L\|A\|^2$, where $\|A\|$ is the spectral norm of A .

Lemma 30. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable, then f is smooth with parameter L if and only if $\forall x, y \in \mathbb{R}^d$:

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

I.e., ∇f is Lipschitz continuous.

Lemma 31 (Sufficient decrease). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and smooth with parameter L . With stepsize

$$\gamma \doteq \frac{1}{L},$$

gradient descent satisfies

$$f(x_{t+1}) \leq f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|^2.$$

Sufficient decrease implies that, every step, we get closer to the minimum.

“Not too curved.”

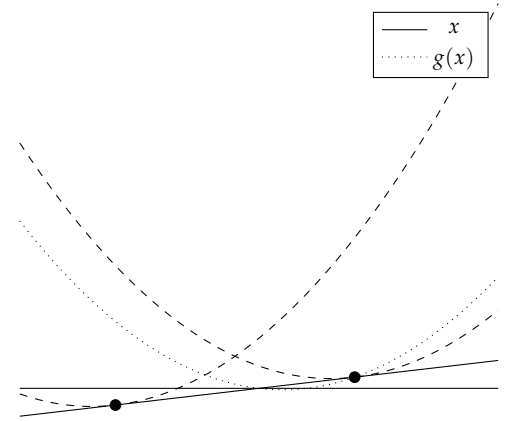


Figure 13. Plot of $f(x) = x$ with the tangent paraboloids at $x = -8, 5$, showing smoothness with parameter $L = 1$. Furthermore, it shows that $g(x) = \frac{1}{2}x^2 - f(x)$ is convex.

Proof.

$$\begin{aligned}
 f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\
 &= f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top \left(-\frac{\nabla f(\mathbf{x}_t)}{L} \right) + \frac{L}{2} \left\| \frac{\nabla f(\mathbf{x}_t)}{L} \right\|^2 \\
 &= f(\mathbf{x}_t) - \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\
 &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.
 \end{aligned}$$

Smoothness.

$$\text{GD: } \mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t).$$

Theorem 32. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable with a global minimum \mathbf{x}^* . Furthermore, suppose that f is smooth with parameter L . Choosing stepsize

$$\gamma \doteq \frac{1}{L},$$

gradient descent yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

Proof. Due to sufficient decrease, we can bound the sum of squared gradients (which will be useful when bounding the vanilla analysis),

$$\begin{aligned}
 \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 &\leq \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) \\
 &= f(\mathbf{x}_0) - f(\mathbf{x}_T).
 \end{aligned}$$

$$\text{SD: } f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

Telescoping sum.

Using the vanilla analysis, we can derive a bound on the average error,

$$\begin{aligned}
 \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\
 &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\
 \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\
 \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.
 \end{aligned}$$

$$\gamma \doteq \frac{1}{L}.$$

From sufficient decrease, we know that the last iterate must be the best, thus

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Using this result, we can compute a bound on T to achieve an error smaller than ϵ ,

$$f(x_T) - f(x^*) \leq \frac{L}{2T} \|x_0 - x^*\|^2 \leq \frac{R^2 L}{2T} \leq \epsilon.$$

The error then becomes,

$$T \geq \frac{R^2 L}{2\epsilon},$$

which is on the order of $\mathcal{O}(1/\epsilon)$. This means that we need at most $50 \cdot R^2 L$ iterations for an error of 0.01 (as opposed to $10000 \cdot R^2 B^2$ in the Lipschitz case).

3.4 Smooth and strongly convex functions

Definition 33 (Strong convexity). Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be a convex and differentiable function, $\mathcal{X} \in \text{dom}(f)$ convex and $\mu \in \mathbb{R}_+, \mu > 0$. f is called *strongly convex* with parameter μ over \mathcal{X} if $\forall x, y \in \mathcal{X}$:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|x - y\|^2.$$

Geometrically, this means that, for any x , the graph of f is above a not too flat tangent paraboloid at $(x, f(x))$; see Figure 14. Furthermore, we can easily check if f is strongly convex if and only if the following function is convex,

$$g(x) \doteq f(x) - \frac{\mu}{2} x^\top x.$$

By assuming that a function f is smooth and strongly convex, we can use a stronger lower bound to derive a bound on the error from the vanilla analysis.

Theorem 34. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable with a global minimum x^* . Furthermore, suppose that f is smooth with parameter L and strongly convex with parameter $\mu > 0$. Choosing $\gamma \doteq 1/L$, gradient descent with arbitrary x_0 satisfies the following two properties,

1. Squared distances to x^* are geometrically decreasing,

$$\|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|x_t - x^*\|^2;$$

2. The absolute error after T iterations is exponentially small in T ,

$$f(x_T) - f(x^*) \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|x_0 - x^*\|^2.$$



Figure 14. Plot of $f(x) = x^2$ with the tangent paraboloids at $x = -6, 4$, showing strong convexity with parameter $\mu = 1$. Furthermore, it shows $g(x) = f(x) - \frac{1}{2}x^2$, which is convex.

Proof of 1. Using the vanilla analysis, we know the following,

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2). \end{aligned}$$

But, using strong convexity, we have a stronger lower bound on the left hand side,

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 && \text{Strong convexity.} \\ &= \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &\quad - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2. \end{aligned}$$

This can be rewritten to the following bound,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \underbrace{2\gamma(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2}_{\text{"noise"}} + (1 - \mu\gamma) \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

Now we need to show that the noise is non-positive,

$$\begin{aligned} 2\gamma(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 &= \frac{2}{L} (f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \gamma \doteq \frac{1}{L}. \\ &\leq \frac{2}{L} (f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \text{Sufficient decrease.} \\ &\leq -\frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \text{SD: } f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2. \\ &= 0. \end{aligned}$$

Hence, the noise is non-positive, and we get the following,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

■

Proof of 2. From (1), we know the following,

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Using smoothness, we can derive a bound on the final iterate error,

$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_T - \mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}_T - \mathbf{x}^*\|^2 \\ &= \frac{L}{2} \|\mathbf{x}_T - \mathbf{x}^*\|^2 && \nabla f(\mathbf{x}^*) = 0, \text{ because it is a stationary point.} \\ &\leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2. && \text{Using (1).} \end{aligned}$$

■

From there, we can derive a lower bound on the number of iterations T to get an error of at most ϵ ,

$$T \geq \frac{L}{\mu} \log\left(\frac{R^2 L}{2\epsilon}\right),$$

This means that we need $\frac{L}{\mu} \log(50 \cdot R^2 L)$ iterations for an error of at most 0.01, as opposed to $50 \cdot R^2 L$ in the smooth case. This bound only depends linearly on $\frac{\mu}{L}$, which might be very high.

4 Projected gradient descent

In constrained optimization, we want to minimize $f(x)$, subject to $x \in \mathcal{X} \subseteq \mathbb{R}^d$. We will assume that $\mathcal{X} \subseteq \mathbb{R}^d$ is a closed convex set. The idea of projected gradient descent is to project onto \mathcal{X} after every step,

$$\Pi_{\mathcal{X}}(\mathbf{y}) \doteq \operatorname{argmin}_{x \in \mathcal{X}} \|\mathbf{x} - \mathbf{y}\|^2.$$

The algorithm thus alternates between the following two steps,

$$\begin{aligned} \mathbf{y}_{t+1} &\doteq \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &\doteq \Pi_{\mathcal{X}}(\mathbf{y}_{t+1}). \end{aligned}$$

Observation. Let $\mathcal{X} \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in \mathcal{X}$, $\mathbf{y} \in \mathbb{R}^d$, then

1. $(\mathbf{x} - \Pi_{\mathcal{X}}(\mathbf{y}))^\top (\mathbf{y} - \Pi_{\mathcal{X}}(\mathbf{y})) \leq 0$;
2. $\|\mathbf{x} - \Pi_{\mathcal{X}}(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_{\mathcal{X}}(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$.

The two properties equivalently say that the vectors $\mathbf{y} - \Pi_{\mathcal{X}}(\mathbf{y})$ and $\mathbf{y} - \Pi_{\mathcal{X}}(\mathbf{y})$ form an obtuse angle; see Figure 15.



Figure 15. Proof by picture of the properties of the projection step made in projected gradient descent.

4.1 Lipschitz convex functions

Since we minimize f over a closed and bounded convex set \mathcal{X} , we get the existence of a minimizer and bound R for free. Furthermore, if we assume that f is continuously differentiable, we also have a bound B for the gradient norms over \mathcal{X} . This makes the following theorem a much more useful result than in the unconstrained case, since we need to make less assumptions about f .

Theorem 35. Let $f : \operatorname{dom}(f) \rightarrow \mathbb{R}$ be convex and differentiable, $\mathcal{X} \subseteq \operatorname{dom}(f)$ closed and convex, \mathbf{x}^* a minimizer of f over \mathcal{X} . Furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ and $\|\nabla f(\mathbf{x})\| \leq B$ for all $\mathbf{x} \in \mathcal{X}$. Choosing the constant stepsize,

$$\gamma \doteq \frac{R}{B\sqrt{T}},$$

with projected gradient descent yields the following bound,

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{RB}{\sqrt{T}}.$$

Proof. We only need to show that the vanilla analysis still holds for the projected \mathbf{x}_{t+1} . We know that for the non-projected iterate, we get the following equality from the vanilla analysis,

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{2\gamma} \left(\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2 \right).$$

By the second property of projected gradient descent, we know $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2$, hence we get

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \leq \frac{1}{2\gamma} \left(\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right).$$

Then, we can prove the rest of the vanilla analysis as for unbounded gradient descent. ■

This is the same bound as in the unconstrained case, thus the number of necessary iteration is of the order $\mathcal{O}(1/\epsilon^2)$.

4.2 Smooth functions

Lemma 36 (Sufficient decrease of projected gradient descent). Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be differentiable and smooth with parameter L over a closed and convex set $\mathcal{X} \subseteq \text{dom}(f)$. Choosing stepsize

$$\gamma \doteq \frac{1}{L},$$

projected gradient descent satisfies

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 + \underbrace{\frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2}_{\text{additional term}}.$$

Proof.

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - L(\mathbf{y}_{t+1} - \mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \frac{L}{2} \left(\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right) \\ &\quad + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2. \end{aligned}$$

Smoothness.

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \Rightarrow \nabla f(\mathbf{x}_t) = -L(\mathbf{y}_{t+1} - \mathbf{x}_t).$$

$$2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2.$$

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t) \Rightarrow \mathbf{y}_{t+1} - \mathbf{x}_t = -\frac{1}{L} \nabla f(\mathbf{x}_t).$$

Thus, we also have a sufficient decrease lemma for this version of gradient descent, which has an additional term in its bound. However, as we will see, this does not matter, because we can compensate for it in the vanilla analysis.

Theorem 37. Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be convex and differentiable. Let $\mathcal{X} \subseteq \text{dom}(f)$ be a closed convex set, and \mathbf{x}^* the minimizer of f over \mathcal{X} . Furthermore, suppose that f is smooth over \mathcal{X} with parameter L . Choosing stepsize

$$\gamma \doteq \frac{1}{L},$$

projected gradient descent yields the following bound,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Proof. From property 2 of gradient descent, we get the following inequality,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \leq \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2.$$

Using this inequality, we get the following upper bound,

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{1}{2\gamma} \left(\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2 \right) \\ &\leq \frac{1}{2\gamma} \left(\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right. \\ &\quad \left. - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right). \end{aligned}$$

See vanilla analysis, where \mathbf{x}_{t+1} is substituted by \mathbf{y}_{t+1} , since that is the next unconstrained iterate.

We use sufficient decrease to bound the sum of gradients,

$$\begin{aligned} \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 &\leq \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2. \end{aligned}$$

Sufficient decrease.

Telescoping sum.

Now, we can bound the summed error and we will see that the additional terms cancel,

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\quad - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \\ &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \\ &\quad + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \end{aligned}$$

which results in

$$\sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Due to sufficient decrease, we get

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

■

This is the same bound as in the unconstrained case, thus the number of necessary iterations is of the order $\mathcal{O}(1/\epsilon)$.

4.3 Smooth and strongly convex functions

Theorem 38. Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be convex and differentiable. Let $\mathcal{X} \subseteq \text{dom}(f)$ be a nonempty closed and convex set and suppose that f is smooth over \mathcal{X} with parameter L and strongly convex over \mathcal{X} with parameter $\mu > 0$. Choosing stepsize

$$\gamma \doteq \frac{1}{L},$$

projected gradient descent satisfies the following two properties,

1. Squared distance to \mathbf{x}^* are geometrically decreasing,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2;$$

2. The absolute error after T iterations is exponentially small in T ,

$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^*) &\leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\quad + \underbrace{\|\nabla f(\mathbf{x}^*)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^*\|}_{\text{additional term}}. \end{aligned}$$

Proof. These proofs are similar to the normal gradient descent case, starting from the constrained vanilla bound,

$$\frac{1}{2\gamma} \left(\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right),$$

which can be strengthened to

$$\frac{1}{2\gamma} \left(\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2,$$

on $f(\mathbf{x}_t) - f(\mathbf{x}^*)$.

■

5 Coordinate descent

A problem with gradient descent in large-scale learning is that we need to compute the full gradient $\nabla f(x_t)$ in every iteration. The idea of *coordinate descent* is to update only one coordinate of the iterate at a time. To do this, we only need to compute one coordinate of $\nabla f(x_t)$ at a time, which we assume to be a factor of d faster to compute.

However, we also expect to pay a price for this in terms of the number of iterations required. It turns out that, in the worst case, the number of iterations will increase by a factor of d , so we only stand to gain. Under additional assumptions about f , coordinate descent can lead to provable speedups.

Definition 39 (Polyak-Łojasiewicz (PL) inequality). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable function with a global minimum x^* . We say that f satisfies the PL inequality if the following holds for some $\mu > 0$,

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

Observation. The PL inequality directly implies that every critical point is a minimizer of f .

There are non-convex functions that satisfy the PL inequality.

Lemma 40. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and strongly convex with parameter $\mu > 0$. Then f satisfies the PL inequality for the same μ .

To analyze coordinate descent, we need the notion of *coordinate-wise smoothness*.

Definition 41 (Coordinate-wise smoothness). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable, and $\mathcal{L} = [L_1, \dots, L_d] \in \mathbb{R}_+^d$. Function f is called coordinate-wise smooth (with parameter \mathcal{L}) if for every coordinate $i = 1, \dots, d$, the following holds,

$$f(x + \lambda e_i) \leq f(x) + \lambda \nabla_i f(x) + \frac{L_i}{2} \lambda^2, \quad \forall x \in \mathbb{R}^d, \lambda \in \mathbb{R}.$$

Compare this to the definition of smoothness; Definition 29. In our new coordinate-wise definition, we define y as $x + \lambda e_i$, since we only want to change coordinate i . Hence, $y - x$ becomes λe_i . From there, it is easy to see that $\nabla f(x)^\top (y - x)$ becomes $\lambda \nabla_i f(x)$, and $\|x - y\|$ becomes λ . Thus, smoothness with parameter L implies coordinate-wise smoothness with parameter $\mathcal{L} = [L, \dots, L]$.

Example 42. $f(x_1, x_2) = x_1^2 + 10x_2^2$ is smooth with parameter $L = 20$, but f is coordinate-wise smooth with parameter $\mathcal{L} = [2, 20]$. Such differences will become important later when showing faster convergence of coordinate descent.

In general, coordinate (gradient) descent algorithms perform the following actions,

$$\begin{aligned} &\text{choose an active coordinate } i \in [d] \\ &\mathbf{x}_{t+1} \doteq \mathbf{x}_t - \gamma_i \nabla_i f(\mathbf{x}_t) \mathbf{e}_i. \end{aligned}$$

We call the coordinate that is currently being update *active*.

Lemma 43 (Coordinate-wise sufficient decrease). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and coordinate-wise smooth with parameter $\mathcal{L} = [L_1, \dots, L_d]$. With active coordinate i in iteration t and stepsize

$$\gamma_i \doteq \frac{1}{L_i},$$

coordinate descent satisfies

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L_i} |\nabla_i f(\mathbf{x}_t)|^2.$$

Proof. Let $\lambda = -\nabla_i f(\mathbf{x}_t)/L_i$, then $\mathbf{x}_{t+1} = \mathbf{x}_t + \lambda \mathbf{e}_i$. Then, we can apply coordinate-wise smoothness,

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \lambda \nabla_i f(\mathbf{x}_t) + \frac{L_i}{2} \lambda^2 \\ &= f(\mathbf{x}_t) - \frac{1}{L_i} |\nabla_i f(\mathbf{x}_t)|^2 + \frac{1}{2L_i} |\nabla_i f(\mathbf{x}_t)|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L_i} |\nabla_i f(\mathbf{x}_t)|^2. \end{aligned}$$

Coordinate-wise smoothness.

■

5.1 Randomized coordinate descent

In *randomized gradient descent*, the active coordinate is chosen uniformly at random,

$$\begin{aligned} &\text{sample } i \in [d] \text{ uniformly at random} \\ &\mathbf{x}_{t+1} \doteq \mathbf{x}_t - \gamma_i \nabla_i f(\mathbf{x}_t) \mathbf{e}_i. \end{aligned}$$

Randomized coordinate descent is at least as fast as gradient descent on smooth functions if we assume that it is d times cheaper to update one coordinate than the full iterate [Nesterov, 2012]. If we additionally assume the PL inequality, we can obtain faster convergence.

Theorem 44. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with global minimum \mathbf{x}^* . Suppose that f is coordinate-wise smooth with parameter L and satisfies the PL inequality with parameter $\mu > 0$. Choosing stepsize

$$\gamma_i \doteq \frac{1}{L},$$

randomized gradient descent with arbitrary \mathbf{x}_0 yields

$$\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{dL}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

Proof. Sufficient decrease yields

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} |\nabla_i f(\mathbf{x}_t)|^2.$$

By taking the expectation of both sides with respect to i , we have

$$\begin{aligned} \mathbb{E}_i[f(\mathbf{x}_{t+1}) \mid \mathbf{x}_t] &\leq \mathbb{E}_i \left[f(\mathbf{x}_t) - \frac{1}{2L} |\nabla_i f(\mathbf{x}_t)|^2 \mid \mathbf{x}_t \right] \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \sum_{i=1}^d \frac{1}{d} |\nabla_i f(\mathbf{x}_t)|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2dL} \|\nabla f(\mathbf{x}_t)\|^2 \\ &\leq f(\mathbf{x}_t) - \frac{\mu}{dL} (f(\mathbf{x}_t) - f(\mathbf{x}^*)). \end{aligned}$$

PL inequality.

Subtracting $f(\mathbf{x}^*)$ from both sides and taking the expectation over \mathbf{x}_t , we obtain

$$\mathbb{E}[f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{dL}\right) (f(\mathbf{x}_t) - f(\mathbf{x}^*)).$$

$\mathbb{E}_i[f(\mathbf{x}_{t+1}) \mid \mathbf{x}_t]$ is a random variable whose expectation is $\mathbb{E}_{\mathbf{x}_t}[f(\mathbf{x}_{t+1})]$.

The statement follows. ■

5.2 Importance sampling

As seen, uniformly random selection of the active coordinate does not yield a better bound than gradient descent. However, we have not made use of the fact that the coordinate-wise smoothness parameters can differ. Intuitively, we would want to sample coordinates with high smoothness more frequently than coordinates with low smoothness, since they change more rapidly, leading to faster convergence to the optimum. This leads us to *importance sampling* [Nesterov, 2012],

$$\begin{aligned} &\text{sample } i \in [d] \text{ with probability } \frac{L_i}{\sum_{j=1}^d L_j} \\ &\mathbf{x}_{t+1} \doteq \mathbf{x}_t - \gamma_i \nabla_i f(\mathbf{x}_t) \mathbf{e}_i. \end{aligned}$$

Theorem 45. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a global minimum \mathbf{x}^* . Suppose that f is coordinate-wise smooth with parameter $\mathcal{L} = [L_1, \dots, L_d]$ and satisfies the PL inequality with parameter $\mu > 0$. Let

$$\bar{L} = \frac{1}{d} \sum_{i=1}^d L_i$$

be the average of coordinate-wise smoothness constants. Then, coordinate descent with importance sampling and arbitrary \mathbf{x}_0 yields

$$\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{d\bar{L}}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

Proof. The proof is nearly identical to the proof of randomized coordinate descent. The difference lies in the expectation over i . Importance sampling yields

$$\begin{aligned} \mathbb{E}_i \left[f(\mathbf{x}_t) - \frac{1}{2L_i} |\nabla_i f(\mathbf{x}_t)|^2 \mid \mathbf{x}_t \right] &= f(\mathbf{x}_t) - \sum_{i=1}^d \frac{L_i}{\sum_{j=1}^d L_j} \frac{1}{2L_i} |\nabla_i f(\mathbf{x}_t)|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2d\bar{L}} \sum_{i=1}^d |\nabla_i f(\mathbf{x}_t)|^2. \end{aligned} \quad d\bar{L} = \sum_{i=1}^d L_i.$$

The rest follows identically. ■

Note that \bar{L} can be much smaller than $L = \max_{i=1}^d L_i$, so coordinate descent with important sampling is potentially faster than randomized gradient descent. In the worst-case, both algorithms are the same.

5.3 Steepest coordinate descent

In contrast to random coordinate descent, *steepest coordinate descent* chooses the active coordinate according to the coordinate with the largest gradient (*Gauss-Southwell* rule),

$$\begin{aligned} \text{choose } i &= \underset{i \in [d]}{\operatorname{argmax}} |\nabla_i f(\mathbf{x}_t)| \\ \mathbf{x}_{t+1} &\doteq \mathbf{x}_t - \gamma_i \nabla_i f(\mathbf{x}_t) \mathbf{e}_i. \end{aligned}$$

The main difference from the previous algorithms is that this algorithm is deterministic, thus we do not need to take the expectation.

Corollary. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a global minimum \mathbf{x}^* . Suppose that f is coordinate-wise smooth with parameter L and satisfies the PL inequality with parameter $\mu > 0$. Choosing stepsize

$$\gamma_i \doteq \frac{1}{L},$$

steepest coordinate descent with arbitrary \mathbf{x}_0 yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \left(1 - \frac{\mu}{dL}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

This is not good. It needs the same amount of iterations as randomized coordinate descent, but each iteration takes as long as in gradient descent.³ However, this algorithm allows for a speedup in certain cases. Furthermore, it may be possible to efficiently maintain the maximum absolute gradient value throughout the iterations, so that the full evaluation of the gradient can be avoided.

Nutini et al. [2015] showed that a better convergence result can be obtained for strongly convex functions, when strong convexity is measured w.r.t. the ℓ_1 norm instead of the standard Euclidean norm, *i.e.*,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu_1}{2} \|\mathbf{y} - \mathbf{x}\|_1^2.$$

Due to $\|\mathbf{y} - \mathbf{x}\|_1 \geq \|\mathbf{y} - \mathbf{x}\|$, f is then also strongly convex with $\mu = \mu_1$. On the other hand, if f is μ -strongly convex w.r.t. the ℓ_2 norm, then f is μ/d -strongly convex w.r.t. the ℓ_1 norm, due to $\|\mathbf{y} - \mathbf{x}\| \geq \|\mathbf{y} - \mathbf{x}\|_1 / \sqrt{d}$.

Lemma 46. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and strongly convex with parameter $\mu_1 > 0$ w.r.t. the ℓ_1 norm. Then, f is μ_1 -strongly convex w.r.t. the Euclidean norm, so a global minimum \mathbf{x}^* exists. Then, f satisfies the PL inequality w.r.t. the ℓ_∞ norm with the same μ_1 ,

$$\frac{1}{2} \|\nabla f(\mathbf{x})\|_\infty^2 \geq \mu_1 (f(\mathbf{x}) - f(\mathbf{x}^*)).$$

Theorem 47. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a global minimum \mathbf{x}^* . Suppose that f is coordinate-wise smooth with parameter L and satisfies the PL inequality with parameter $\mu_1 > 0$. Choosing stepsize

$$\gamma_i \doteq \frac{1}{L},$$

steepest coordinate descent with arbitrary \mathbf{x}_0 yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \left(1 - \frac{\mu_1}{L}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

³Note that a function may be coordinate-wise smooth with an L for all coordinates that is smaller than smoothness, so it is not completely fair to compare this to gradient descent.

Proof.

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) - \frac{1}{2L} |\nabla_i f(\mathbf{x}_t)|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|_\infty^2 \\ &\leq f(\mathbf{x}_t) - \frac{\mu_1}{L} (f(\mathbf{x}_t) - f(\mathbf{x}^*)). \end{aligned}$$

Subtracting $f(\mathbf{x}^*)$ from both sides yields

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{\mu_1}{L}\right) (f(\mathbf{x}_t) - f(\mathbf{x}^*)).$$

The statement follows. ■

We see that if $\mu_1 = \mu/d$, we do not gain anything. However, this is not the case in general. If, for the worst case \mathbf{x}, \mathbf{y} , which satisfy $\|\mathbf{y} - \mathbf{x}\| = \|\mathbf{y} - \mathbf{x}\|_1 / \sqrt{d}$, strong convexity holds with $\mu' > \mu$, then we can achieve $\mu_1 = \mu'/d > \mu/d$, resulting in better convergence.

5.4 Greedy coordinate descent

Greedy coordinate descent is a variant that does not even require f to be differentiable. In each iteration, we make the step that maximizes the progress in the chosen coordinate. This requires performing a *line search* by solving a one-dimensional optimization problem,

$$\begin{aligned} &\text{choose } i \in [d] \\ &\mathbf{x}_{t+1} \doteq \operatorname{argmin}_{\lambda \in \mathbb{R}} f(\mathbf{x}_t + \lambda \mathbf{e}_i). \end{aligned}$$

However, greedy coordinate descent can get stuck in non-optimal points; see Figure 16. Thus, we need some additional conditions to make sure this does not occur.

Theorem 48. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be of the form

$$f(\mathbf{x}) \doteq g(\mathbf{x}) + h(\mathbf{x}), \quad h(\mathbf{x}) = \sum_{i=1}^d h_i(x_i), \quad \mathbf{x} \in \mathbb{R}^d,$$

with g convex and differentiable, and the h_i convex.

Let \mathbf{x} be a point such that greedy coordinate descent cannot make progress in any coordinate. Then \mathbf{x} is a global minimum of f .

In the context of machine learning, an important class of functions that satisfies the conditions of Theorem 48 is the following form,

$$f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1,$$

where $\lambda \|\mathbf{x}\|_1$ is a separable ℓ_1 -regularization term used in for example LASSO [Tibshirani, 1996].

Sufficient decrease.

Active coordinate of steepest coordinate descent is the maximum gradient.

PL inequality.

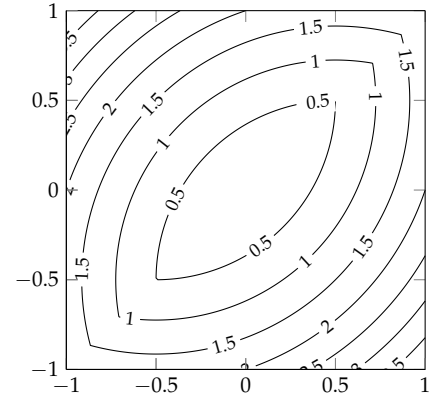


Figure 16. Level set plot of $f(\mathbf{x}) = \|\mathbf{x}\|^2 + |x_1 - x_2|$. The global minimum is $[0, 0]$, but greedy coordinate descent cannot escape any point $[\mathbf{x}, \mathbf{x}]$, s.t. $|\mathbf{x}| \leq 1/2$.

6 Nonconvex functions

So far, all convergence results that we have proved have been for variants of gradient descent on convex functions. The reason for this is that, in general, we cannot expect gradient descent to come close to the global minimum x^* of nonconvex functions. Figures 17 to 19 show what can go wrong in nonconvex functions, under the assumption that we have set the γ such that we do not overshoot. These figures show points that gradient descent cannot escape; local minima and saddle points. Furthermore, it might even be that gradient descent never converges to a critical point; see Figure 19.

In practice, gradient descent works well on the nonconvex functions that we care about. But, theoretical explanations for this are mostly missing. Despite this, we will show that under favorable conditions, we can still say something useful about the behavior of gradient descent on nonconvex functions.

We can easily make an analysis of gradient descent on smooth functions. A useful property that we will use is that functions with bounded Hessians are smooth, as shown in the following lemma.

Lemma 49. Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be twice differentiable with $\mathcal{X} \in \text{dom}(f)$ a convex set and $\|\nabla^2 f(x)\| \leq L$ for all $x \in \mathcal{X}$. Then f is smooth with parameter L over \mathcal{X} .

Proof. Bounded Hessians imply Lipschitz continuity of the gradient,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathcal{X}.$$

We will use the fundamental theorem of calculus with

$$h(t) \doteq f(x + t(y - x)), \quad t \in [0, 1].$$

The derivative can be calculated by chain rule,

$$h'(t) = \nabla f(x + t(y - x))^\top (y - x).$$

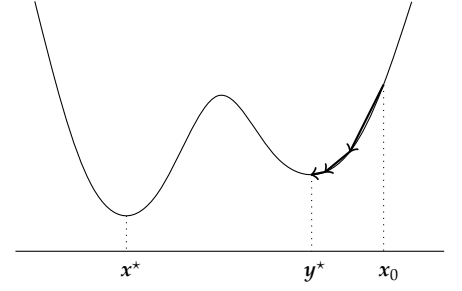


Figure 17. Gradient descent may get stuck in a local minimum $y^* \neq x^*$, in nonconvex functions.

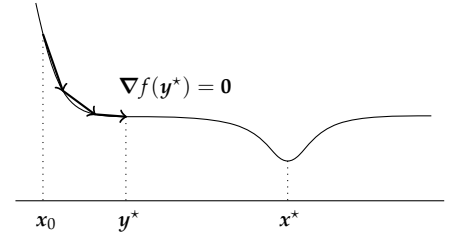


Figure 18. Gradient may get stuck in a flat region (saddle point) in nonconvex functions.

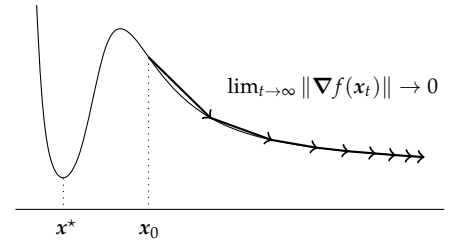


Figure 19. Gradient descent may never even reach a critical point in nonconvex functions.

Now, we can show smoothness,

$$\begin{aligned}
& f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \\
&= h(1) - h(0) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) && \text{Definition of } h. \\
&= \int_0^1 h'(t) dt - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) && \text{Fundamental theorem of calculus.} \\
&= \int_0^1 \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) dt - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) && \text{Fill in } h'(t). \\
&= \int_0^1 \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) dt && \text{Integral from 0 to 1.} \\
&= \int_0^1 (\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) dt \\
&\leq \int_0^1 |(\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x})| dt \\
&\leq \int_0^1 \|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| dt && \text{Cauchy-Schwarz inequality.} \\
&\leq \int_0^1 L \|t(\mathbf{y} - \mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| dt && \text{Lipschitz continuous gradient.} \\
&= \int_0^1 Lt \|\mathbf{x} - \mathbf{y}\|^2 dt \\
&= \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.
\end{aligned}$$

Thus, we have smoothness,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

■

Now, we can use this fact and sufficient decrease⁴ to prove that the gradients of smooth functions are bounded and approach 0, as we increase the number of iterations.

⁴ Recall that sufficient decrease did not require convexity.

Theorem 50. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable with a global minimum \mathbf{x}^* . Furthermore, suppose that f is smooth with parameter L . Choosing stepsize

$$\gamma \doteq \frac{1}{L},$$

gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{2L}{T} (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

Remark. Note that concave functions are not a counter example to this theorem, despite their gradients growing, because they have no global minimum \mathbf{x}^* .

Proof. Recall that sufficient decrease does not require convexity,

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

Rewriting this, we get

$$\|\nabla f(\mathbf{x}_t)\|^2 \leq 2L(f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})).$$

Then, by telescoping sum, we get

$$\sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq 2L(f(\mathbf{x}_0) - f(\mathbf{x}_T)) \leq 2L(f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

The statement follows by dividing both sides by T . ■

This has the result that

$$\lim_{t \rightarrow \infty} \|\nabla f(\mathbf{x}_t)\|^2 = 0.$$

It might seem that convergence of the gradients to 0 is the same as convergence to a critical point. But, this interpretation does not hold in general; see Figure 19. In this case, the gradient converges to 0, but the iterates only move further away from the critical point. So, this is not a very strong result.

6.1 Trajectory analysis

Despite the fact that a nonconvex function may contain local minima, saddle points, and flat parts, gradient descent may avoid them and still converge to a global minimum. For this, you need a good starting point and do a trajectory analysis. As an example, we will do a trajectory analysis for a simplified deep linear neural network. It turns out that this function is smooth along the trajectories that we analysis, and this is the most important ingredient of the analysis.

Let

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_\ell\}$$

be the weights of the deep linear network. And, in general, we want to approximate a matrix \mathbf{Y} , given input matrix \mathbf{X} . Thus, we want to minimize

$$\|\mathbf{W}_\ell \mathbf{W}_{\ell-1} \cdots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2.$$

Arora et al. [2018] consider this general framework, but we will only consider the case where all matrices are 1×1 , *i.e.*, scalars. Assume we have training input $x = 1$ and output $y = 1$, then we have the following function to optimize,

$$f(\mathbf{x}) \doteq \frac{1}{2} \left(\prod_{k=1}^d x_k - 1 \right)^2.$$

We can immediately see that setting $x_k = 1$ for all k minimizes the function at 0. However, we want to know whether gradient descent will also be able to find this set of weights.

Note that stacking linear layers has no benefit, since any stacking of linear layers can be represented by a single linear layer. However, the reason for doing this is that it gives us a simple playground in which we can try to understanding why training deep neural networks with gradient descent works, despite the fact that the objective function is nonconvex.

We rewrite w as x and ℓ as d to be more in line with the notation used here.

The gradient of this function is computed by

$$\nabla_i f(\mathbf{x}) = \left(\prod_{k=1}^d x_k - 1 \right) \prod_{k \neq i}^d x_k$$

Whenever at least two dimensions are zero, the gradient vanishes. Thus, any \mathbf{x} with two zero entries are critical points, despite not being global minima, since then the product of all entries must be 1, which is not possible if at least two are zero.⁵ We know that the value of all such saddle points is $1/2$.

We now want to show that for any number of layers (*i.e.* dimensionality of \mathbf{x}), anywhere in $\mathcal{X} = \{\mathbf{x} \mid \mathbf{x} > \mathbf{0}, \prod_{k=1}^d x_k \leq 1\}$, despite that f is not smooth over \mathcal{X} . However, we only need to show that f is smooth along the trajectory of gradient descent for suitable L , so that we get sufficient decrease. We will now show this by showing that the Hessians over the trajectory are bounded. The Hessian is given by

$$\nabla_{ij}^2 f(\mathbf{x}) = \begin{cases} \left(\prod_{k \neq i}^d x_k \right)^2, & j = i \\ 2 \prod_{k \neq i}^d x_k \prod_{k \neq j}^d x_k - \prod_{k \neq i, j}^d x_k, & j \neq i. \end{cases}$$

Definition 51 (*c*-balanced.). Let $\mathbf{x} > \mathbf{0}$ and $c \geq 1$. \mathbf{x} is called *c*-balanced if $x_i \leq cx_j$ for all $1 \leq i, j \leq d$.

Lemma 52. Let $\mathbf{x} > \mathbf{0}$ be *c*-balanced with $\prod_k x_k \leq 1$, then for any stepsize $\gamma > 0$, $\mathbf{x}' \doteq \mathbf{x} - \gamma \nabla f(\mathbf{x})$ satisfies $\mathbf{x}' \geq \mathbf{x}$ componentwise, and is also *c*-balanced.

Proof. Let

$$\Delta \doteq -\gamma \left(\prod_{k=1}^d x_k - 1 \right) \left(\prod_{k=1}^d x_k \right) \geq 0.$$

Then,

$$-\gamma \nabla_k f(\mathbf{x}) = \frac{\Delta}{x_k}.$$

Thus, the gradient descent update has the following form,

$$x'_k = x_k + \frac{\Delta}{x_k} \geq x_k, \quad k \in [d].$$

For all i, j , we thus get

$$x'_i = x_i + \frac{\Delta}{x_i} \leq cx_j + \frac{c\Delta}{x_j} = cx'_j.$$

■

So, we know now that all iterates are *c*-balanced if \mathbf{x}_0 is balanced with *c*. We can use this to compute a bound on the Hessian by bounding the products.

⁵ This shows that f is nonconvex, since local minima are global minima in convex functions.

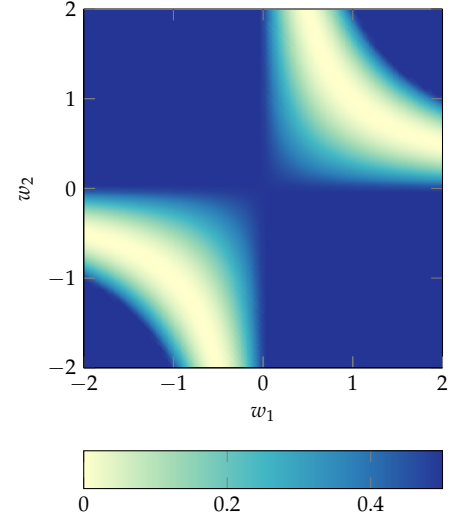


Figure 20. $f(\mathbf{x}) = \frac{1}{2} (\prod_k x_k - 1)^2$ for $d = 2$, where the loss is clipped to be at most $1/2$.

$$x_j \leq cx_i \iff \frac{1}{x_i} \leq \frac{c}{x_j}.$$

Lemma 53. Suppose $x > \mathbf{0}$ is c -balanced. Then, for any $I \subseteq [d]$, we have

$$\prod_{k \notin I} x_k \leq c^{|I|} \left(\prod_{k=1}^d x_k \right)^{1-|I|/d} \leq c^{|I|}.$$

Proof. For any $i \in [d]$, we have $c^d \cdot x_i^d \geq \prod_k x_k$ by c -balancedness, hence $x_i \geq \frac{1}{c} (\prod_k x_k)^{1/d}$. Thus,

$$\prod_{k \notin I} x_k = \frac{\prod_k x_k}{\prod_{i \in I} x_i} \leq \frac{\prod_k x_k}{(1/c)^{|I|} (\prod_k x_k)^{|I|/d}} = c^{|I|} \left(\prod_{k=1}^d x_k \right)^{1-|I|/d}.$$

Since $\prod_k x_k \leq 1$, we can bound the above by $c^{|I|}$. ■

Lemma 54. Let $x > \mathbf{0}$ be c -balanced with $\prod_k x_k \leq 1$, then

$$\|\nabla^2 f(x)\| \leq \|\nabla^2 f(x)\|_F \leq 3dc^2.$$

Proof. The fact that $\|A\| \leq \|A\|_F$ is well known. To bound the Frobenius norm, we use the previous lemma to compute

$$|\nabla_{ii}^2 f(x)| = \left| \left(\prod_{k \neq i} x_k \right)^2 \right| \leq c^2,$$

and for $i \neq j$, we get

$$|\nabla_{ij}^2 f(x)| \leq \left| 2 \prod_{k \neq i} x_k \prod_{k \neq j} x_k \right| + \left| \prod_{k \neq i, j} x_k \right| \leq 3c^2.$$

Thus,

$$\|\nabla^2 f(x)\|_F^2 \leq 9d^2 c^4.$$

Taking the square root, the statement follows. ■

This lemma implies smoothness of f with parameter $L = 3dc^2$ along the whole trajectory of gradient descent, under the “smooth stepsize” $\gamma \doteq 1/L = 1/3dc^2$. Now, we can use this to prove convergence.

Theorem 55. Let $c \geq 1$ and $\delta < 0$ such that $x_0 > \mathbf{0}$ is c -balanced with $\delta \leq \prod_k (x_0)_k < 1$. Choosing stepsize

$$\gamma \doteq \frac{1}{3dc^2},$$

gradient descent satisfies

$$f(x_T) \leq \left(1 - \frac{\delta^2}{3c^4} \right)^T f(x_0).$$

Proof. For each $t \geq 0$, f is smooth over $\text{conv}(\{\mathbf{x}_t, \mathbf{x}_{t+1}\})$ with parameter $L = 3dc^2$, hence we have sufficient decrease,

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{6dc^2} \|\nabla f(\mathbf{x}_t)\|^2.$$

For every c -balanced \mathbf{x} with $\delta \leq \prod_k x_k \leq 1$, we have

$$\begin{aligned} \|\nabla f(\mathbf{x}_t)\|^2 &= 2f(\mathbf{x}) \sum_{i=1}^d \left(\prod_{k \neq i} x_k \right)^2 \\ &\geq 2f(\mathbf{x}) \frac{d}{c^2} \left(\prod_{k=1}^d x_k \right)^{2-2/d} \\ &\geq 2f(\mathbf{x}) \frac{d}{c^2} \left(\prod_{k=1}^d x_k \right)^2 \\ &\geq 2f(\mathbf{x}) \frac{d}{c^2} \delta^2. \end{aligned}$$

Hence,

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{6dc^2} 2f(\mathbf{x}_t) \frac{d}{c^2} \delta^2 = f(\mathbf{x}_t) \left(1 - \frac{\delta^2}{3c^4} \right).$$

■

Thus, we seem to have fast convergence, since the function value goes down by a constant factor in each step. However, there is a catch. Consider the $\mathbf{x}_0 = [1/2, \dots, 1/2]$, which is c -balanced with $c = 1$, and $\delta = 1/2^d$. Hence, the constant factor is

$$1 - \frac{1}{3 \cdot 4^d}.$$

This means that we would need $T \approx 4^d$ iterations to reduce the initial error by a constant factor not depending on d . Hence, for this starting value, the gradient is exponentially small. In order to get polynomial convergence, we need to start with a δ that decays at most polynomially with d . For large d , this has the consequence that we must start very close to optimality. In particular, we need to start at a distance $\mathcal{O}(1/\sqrt{d})$ from the optimal solution $[1, \dots, 1]$.

7 The Frank-Wolfe algorithm

Projected gradient descent is the only algorithm we have seen that dealt with constrained optimization problems. However, that algorithm came with the clear disadvantage that projections can be very expensive, even when \mathcal{X} is convex. The Frank-Wolfe algorithm solves constrained optimization problems without projection steps. Instead, it makes use of a *linear minimization oracle* (LMO). For the feasible region $\mathcal{X} \subseteq \mathcal{R}^d$ and an arbitrary vector $g \in \mathcal{R}^d$,⁶

$$\text{LMO}_{\mathcal{X}}(g) \doteq \underset{z \in \mathcal{X}}{\operatorname{argmin}} g^{\top} z.$$

Notice that this is the minimization of a linear function.

The Frank-Wolfe algorithm iteratively updates by calling the oracle in the direction of the gradient,

$$\begin{aligned} s &\doteq \text{LMO}_{\mathcal{X}}(\nabla f(x_t)) \\ x_{t+1} &\doteq (1 - \gamma_t)x_t + \gamma_t s. \end{aligned}$$

The algorithm reduces non-linear constrained optimization to linear optimization over the same set \mathcal{X} ; it is able to solve general non-linear constrained optimization problems by only solving a simpler linear constrained optimization problem over the same set \mathcal{X} in each iteration, by calling the oracle. We solve this linear optimization problem in the direction of the gradient, since that is the best linear approximation of f at x_t .

A nice property of the oracle is that if $\mathcal{X} = \operatorname{conv}(\mathcal{A})$, then $\text{LMO}_{\mathcal{X}}(x) \in \mathcal{A}$. So, if we have a set \mathcal{X} that is the convex hull of a small number of points, such as the ℓ_1 -ball, we have an easy optimization problem with runtime $\mathcal{O}(|\mathcal{A}|)$.

The advantages of this method are

- Iterates are always feasible if \mathcal{X} is convex;
- No projections, which are often harder to compute than linear optimization problems;
- Iterates always have a simple sparse representations: x_t is a convex combination of x_0 and all s used so far.

7.1 Linear minimization oracles

LASSO. The LASSO problem in its standard form is given by

$$\begin{aligned} &\text{minimize} \quad \|Ax - b\|^2 \\ &\text{subject to} \quad \|x\|_1 \leq 1. \end{aligned}$$

The constraint set $\mathcal{X} = \{x \in \mathbb{R}^d \mid \|x\|_1 \leq 1\}$ is the unit ℓ_1 -ball. This is the convex hull of the unit basis vectors; $\mathcal{X} = \operatorname{conv}(\{\pm e_1, \dots, \pm e_d\})$. The



Figure 21. Illustration of a Frank-Wolfe step.

⁶ g can be thought of as an “optimization direction”.

Move toward minimizer.

LMO for this set is easy to compute,

$$\begin{aligned}\text{LMO}_{\mathcal{X}}(g) &= \underset{z \in \mathcal{X}}{\operatorname{argmin}} z^\top g \\ &= \underset{z \in \{\pm e_1, \dots, \pm e_d\}}{\operatorname{argmin}} z^\top g \\ &= -\operatorname{sign}(g_i) e_i, \quad i \doteq \underset{i \in [d]}{\operatorname{argmax}} |g_i|.\end{aligned}$$

So, we only have to identify g 's largest coordinate, which is much more efficient than projection onto an ℓ_1 -ball.

7.2 Duality gap

We define the duality gap of $x \in \mathcal{X}$ as

$$g(x) \doteq \nabla f(x)^\top (x - s), \quad s \doteq \text{LMO}_{\mathcal{X}}(\nabla f(x)).$$

This can be interpreted as the optimality gap $\nabla f(x)^\top x - \nabla f(x)^\top s$ of the linear subproblem.

Lemma 56. Suppose that the constrained minimization problem has a minimizer x^* . Let $x \in \mathcal{X}$, then

$$g(x) \geq f(x) - f(x^*),$$

meaning that the duality gap is an upper bound for the optimality gap.

Proof.

$$\begin{aligned}g(x) &= \nabla f(x)^\top (x - s) \\ &\geq \nabla f(x)^\top (x - x^*) \\ &\geq f(x) - f(x^*).\end{aligned}$$

s is the minimizer of $f(z) \doteq \nabla f(x)^\top z$.

First-order characterization of convexity:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x), \forall x, y.$$

■

Thus, we always have a computable upper bound $g(x_t)$ on the unknown error $f(x_t) - f(x^*)$, which is not the case in general for *unconstrained* optimization. Furthermore, at an optimal point x^* , $g(x^*) = 0$, which follows from the optimality conditions for constrained convex optimization.⁷

⁷ $\nabla f(x)^\top (x - x^*) \geq 0, \forall x \in \mathcal{X}$.

7.3 Convergence analysis

To prove convergence, we need the following descent lemma,

Lemma 57. For a step $\mathbf{x}_{t+1} \doteq \mathbf{x}_t + \gamma_t(\mathbf{s} - \mathbf{x}_t)$ with stepsize $\gamma_t \in [0, 1]$, it holds that

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma_t g(\mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s} - \mathbf{x}_t\|^2,$$

where $\mathbf{s} = \text{LMO}_{\mathcal{X}}(\nabla f(\mathbf{x}_t))$.

Proof.

$$\begin{aligned} f(\mathbf{x}_{t+1}) &= f(\mathbf{x}_t + \gamma_t(\mathbf{s} - \mathbf{x}_t)) \\ &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top \gamma_t(\mathbf{s} - \mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s} - \mathbf{x}_t\|^2 \\ &= f(\mathbf{x}_t) - \gamma_t g(\mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s} - \mathbf{x}_t\|^2. \end{aligned}$$

Smoothness:

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Duality gap: $g(\mathbf{x}) \doteq \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{s})$.

Now, we can prove the main convergence theorem of the Frank-Wolfe algorithm.

Theorem 58. Consider the constrained minimization problem where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and smooth with parameter L , and \mathcal{X} is convex, closed and bounded. (This means that the minimizer \mathbf{x}^* of f over \mathcal{X} exists and that all minimization oracles have minimizers.) With any $\mathbf{x}_0 \in \mathcal{X}$ and stepsizes

$$\gamma_t \doteq \frac{2}{t+2},$$

the Frank-Wolfe algorithm yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2L}{T+1} \text{diam}(\mathcal{X})^2,$$

where $\text{diam}(\mathcal{X}) \doteq \max_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \|\mathbf{x} - \mathbf{y}\|$ is the diameter of \mathcal{X} .

Proof. Let $h(\mathbf{x}) \doteq f(\mathbf{x}) - f(\mathbf{x}^*)$ and $C \doteq \frac{L}{2} \text{diam}(\mathcal{X})^2$.

$$\begin{aligned} f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) &\leq f(\mathbf{x}_t) - \gamma_t g(\mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s} - \mathbf{x}_t\|^2 - f(\mathbf{x}^*) \\ h(\mathbf{x}_{t+1}) &\leq h(\mathbf{x}_t) - \gamma_t h(\mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s} - \mathbf{x}_t\|^2 \\ &= (1 - \gamma_t) h(\mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s} - \mathbf{x}_t\|^2 \\ &\leq (1 - \gamma_t) h(\mathbf{x}_t) + \gamma_t^2 C. \end{aligned}$$

Lemma 57 and subtract $f(\mathbf{x}^*)$ from both sides.

Duality gap, $g(\mathbf{x}) \geq h(\mathbf{x})$.

Using our new definitions, we want to prove

$$h(\mathbf{x}_T) \leq \frac{4C}{T+1}, \quad T \geq 1.$$

We can prove this by induction. The base case is $T = 1$,

$$\begin{aligned} h(\mathbf{x}_1) &\leq (1 - \gamma_0)h(\mathbf{x}_0) + \gamma_0^2 C \\ &= C \\ &\leq 2C. \end{aligned} \quad \gamma_0 = 1.$$

Suppose it holds for $T = k$,

$$h(\mathbf{x}_k) \leq \frac{4C}{k+1}.$$

then we need to show that it also holds for $T = k + 1$,

$$\begin{aligned} h(\mathbf{x}_{k+1}) &\leq (1 - \gamma_k)h(\mathbf{x}_k) + \gamma_k^2 C \\ &\leq \left(1 - \frac{2}{k+2}\right) \frac{4C}{k+1} + \left(\frac{2}{k+2}\right)^2 C \\ &= \frac{k}{k+2} \frac{4C}{k+1} + \frac{4}{(k+2)^2} C \\ &= \frac{4C}{k+2} \left(\frac{k}{k+1} + \frac{1}{k+2}\right) \\ &\leq \frac{4C}{k+2}. \end{aligned} \quad \text{Induction step and } \gamma_k = 2/(k+2).$$

Thus, it holds for all $T \geq 1$. ■

Affine invariance. Consider the problem of minimizing $f(x_1, x_2) \doteq x_1^2 + x_2^2$ over the unit square $\mathcal{X} = \{[x_1, x_2] \mid 0 \leq x_1 \leq 1, 0 \leq x_2 \leq 1\}$. This function is smooth with $L = 2$ and $\text{diam}(\mathcal{X})^2 = 2$. This has the following error bound,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{8}{T+1}.$$

Next consider $f'(x_1, x_2) \doteq x_1^2 + (10x_2)^2$ over the rectangle $\mathcal{X}' = \{[x_1, x_2] \mid 0 \leq x_1 \leq 1, 0 \leq x_2 \leq 1/10\}$. This function is smooth with $L' = 200$ and $\text{diam}(\mathcal{X}')^2 = 1 + 1/100$. Thus, f' has the error bound

$$f'(\mathbf{x}_T) - f'(\mathbf{x}^*) \leq \frac{404}{T+1}.$$

Hence, according to our analysis, it seems that the error of the Frank-Wolfe algorithm on f' over \mathcal{X}' is roughly 50 times larger than on f over \mathcal{X} .

However, when we look more closely at the function, the two problems (f, \mathcal{X}) and (f', \mathcal{X}') are equivalent under a rescaling of x_2 . The Frank-Wolfe algorithm is invariant under affine transformations, thus there should be no difference in the analysis of the two considered problems.

Formally, two problems (f, \mathcal{X}) and (f', \mathcal{X}') are *affinely equivalent* if $f'(x) = f(Ax + b)$ and $\mathcal{X}' = \{A^{-1}(x - b) \mid x \in \mathcal{X}\}$ for some invertible matrix A and some vector b . The consequence is that $f(x) = f'(x')$ if $x' = A^{-1}(x - b)$.⁸

⁸ In our example problem, this means that we have

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 10 \end{bmatrix}, \quad b = 0.$$

By the chain rule, we get

$$\nabla f'(x') = A^\top \nabla f(Ax' + b) = A^\top f(x).$$

Consider the iterate x_k and its corresponding iterate $x'_k = A^{-1}(x_k - b)$, in their respective problems. We can compute the their oracle calls by

$$\begin{aligned} \text{LMO}_{\mathcal{X}}(\nabla f(x_k)) &= \underset{z \in \mathcal{X}}{\operatorname{argmin}} \nabla f(x_k)^\top z \\ &\doteq s \\ \text{LMO}_{\mathcal{X}'}(\nabla f'(x'_k)) &= \underset{z' \in \mathcal{X}'}{\operatorname{argmin}} \nabla f'(x'_k)^\top z' \\ &= \underset{A^{-1}(z-b) \in \mathcal{X}'}{\operatorname{argmin}} \nabla f(x_k)^\top A A^{-1}(z-b) \\ &= \underset{A^{-1}(z-b) \in \mathcal{X}'}{\operatorname{argmin}} \nabla f(x_k)^\top z - \nabla f(x_k)^\top b \\ &= A^{-1} \left(\left(\underset{z \in \mathcal{X}}{\operatorname{argmin}} \nabla f(x_k)^\top z \right) - b \right) \\ &= A^{-1}(s - b). \end{aligned}$$

Thus, the step directions s' and s also correspond to each other under the affine transformation. As a consequence, the next iterates will also correspond to each other,

$$\begin{aligned} x_{k+1} &= (1 - \gamma)x_k + \gamma s \\ x'_{k+1} &= (1 - \gamma)x'_k + \gamma s' \\ &= (1 - \gamma)A^{-1}(x_k - b) + \gamma A^{-1}(s - b) \\ &= A^{-1}((1 - \gamma)x_k + \gamma s) - b \\ &= A^{-1}x_{k+1} - b. \end{aligned}$$

In particular, after any number of steps, both problems will incur the same optimization error. Thus, we need a better analysis that provides a bound that is invariant under affine transformations. For this, we define a *curvature constant*,

$$C_{(f, \mathcal{X})} \doteq \sup_{\substack{x, s \in \mathcal{X}, \gamma \in (0, 1] \\ y = (1 - \gamma)x + \gamma s}} \frac{1}{\gamma^2} \left(f(y) - f(x) - \nabla f(x)^\top (y - x) \right).$$

This quantity serves as a notion of complexity of both the objective function f and the constraint set \mathcal{X} . It is essentially the supremum of the normalized pointwise vertical distance between the graph of f , $f(y)$ and its linear approximation at x , $f(x) + \nabla f(x)^\top (y - x)$.

Theorem 59. Consider the constrained minimization problem, where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, and \mathcal{X} is convex, closed and bounded. Let $C_{(f,\mathcal{X})}$ be the curvature constant of f over \mathcal{X} . With any $x_0 \in \mathcal{X}$ and with stepsizes

$$\gamma_t = \frac{2}{t+2},$$

the Frank-Wolfe algorithm yields

$$f(x_T) - f(x^*) \leq \frac{4C_{(f,\mathcal{X})}}{T+1}.$$

Proof. We can regain the descent lemma by rewriting the curvature constant. We know by the definition of the supremum,

$$\frac{1}{\gamma^2} \left(f(y) - f(x) - \nabla f(x)^\top (y - x) \right) \leq C_{(f,\mathcal{X})}.$$

$$\forall x, s \in \mathcal{X}, \gamma \in (0,1], y = (1-\gamma)x + \gamma s.$$

Setting the variables,

$$x \doteq x_t, \quad y \doteq x_{t+1} = (1-\gamma_t)x_t + \gamma_t s, \quad y - x = -\gamma_t(x - s),$$

we get

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \nabla f(x_t)^\top (-\gamma_t(x - s)) + \gamma_t^2 C_{(f,\mathcal{X})} \\ &= f(x_t) - \gamma_t g(x_t) + \gamma_t^2 C_{(f,\mathcal{X})}. \end{aligned}$$

The rest of the proof follows as in the previous analysis. ■

You might suspect that this bound is worse than the best bound obtainable from the previous analysis. However, one can show

$$C_{(f,\mathcal{X})} \leq \frac{L}{2} \text{diam}(\mathcal{X})^2.$$

Furthermore, we can prove a convergence of the duality gap.

Theorem 60. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$, \mathcal{X} convex with $C_{(f,\mathcal{X})} < \infty$, $x_0 \in \mathcal{X}$, and $T \geq 2$. Then, choosing stepsize

$$\gamma_t = \frac{2}{t+2},$$

the Frank-Wolfe algorithm yields a t with $1 \leq t \leq T$, such that

$$g(x_t) \leq \frac{27/2 \cdot C_{(f,\mathcal{X})}}{T+1}.$$

8 Newton's method

The *Newton-Raphson method* is an iterative method for finding a zero of a differentiable univariate function $f : \mathbb{R} \rightarrow \mathbb{R}$. Starting from some x_0 , it iteratively computes

$$x_{t+1} \doteq x_t - \frac{f(x_t)}{f'(x_t)}.$$

In formulas, x_{t+1} is the solution to the linear equation

$$f(x_t) + f'(x_t)(x - x_t) = 0,$$

yielding the above update formula. The Newton step fails if $f'(x_t) = 0$ or gets out of control if $|f'(x_t)|$ is very small. Thus, we need to keep this in mind when making a theoretical analysis.

We can use this method for optimization as well, called *Newton's method*, where we can find critical points $f'(x) = 0$ by applying the method to the derivative of f ,

$$x_{t+1} \doteq x_t - \frac{f'(x_t)}{f''(x_t)}.$$

We can further generalize this to any dimensionality,

$$x_{t+1} \doteq x_t - \nabla^2 f(x_t)^{-1} \nabla f(x_t).$$

As before, we need to keep in mind that the Hessian must be invertible and may get out of control if the Hessian has small norm.

A second interpretation of Newton's method is that it is a special case of the general update scheme,

$$x_{t+1} \doteq x_t - H(x_t) \nabla f(x_t),$$

where $H(x_t) \in \mathbb{R}^{d \times d}$ is some matrix, like gradient descent with $H(x_t) = \gamma_t I$. Hence, we can think of Newton's method of an "adaptive gradient descent" that adapts to the local curvature of the function at x_t .⁹

We will not prove any general convergence guarantees for Newton's method. We will prove that, under suitable conditions, and starting close to a critical point, we will reach distance at most ϵ to this critical point in $\mathcal{O}(\log \log(1/\epsilon))$ steps. This also holds for non-convex functions. However, this is quite weak, since we assume that we are already close to the critical point. The proof will rely on the assumption that the local curvature in the small space around the critical point is near constant.

The *Babylonian method* to compute square roots is an application of the Newton-Raphson method. It finds zeros of $f(x) = x^2 - R$, which is equal to zero at \sqrt{R} and $-\sqrt{R}$. It takes $\mathcal{O}(\log R)$ steps to get within $1/2$ of a square root. Then, to get within ϵ , it takes $\log \log(1/\epsilon)$ steps. Thus, once we are close, we get very close very quickly.



Figure 22. A step of the Newton-Raphson method.

⁹ This is very apparent in the case of optimizing a quadratic function of the form $f(x) = \frac{1}{2}x^\top Mx - q^\top x + c$, which has the same curvature $\nabla^2 f(x) = M$ everywhere. In this case, Newton's method yields the optimum in a single step, $x_1 = x^*$.

Theorem 61. Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be twice differentiable with a critical point \mathbf{x}^* . Suppose there is a ball $\mathcal{X} \subseteq \text{dom}(f)$ with center \mathbf{x}^* such that the inverse Hessians are bounded,

$$\exists \mu > 0, \quad \left\| \nabla^2 f(\mathbf{x})^{-1} \right\| \leq \frac{1}{\mu}, \quad \forall \mathbf{x} \in \mathcal{X},$$

and the Hessian is Lipschitz continuous,

$$\exists B \geq 0, \quad \left\| \nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y}) \right\| \leq B \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}.$$

Then, for $\mathbf{x}_t \in \mathcal{X}$ and \mathbf{x}_{t+1} , resulting from the Newton step, we have

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\| \leq \frac{B}{2\mu} \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

Proof. Let $H(\mathbf{x}) = \nabla^2 f(\mathbf{x})$, $\mathbf{x} = \mathbf{x}_t$, $\mathbf{x}' = \mathbf{x}_{t+1}$. Then, subtracting \mathbf{x}^* from both sides of the Newton step yields

$$\begin{aligned} \mathbf{x}' - \mathbf{x}^* &= \mathbf{x} - \mathbf{x}^* - H(\mathbf{x})^{-1} \nabla f(\mathbf{x}) \\ &= \mathbf{x} - \mathbf{x}^* + H(\mathbf{x})^{-1} (\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x})). \end{aligned} \quad \nabla f(\mathbf{x}^*) = \mathbf{0}.$$

Let $h(t) \doteq \nabla f(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))$. Then, using the fundamental theorem of calculus, we get

$$\mathbf{x} - \mathbf{x}^* + H(\mathbf{x})^{-1} \int_0^1 H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) (\mathbf{x}^* - \mathbf{x}) dt \quad h'(t) = H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) (\mathbf{x}^* - \mathbf{x}).$$

Using $\mathbf{x} - \mathbf{x}^* = H(\mathbf{x})^{-1} H(\mathbf{x}) (\mathbf{x} - \mathbf{x}^*) = H(\mathbf{x})^{-1} \int_0^1 -H(\mathbf{x}) (\mathbf{x}^* - \mathbf{x}) dt$, we get

$$= H(\mathbf{x})^{-1} \int_0^1 (H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - H(\mathbf{x})) (\mathbf{x}^* - \mathbf{x}) dt$$

Taking norm of both sides yields

$$\begin{aligned} \|\mathbf{x}' - \mathbf{x}\| &= \left\| H(\mathbf{x})^{-1} \int_0^1 (H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - H(\mathbf{x})) (\mathbf{x}^* - \mathbf{x}) dt \right\| \\ &\leq \left\| H(\mathbf{x})^{-1} \right\| \cdot \left\| \int_0^1 (H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - H(\mathbf{x})) (\mathbf{x}^* - \mathbf{x}) dt \right\| \quad \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|. \\ &\leq \left\| H(\mathbf{x})^{-1} \right\| \cdot \int_0^1 \| (H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - H(\mathbf{x})) (\mathbf{x}^* - \mathbf{x}) \| dt \\ &\leq \left\| H(\mathbf{x})^{-1} \right\| \cdot \int_0^1 \| H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - H(\mathbf{x}) \| \cdot \|\mathbf{x}^* - \mathbf{x}\| dt \quad \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|. \\ &\leq \left\| H(\mathbf{x})^{-1} \right\| \cdot \|\mathbf{x}^* - \mathbf{x}\| \int_0^1 \| H(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - H(\mathbf{x}) \| dt \\ &\leq \frac{1}{\mu} \|\mathbf{x}^* - \mathbf{x}\| \int_0^1 B \|t(\mathbf{x}^* - \mathbf{x})\| dt \quad \text{Assumptions of theorem.} \\ &\leq \frac{B}{\mu} \|\mathbf{x}^* - \mathbf{x}\|^2 \int_0^1 t dt \\ &\leq \frac{B}{2\mu} \|\mathbf{x}^* - \mathbf{x}\|^2. \end{aligned}$$



Corollary. With the assumptions of Theorem 61, if $x_0 \in \mathcal{X}$ satisfies

$$\|x_0 - x^*\| \leq \frac{\mu}{B},$$

then Newton's method yields

$$\|x_T - x^*\| \leq \frac{\mu}{B} \left(\frac{1}{2}\right)^{2^T - 1}.$$

Hence, we get the $\mathcal{O}(\log \log(1/\epsilon))$ bound, but only if we are μ/B -close to x^* . Thus, we only converge fast to x^* if we are already close to it. For this to hold, it is of course necessary that x^* is the *only* close critical point to x_0 . However, this necessarily follows from the assumptions, since the Hessians are almost constant this close to x^* under the Lipschitz continuity and inverse Hessian bound. Thus, locally, the function behaves like a quadratic function, which converges to its unique critical point in one step.

9 Quasi-Newton methods

The problem with Newton's method is that it has a high computational complexity, due to the Hessian and inverse of it, which both have $\mathcal{O}(d^3)$ runtime complexity. *Quasi-Newton methods* are optimization methods that approximate the Hessian by a matrix $H_t \approx \nabla^2 f(x_t)$, which is a function of x_t , x_{t-1} , and H_{t-1} . We then iteratively update by

$$x_{t+1} \doteq x_t - H_t^{-1} \nabla f(x_t),$$

where $H_t \in \mathbb{R}^{d \times d}$ must be symmetric and satisfy the *secant condition*,

$$\nabla f(x_t) - \nabla f(x_{t-1}) = H_t(x_t - x_{t-1}).$$

In general, there are many matrices that satisfy these conditions. Thus, we must choose which H_t^{-1} to pick, based on x_{t-1} , x_t , and H_{t-1} .¹⁰

Recall from Newton's method that we wanted $\nabla^2 f(x_t)$ to fluctuate very little in regions of fast convergence. Hence, in Quasi-Newton methods, it makes sense if $H_{t-1}^{-1} \approx H_t^{-1}$. This intuition yields the approach by Greenstadt [1970], where we update H_{t-1}^{-1} by an error matrix E_t ,

$$H_t^{-1} = H_{t-1}^{-1} + E_t,$$

and we want this error to be as small as possible, *i.e.* minimize $\|E\|_F^2$, subject to its constraints. Greenstadt [1970] found this method “too specialized”, which lead him to minimize the following error term instead,

$$\|AEA^\top\|_F^2,$$

where $A \in \mathbb{R}^{d \times d}$ is a fixed invertible transformation matrix.

Let's now use the following notation to develop further algorithms,

$$\begin{aligned} H &\doteq H_{t-1}^{-1} \\ H' &\doteq H_t^{-1} \\ E &\doteq E_t \\ \sigma &\doteq x_t - x_{t-1} \\ y &\doteq \nabla f(x_t) - \nabla f(x_{t-1}) \\ r &\doteq \sigma - Hy. \end{aligned}$$

We then have the following convex constrained minimization problem in d^2 variables,

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|AEA^\top\|_F^2 \\ \text{subject to} \quad & Ey = r \\ & E^\top - E = 0, \end{aligned}$$

where the first condition is the secant condition,

$$H'y = \sigma \Leftrightarrow (H + E)y = \sigma \Leftrightarrow Ey = \sigma - Hy \Leftrightarrow Ey = r,$$

¹⁰ We directly work with H_t^{-1} , instead of H_t , since computing the inverse would again result in a $\mathcal{O}(d^3)$ runtime complexity.

and the second condition ensures symmetry, since if H_{t-1}^{-1} and E_t are symmetric, then H_t^{-1} is as well.

Let

$$f(E) = \frac{1}{2} \|AEA^\top\|_F^2.$$

Because the conditions are all linear, we can summarize them in one equation as $CE = B$ for some matrices C and B . Furthermore, due to this convex program only having equality constraints, the Slater point condition for strong duality becomes void. Thus, we obtain strong duality “for free”. Thus, the Karush-Kuhn-Tucker conditions hold, which imply there exists a vector $\lambda \in \mathbb{R}^m$ such that

$$\nabla f(E^*)^\top = \lambda^\top C.$$

Directly follows from the vanishing gradient condition in KKT.

Let $W = A^\top A$ and $M = W^{-1}$, then the gradient of f can be computed by

$$\nabla f(E) = A^\top AEA^\top A = WEW = M^{-1}EM^{-1}.$$

Now, since the objective is quadratic, we can obtain the minimizer x^* and the Lagrange multipliers λ by solving the following system of linear equations,

$$\begin{aligned} CE &= B \\ E &= M^\top \lambda^\top CM^\top. \end{aligned}$$

Solving this system yields

$$\begin{aligned} E^* &= \frac{1}{y^\top My} \left(\sigma y^\top M + My\sigma^\top - Hyy^\top M - Myy^\top H \right. \\ &\quad \left. - \frac{1}{y^\top My} (y^\top \sigma - y^\top Hy) Myy^\top M \right). \end{aligned}$$

This is called the *Greenstadt method* with parameter M .

Now, we need to decide which M to use. Greenstadt [1970] suggested $M = I$ and $M = H \doteq H_{t-1}^{-1}$. Goldfarb [1970] suggested $M = H' \doteq H_t^{-1}$. Because of the secant condition, we get the following,

$$My = H'y = \sigma.$$

Hence, despite not knowing this value yet, we can still use it, since it will cancel out all terms, containing $M = H'$. This is called the *BFGS method*, and the optimal error matrix becomes

$$E^* = \frac{1}{y^\top \sigma} \left(-Hy\sigma^\top - \sigma y^\top H + \left(1 + \frac{y^\top Hy}{y^\top \sigma} \right) \sigma \sigma^\top \right).$$

With this error matrix, we get the following update,

$$H' = \left(I - \frac{\sigma y^\top}{y^\top \sigma} \right) H \left(I - \frac{y \sigma^\top}{y^\top \sigma} \right) + \frac{\sigma \sigma^\top}{y^\top \sigma}.$$

The cost per step of this algorithm is $\mathcal{O}(d^2)$, which is a big upgrade over $\mathcal{O}(d^3)$ that we had for Newton's method. However, we can make it even faster by making another approximation, which will yield the *L-BFGS* algorithm.

Recall the Quasi-Newton update step,

$$\mathbf{x}_{t+1} \doteq \mathbf{x}_t - \mathbf{H}_t^{-1} \nabla f(\mathbf{x}_t).$$

Observe that we do not necessarily need the $d \times d$ matrix \mathbf{H}_t^{-1} ; we only need the d -dimensional vector $\mathbf{H}_t^{-1} \nabla f(\mathbf{x}_t)$. Let $\mathbf{g}' \in \mathbb{R}^d$. Suppose that we have an oracle to compute $\mathbf{s} = \mathbf{H}\mathbf{g}$ for any vector \mathbf{g} , then we can compute $\mathbf{s}' = \mathbf{H}'\mathbf{g}'$ with one oracle call and $\mathcal{O}(d)$ additional operations, assuming that \mathbf{y} and σ are known.

We can implement the oracle recursively,

$$\begin{aligned}\sigma_k &\doteq \mathbf{x}_k - \mathbf{x}_{k-1} \\ \mathbf{y}_k &\doteq \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k+1}).\end{aligned}$$

This allows us to compute the BFGS-step $\mathbf{H}_t^{-1} \nabla f(\mathbf{x}_t)$ recursively. However, this would result in $\mathcal{O}(td)$ runtime complexity per step, since we would have to go down all steps, and generally $t > d$. Thus, we have a worse algorithm if we want to compute the next vector exactly. But, if we only go down m steps of recursion for some small m , we get $\mathcal{O}(md)$ complexity, which is linear if m is constant; see Algorithm 1. Intuitively, this should give a good approximation, since the earlier steps should not be so relevant anymore, since we are likely in a different landscape at the current timestep.

```

function LBFGSSTEP( $k, \ell, \mathbf{g}'$ )
  if  $\ell = 0$  then
    return  $\mathbf{H}_0^{-1} \mathbf{g}'$ 
  end if
   $\mathbf{h} = \sigma \frac{\sigma_k^\top \mathbf{g}'}{\mathbf{y}_k^\top \sigma_k}$ 
   $\mathbf{g} = \mathbf{g}' - \mathbf{y} \frac{\sigma_k^\top \mathbf{g}'}{\mathbf{y}_k^\top \sigma_k}$ 
   $\mathbf{s} = \text{LBFGSSTEP}(k-1, \ell-1, \mathbf{g})$ 
   $\mathbf{w} = \mathbf{s} - \sigma_k \frac{\mathbf{y}_k^\top \mathbf{s}}{\mathbf{y}_k^\top \sigma_k}$ 
   $\mathbf{z} = \mathbf{w} + \mathbf{h}$ 
  return  $\mathbf{z}$ 
end function

```

Algorithm 1. The L-BFGS algorithm. The outer products can be computed as inner products, giving $\mathcal{O}(d)$ runtime complexity to all the products.

10 Subgradient methods

Until now, we have mostly assumed all functions to be smooth and hence differentiable. However, in general this is not the case. In machine learning, non-smoothness arises everywhere:

- Loss functions, such as the Hinge loss, $\max\{0, 1 - x\}$ (SVM);
- Regularization, such as the ℓ_1 -norm (LASSO);
- Activation functions, such as ReLU.

This motivates the need for a more general notion of the gradient that can be applied to more functions.

Definition 62 (Subgradient). $g \in \mathbb{R}^d$ is a subgradient of f at x if

$$f(y) \geq f(x) + g^\top(y - x), \quad \forall y \in \text{dom}(f).$$

We call $\partial f(x) \subseteq \mathbb{R}^d$ the subdifferential, which is the set of subgradients of f at x .

Example 63. Consider $f(x) = |x|$, then $\partial f(0) = [-1, 1]$.

Lemma 64. If f is differentiable at $x \in \text{dom}(f)$, then $\partial f(x) \subseteq \{\nabla f(x)\}$.

Lemma 64 means that if f is differentiable at x , then this is either the only subgradient or there is no subgradient at all. There might not be any subgradient at all in this case because it might be that the hyperplane is not below the entire function if the function is non-convex.

Lemma 65. The following characterizes convexity with the subgradient:

- If f is convex, then $\partial f(x) \neq \emptyset$ for all x in the relative interior of $\text{dom}(f)$;
- If $\text{dom}(f)$ is convex and $\partial f(x) \neq \emptyset$ for all $x \in \text{dom}(f)$, then f is convex.

Lemma 66 (Subgradient optimality condition). Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ and $x \in \text{dom}(f)$. If $0 \in \partial f(x)$, then x is a global minimum.

Proof. By definition of the subgradient with $g = 0 \in \partial f(x)$ gives

$$f(y) \geq f(x) + g^\top(y - x) = f(x), \quad \forall y \in \text{dom}(f).$$

Thus, x is a global minimum. ■



Figure 23. g is a subgradient of f at x if the whole graph is above x 's supporting hyperplane, parametrized by g .

Lemma 67 (Subgradient calculus). We can use the following operations to work with subgradients:

- Conic combination: Let $h(\mathbf{x}) = \alpha f(\mathbf{x}) + \beta g(\mathbf{x})$ with $\alpha, \beta \geq 0$, then

$$\partial h(\mathbf{x}) = \alpha \partial f(\mathbf{x}) + \beta \partial g(\mathbf{x});$$

- Affine transformation: Let $h(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$, then

$$\partial h(\mathbf{x}) = \mathbf{A}^\top \partial f(\mathbf{A}\mathbf{x} + \mathbf{b});$$

- Pointwise maximum: Let $h(\mathbf{x}) = \max_{i \in [m]} f_i(\mathbf{x})$, then

$$\partial h(\mathbf{x}) = \text{conv}(\{\partial f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = h(\mathbf{x})\}).$$

Thus, at each point where we transition from one function to another, we get the convex hull of subgradients of the functions that transition. At all other points, we take the maximum function's subgradient.

10.1 Subgradient method

In the subgradient method, the general update rule becomes

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}(\mathbf{x}_t - \gamma_t \mathbf{g}_t), \quad \mathbf{g}_t \in \partial f(\mathbf{x}_t),$$

which can be rewritten as an optimization problem,

$$\begin{aligned} &= \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_t - \gamma_t \mathbf{g}_t)\|^2 \\ &= \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_t - (-\gamma_t \mathbf{g}_t)\|^2 \\ &= \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \left(\|\mathbf{x} - \mathbf{x}_t\|^2 + \|\gamma_t \mathbf{g}_t\|^2 + 2\langle \gamma_t \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle \right) \\ &= \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 + \langle \gamma_t \mathbf{g}_t, \mathbf{x} \rangle \right\}. \end{aligned}$$

Cosine theorem.

Remove terms that do not depend on \mathbf{x} .

If f is differentiable, gradient descent and projected gradient descent are special cases of this update rule, where $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{X} \subset \mathbb{R}^d$, respectively. However, if f is non-differentiable, we will see that this is technically not a descent method, because the subgradient is not a descent direction in general.

Lemma 68 (Subgradient method “descent” lemma). If f is convex, then for any optimal solution \mathbf{x}^* ,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) + \gamma_t^2 \|\mathbf{g}_t\|^2.$$

Proof.

$$\begin{aligned}
\|x_{t+1} - x^*\|^2 &= \|\Pi_{\mathcal{X}}(x_t - \gamma_t g_t) - x^*\|^2 \\
&= \|x_t - \gamma_t g_t - x^*\|^2 \\
&= \|x_t - x^*\|^2 - 2\gamma_t g_t^\top (x_t - x^*) + \gamma_t^2 \|g_t\|^2 \\
&\leq \|x_t - x^*\|^2 - 2\gamma_t (f(x_t) - f(x^*)) + \gamma_t^2 \|g_t\|^2.
\end{aligned}$$

■

Theorem 69 (Convergence of the subgradient method). If f is convex, then the subgradient method satisfies

$$\min_{t \in [T]} f(x_t) - f(x^*) \leq \frac{\|x_0 - x^*\|^2 + \sum_{t=0}^{T-1} \gamma_t^2 \|g_t\|^2}{2 \sum_{t=0}^{T-1} \gamma_t}.$$

Proof. This can easily be shown by telescoping sum of Lemma 68 and invoking convexity. ■

Assuming bounded subgradient $\|g_t\| \leq B$ for all steps t , we get the following convergence rates under various stepsizes,

- Constant stepsize ($\gamma_t = \gamma$):

$$\lim_{t \rightarrow \infty} f(x_t^{\text{best}}) \leq f(x^*) + \frac{B^2 \gamma}{2};$$

- Scaled stepsize ($\gamma_t = \gamma / \|g_t\|$):

$$\lim_{t \rightarrow \infty} f(x_t^{\text{best}}) \leq f(x^*) + \frac{B \gamma}{2};$$

- Square-summable stepsize ($\sum_{t=0}^{\infty} \gamma_t^2 < +\infty$, $\sum_{t=0}^{\infty} \gamma_t = +\infty$):

$$\lim_{t \rightarrow \infty} f(x_t^{\text{best}}) = f(x^*);$$

- Diminishing stepsize ($\gamma_t \rightarrow 0$ and $\sum_{t=0}^{\infty} \gamma_t = +\infty$):

$$\lim_{t \rightarrow \infty} f(x_t^{\text{best}}) = f(x^*).$$

Corollary. Let f be convex and B -Lipschitz continuous. Let \mathcal{X} be convex compact with $R^2 = \max_{x, y \in \mathcal{X}} \|x - y\|_2^2 < +\infty$. Setting

$$\gamma \doteq \frac{R}{B\sqrt{T}},$$

then the subgradient method satisfies

$$\min_{t \in [T]} f(x_t) - f(x^*) \leq \frac{BR}{\sqrt{T}}.$$

Subgradient descent update rule.

\mathcal{X} is convex, so we get closer to x^* after projection.

Cosine theorem: $\|x - y\|^2 = \|x\|^2 + \|y\|^2 - 2x^\top y$.

Convexity: $f(x^*) \geq f(x_t) + g_t^\top (x_t - x^*)$.

To achieve ϵ -optimality, the subgradient method requires $\mathcal{O}\left(\frac{B^2 R^2}{\epsilon^2}\right)$ iterations.

10.2 Strong convexity

Theorem 70. Let f be μ -strongly convex and B -Lipschitz continuous on \mathcal{X} . Setting

$$\gamma_t \doteq \frac{2}{\mu(t+1)},$$

then the subgradient method satisfies

$$\min_{t \in [T]} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2B^2}{\mu(T+1)}.$$

In the case of strong convexity, to achieve ϵ -optimality, the subgradient method requires $\mathcal{O}\left(\frac{B^2}{\mu\epsilon}\right)$ iterations.

Proof. Adapting the proof of Lemma 68 to use strong convexity in its last step, we get

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq (1 - \mu\gamma_t)\|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) + \gamma_t^2\|\mathbf{g}_t\|^2.$$

Using this, we get the following,

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{1 - \mu\gamma_t}{2\gamma_t}\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{2\gamma_t}\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{\gamma_t}{2}\|\mathbf{g}_t\|^2 \\ &= \frac{\mu(t-1)}{4}\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{\mu(t+1)}{4}\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{1}{\mu(t+1)}\|\mathbf{g}_t\|^2. \gamma_t \doteq 2/\mu(t+1) \end{aligned}$$

Now, it is easy to show the result by a telescoping sum. ■

11 Mirror descent

Like in subgradient descent, we continue to assume that f is non-smooth. In practice, we often have additional information about set \mathcal{X} that we might be able to exploit. Specifically, we will explore how we can exploit non-Euclidean geometry of a convex set \mathcal{X} .¹¹

¹¹ Until this point, we have only made use of Euclidean geometry by way of using the $\|\cdot\|_2$ -norm.

11.1 Norm and Bregman divergence

Definition 71 (Norm). A function $\|\cdot\| : \mathcal{X} \rightarrow \mathbb{R}_+$ is a norm if it satisfies the following properties,

1. (Positive definiteness) $\|x\| = 0$ if and only if $x = \mathbf{0}$;
2. (Positive homogeneity) $\|\alpha x\| = |\alpha| \|x\|$;
3. (Subadditivity): $\|x + y\| \leq \|x\| + \|y\|$.

Definition 72 (Dual norm). The dual norm $\|\cdot\|_*$ of a norm $\|\cdot\|$ satisfies the properties of a norm and

$$\|y\|_* \doteq \max_{\|x\| \leq 1} \langle x, y \rangle.$$

Example 73. For $p \geq 1$ and $1/p + 1/q = 1$, we have the following norms with their dual norms,

$$\|x\|_p \doteq \left(\sum_{i=1}^d |x_i|^p \right)^{1/p}, \quad \|\cdot\|_{p,*} = \|\cdot\|_q.$$

Lemma 74.

$$\frac{1}{\sqrt{d}} \|x\|_2 \leq \|x\|_\infty \leq \|x\|_2 \leq \|x\|_1 \leq \sqrt{d} \|x\|_2.$$

The nice thing about smoothness, Lipschitz continuity, and strong convexity is that they can be defined for any norm.

Definition 75 (Bregman divergence). Let $\omega : \Omega \rightarrow \mathbb{R}$ be continuously differentiable on Ω and 1-strongly convex w.r.t. some norm $\|\cdot\|$,

$$\omega(x) \geq \omega(y) + \nabla \omega(y)^\top (x - y) + \frac{1}{2} \|x - y\|^2, \quad \forall x, y \in \Omega.$$

The Bregman divergence V_ω is defined as

$$V_\omega(x, y) \doteq \omega(x) - \omega(y) - \nabla \omega(y)^\top (x - y), \quad \forall x, y \in \Omega.$$

Note that the Bregman divergence V_ω is thus the first-order Taylor approximation of ω , evaluated at x .

Example 76. We have the following examples of Bregman divergences,

1. (Euclidean distance) $\Omega = \mathbb{R}^d$, $\omega(x) = \frac{1}{2}\|x\|_2^2$, and $\|\cdot\| = \|\cdot\|_2$. Then

$$V_\omega(x, y) = \frac{1}{2}\|x - y\|_2^2.$$

2. (Mahalanobis distance) $\Omega = \mathbb{R}^d$, $\omega(x) = \frac{1}{2}x^\top Qx$ with $Q \succeq I$, and $\|\cdot\| = \|\cdot\|_2$. Then,

$$V_\omega(x, y) = \frac{1}{2}(x - y)^\top Q(x - y).$$

3. (Kullback-Leibler divergence) $\Omega = \Delta^{d-1}$, $\omega(x) = \sum_{i=1}^d x_i \log x_i$, and $\|\cdot\| = \|\cdot\|_1$. Then,

$$V_\omega(x, y) = \text{KL}(x; y) \doteq \sum_{i=1}^d x_i \log \frac{x_i}{y_i}.$$

Lemma 77. Any Bregman divergence satisfies the following properties:

1. (Non-negativity) $V_\omega(x, y) \geq 0$;
2. (Convexity) $V_\omega(x, y)$ is convex in x ;
3. (Positivity) $V_\omega(x, y) = 0$ if and only if $x = y$;
4. $V_\omega(x, y) \geq \frac{1}{2}\|x - y\|^2$.

The following Lemma is a key property of the Bregman divergence and is used extensively in this course.

Lemma 78 (Three-point identity). $\forall x, y, z \in \Omega$:

$$V_\omega(x, z) = V_\omega(x, y) + V_\omega(y, z) - \langle \nabla \omega(z) - \nabla \omega(y), x - y \rangle.$$

11.2 Mirror descent algorithm

The mirror descent algorithm is a generalization of the subgradient method, where the norm is replaced by a Bregman divergence,¹²

$$x_{t+1} = \underset{x \in \mathcal{X}}{\operatorname{argmin}} \{V_\omega(x, x_t) + \langle \gamma_t g_t, x \rangle\}, \quad g_t \in \partial f(x_t).$$

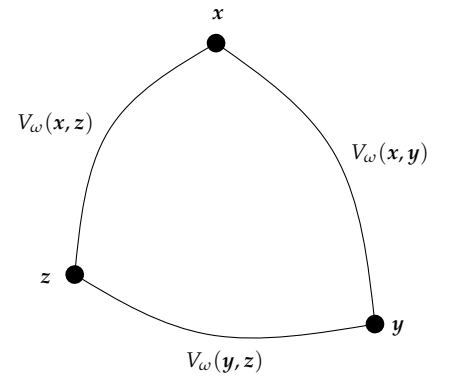


Figure 24. Illustration of the three-point identity of a non-Euclidean Bregman divergence.

¹² Using Item 1 of Example 76 recovers the subgradient method, and Item 3 recovers Entropic descent,

$$x_{t+1} \propto x_t \odot \exp(-\gamma_t g_t).$$

Lemma 79. Let f be convex and ω be 1-strongly convex on \mathcal{X} w.r.t. norm $\|\cdot\|$. Then, the following inequality holds,

$$\gamma_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq V_\omega(\mathbf{x}^*, \mathbf{x}_t) - V_\omega(\mathbf{x}^*, \mathbf{x}_{t+1}) + \frac{\gamma_t^2}{2} \|\mathbf{g}_t\|_*^2.$$

Proof. We have the following update rule,

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \{V_\omega(\mathbf{x}, \mathbf{x}_t) + \langle \gamma_t \mathbf{g}_t, \mathbf{x} \rangle\}.$$

Thus, by the optimality condition, we have

$$\langle \nabla \omega(\mathbf{x}_{t+1}) + \gamma_t \mathbf{g}_t - \nabla \omega(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_{t+1} \rangle \geq 0, \quad \forall \mathbf{x} \in \mathcal{X},$$

which can be equivalently written as $\forall \mathbf{x} \in \mathcal{X}$:

$$\begin{aligned} \langle \gamma_t \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{x} \rangle &\leq \langle \nabla \omega(\mathbf{x}_{t+1}) - \nabla \omega(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_{t+1} \rangle \\ &= V_\omega(\mathbf{x}, \mathbf{x}_t) - V_\omega(\mathbf{x}, \mathbf{x}_{t+1}) - V_\omega(\mathbf{x}_{t+1}, \mathbf{x}_t) \\ &\leq V_\omega(\mathbf{x}, \mathbf{x}_t) - V_\omega(\mathbf{x}, \mathbf{x}_{t+1}) - \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2. \end{aligned}$$

Three-point identity.

Fourth property of Bregman divergence.

As a result,

$$\begin{aligned} \gamma_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \langle \gamma_t \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \\ &= \langle \gamma_t \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{x}^* \rangle + \langle \gamma_t \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_{t+1} \rangle \\ &\leq V_\omega(\mathbf{x}^*, \mathbf{x}_t) - V_\omega(\mathbf{x}^*, \mathbf{x}_{t+1}) - \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &\quad + \langle \gamma_t \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_{t+1} \rangle \\ &\leq V_\omega(\mathbf{x}^*, \mathbf{x}_t) - V_\omega(\mathbf{x}^*, \mathbf{x}_{t+1}) - \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &\quad + \frac{1}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \frac{1}{2} \|\gamma_t \mathbf{g}_t\|_*^2 \\ &\leq V_\omega(\mathbf{x}^*, \mathbf{x}_t) - V_\omega(\mathbf{x}^*, \mathbf{x}_{t+1}) + \frac{\gamma_t^2}{2} \|\mathbf{g}_t\|_*^2. \end{aligned}$$

By definition of the subgradient.

Young's inequality.

■

Theorem 80 (Convergence of mirror descent). Let f be convex and ω be 1-strongly convex on \mathcal{X} w.r.t. norm $\|\cdot\|$. Then, mirror descent satisfies

$$\min_{t \in [T]} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{V_\omega(\mathbf{x}^*, \mathbf{x}_0) + \frac{1}{2} \sum_{t=0}^{T-1} \gamma_t^2 \|\mathbf{g}_t\|_*^2}{\sum_{t=0}^{T-1} \gamma_t}.$$

Note that this generalizes the convergence result of the subgradient method.

Suppose f is B -Lipschitz continuous such that $|f(\mathbf{x}) - f(\mathbf{y})| \leq B\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$. Namely, we then have $\|\mathbf{g}\|_* \leq B, \forall \mathbf{g} \in \partial f(\mathbf{x}), \mathbf{x} \in \mathcal{X}$.

Furthermore, let $R^2 \doteq \sup_{\mathbf{x} \in \mathcal{X}} V_\omega(\mathbf{x}, \mathbf{x}_0)$ and set

$$\gamma_t \doteq \frac{\sqrt{2}R}{B\sqrt{T}}.$$

Then, we have the following convergence rate,

$$\min_{t \in [T]} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{BR}{\sqrt{T}}\right).$$

This is equivalent to the convergence rate of the subgradient method, but then for a more general notions of norm.

In practice, if we optimize over the simplex Δ^{d-1} with $\|\mathbf{g}\|_\infty \leq 1, \forall \mathbf{g} \in \partial f(\mathbf{x})$ and $\mathbf{x}_0 = [1/d, \dots, 1/d]$. Then, we have the following convergence rate for the subgradient method, $\mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{T}}\right)$, because $B \in \mathcal{O}(\sqrt{d})$ and $R \in \mathcal{O}(1)$. On the other hand, we have the following convergence rate for mirror descent, $\mathcal{O}\left(\frac{\sqrt{\log d}}{\sqrt{T}}\right)$, since $B \in \mathcal{O}(1)$ and $R \in \mathcal{O}(\sqrt{\log d})$. This is a considerable speedup.

12 Smoothing and proximal algorithms

Often, we want to optimize non-smooth functions. However, most of the time, we assume functions to be smooth. The question is thus whether we can exploit additional structure of non-smooth functions, instead of treating them as black boxes. The idea behind smoothing is to optimize a smooth approximation f_μ of the non-smooth function f .

12.1 Nesterov smoothing

Definition 81 (Conjugate function). The conjugate function of f is

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom}(f)} \left\{ \mathbf{x}^\top \mathbf{y} - f(\mathbf{x}) \right\}.$$

It is also called the Legendre-Fenchel transformation.

Lemma 82 (Convex conjugate properties). The following holds for conjugate functions,

1. (Duality) If f is continuous and convex, then $f^{**} = f$;
2. (Fenchel's inequality)

$$f(\mathbf{x}) + f^*(\mathbf{y}) \geq \mathbf{x}^\top \mathbf{y}, \quad \forall \mathbf{x}, \mathbf{y};$$

3. If f and g are continuous and convex, then

$$(f + g)^*(\mathbf{x}) = \inf_{\mathbf{y}} \{f^*(\mathbf{y}) + g^*(\mathbf{x} - \mathbf{y})\};$$

4. If f is μ -strongly convex, then f^* is differentiable and $1/\mu$ -smooth.

Nesterov smoothing approximates a non-smooth function f by

$$f_\mu(\mathbf{x}) = \max_{\mathbf{y} \in \text{dom}(f^*)} \left\{ \mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y}) - \mu \cdot d(\mathbf{y}) \right\},$$

where $d(\mathbf{y})$ is a proximity function. A proximity function is 1-strongly convex and non-negative. The function f_μ is $1/\mu$ -smooth and approximates a convex f by

$$f(\mathbf{x}) - \mu D^2 \leq f_\mu(\mathbf{x}) \leq f(\mathbf{x}), \quad D^2 \doteq \max_{\mathbf{y} \in \text{dom}(f^*)} d(\mathbf{y}).$$

High μ results in a bad approximation.

Thus, we have a trade-off between approximation error and optimization efficiency. Specifically,

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \underbrace{f(\mathbf{x}) - f_\mu(\mathbf{x})}_{\text{approximation error}} + \underbrace{f_\mu(\mathbf{x}) - \min_{\mathbf{x}} f_\mu(\mathbf{x})}_{\text{optimization error}}.$$

The approximation error is on the order $\mathcal{O}(\mu)$, while the optimization error is on the order $\mathcal{O}(1/\mu t)$ using gradient descent.

If we apply accelerated gradient descent to solve the smoothed problem, we get an error of the following order,

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathcal{O}\left(\mu D^2 + \frac{R^2}{\mu t^2}\right).$$

Note that this is faster than applying subgradient methods.

12.2 Moreau-Yosida smoothing

Moreau-Yosida regularization smooths f by

$$f_\mu(\mathbf{x}) = \min_{\mathbf{y} \in \text{dom}(f^*)} \left\{ f(\mathbf{y}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{y}\|_2^2 \right\}.$$

This function is called the *Moreau envelope* of $f(\mathbf{x})$. For example, the Huber function is the Moreau envelope of $f(x) = |x|$,

$$f_\mu(x) = \begin{cases} \frac{x^2}{2\mu}, & |x| \leq \mu \\ |x| - \frac{\mu}{2}, & |x| > \mu. \end{cases}$$

As in Nesterov smoothing, f_μ is $1/\mu$ -smooth. However, the advantage is that it minimizes exactly, i.e., $\min_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{x}} f_\mu(\mathbf{x})$.

12.3 Proximal operators

Definition 83 (Proximal operator). The proximal operator of a convex function f at \mathbf{x} is defined as

$$\text{prox}_f(\mathbf{x}) \doteq \underset{\mathbf{y} \in \text{dom}(f)}{\text{argmin}} \left\{ f(\mathbf{y}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \right\}.$$

Note that this is a special case of Moreau-Yosida regularization with $\mu = 1$. Furthermore, for many non-smooth functions, their proximal operator can be computed efficiently in a closed form.

12.4 Proximal point algorithm

The proximal point algorithm (PPA) repeatedly applies the proximal operator,

$$\mathbf{x}_{t+1} = \text{prox}_{\lambda_t f}(\mathbf{x}_t).$$

Theorem 84 (Convergence of PPA). If f is convex, then for any $T \geq 1$, we have

$$f(\mathbf{x}_{T+1}) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2 \sum_{t=1}^T \lambda_t}.$$

If we set $\lambda_t = \lambda$ to be constant, we get a $\mathcal{O}(1/t)$ convergence rate.

12.5 Proximal gradient method

Assume the following convex composite optimization problem,

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

where f and g are convex.¹³ The proximal gradient method (PGM) has the following update rule,

$$\mathbf{x}_{t+1} = \text{prox}_{\gamma_t g}(\mathbf{x}_t - \gamma_t \nabla f(\mathbf{x}_t)).$$

Note that it alternates between a gradient update on f and a proximal operator on g .

Theorem 85 (Convergence of PGM). Let $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$. Assume f is convex and L -smooth, g is convex and possibly non-smooth. Proximal gradient method with fixed stepsize $\gamma_t = 1/L$ satisfies

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

This is nearly the same convergence rate as gradient descent, despite F being possibly non-smooth.

¹³ Most supervised learning problems can be cast into this form,

$$\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \ell(h_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i) + g(\boldsymbol{\theta}),$$

where $h_{\boldsymbol{\theta}}$ is the predictor and g is a regularization function.

13 Stochastic optimization

Stochastic optimization involves decision-making in the presence of randomness. The optimization problem is formalized by a random vector $\xi \sim P$,

$$\min_{x \in \mathbb{R}^d} F(x) \doteq \mathbb{E}_{\xi}[f(x, \xi)].$$

P is unknown and can only be accessed through data, which make F and ∇F hard to compute. Furthermore, we assume unbiasedness of the stochastic gradient,

$$\mathbb{E}_{\xi_t}[\nabla f(x_t, \xi_t) \mid x_t] = \nabla F(x_t).$$

In this setting, we use stochastic gradient descent (SGD), which has the following update rule,

$$\begin{aligned} \xi_t &\stackrel{\text{iid}}{\sim} P \\ x_{t+1} &= x_t - \gamma_t \nabla f(x_t, \xi_t). \end{aligned}$$

13.1 Convergence analysis

In the non-convex case, we can show that SGD finds a stationary point with $\mathbb{E}\|\nabla F(\bar{x})\| \leq \epsilon$ in $\mathcal{O}(1/\epsilon^4)$ gradient evaluations.

Theorem 86 (Non-convex, random output). Suppose F is L -smooth and the stochastic gradient has bounded variance, i.e., $\mathbb{E}\|\nabla f(x, \xi) - \nabla F(x)\|_2^2 \leq \sigma^2$. Then, SGD with $\gamma \doteq \min\{1/L, \gamma_0/\sigma\sqrt{T}\}$ achieves

$$\begin{aligned} \mathbb{E}\|\nabla F(\hat{x}_T)\|^2 &\leq \frac{\sigma}{\sqrt{T}} \left(\frac{2(F(x_1) - F(x^*))}{\gamma_0} + L\gamma_0 \right) \\ &\quad + \frac{2L(F(x_1) - F(x^*))}{T}, \end{aligned}$$

where $\hat{x}_T \sim \text{Unif}(\{x_1, \dots, x_T\})$.

Proof.

$$\begin{aligned} \mathbb{E}[F(x_{t+1}) - F(x_t)] &\leq \mathbb{E}\left[\nabla F(x_t)^\top (x_{t+1} - x_t) + \frac{L}{2}\|x_{t+1} - x_t\|^2\right] && \text{Smoothness of } F. \\ &= \mathbb{E}\left[-\gamma_t \nabla F(x_t)^\top \nabla f(x_t, \xi_t) + \frac{L\gamma_t^2}{2}\|\nabla f(x_t, \xi_t)\|^2\right] && \text{SGD update rule.} \\ &= -(\gamma_t - \frac{L\gamma_t^2}{2})\mathbb{E}\|\nabla F(x_t)\|^2 + \frac{L\sigma^2\gamma_t^2}{2} && \mathbb{E}[X^2] = \mathbb{E}[X]^2 + \mathbb{V}[X]. \\ &\leq -\frac{\gamma_t}{2}\mathbb{E}\|\nabla F(x_t)\|^2 + \frac{L\sigma^2\gamma_t^2}{2}. && \gamma_t \leq 1/L. \end{aligned}$$

We can rewrite this as

$$\mathbb{E}\|\nabla F(x_t)\|^2 \leq \frac{2\mathbb{E}[F(x_t) - F(x_{t+1})]}{\gamma_t} + \gamma_t\sigma^2L.$$

By definition of $\hat{\mathbf{x}}_T$, we have

$$\begin{aligned}
\mathbb{E} \|\nabla F(\hat{\mathbf{x}}_T)\|^2 &= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla F(\mathbf{x}_t)\|^2 \\
&\leq \frac{1}{T} \left(\sum_{t=1}^T \frac{2\mathbb{E}[F(\mathbf{x}_t) - F(\mathbf{x}_{t+1})]}{\gamma_t} + \gamma_t \sigma^2 L \right) \\
&= \frac{2}{\gamma T} \left(\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}_{t+1}) \right) + \gamma \sigma^2 L && \text{Constant stepsize.} \\
&= \frac{2(F(\mathbf{x}_1) - F(\mathbf{x}_{T+1}))}{\gamma T} + \gamma \sigma^2 L && \text{Telescoping sum.} \\
&\leq \frac{2(F(\mathbf{x}_1) - F(\mathbf{x}^*))}{\gamma T} + \gamma \sigma^2 L \\
&\leq \frac{2(F(\mathbf{x}_1) - F(\mathbf{x}^*))}{T} \max \left\{ L, \frac{\sigma \sqrt{T}}{\gamma_0} \right\} + \frac{\gamma_0 \sigma L}{\sqrt{T}} \\
&\leq \frac{2L(F(\mathbf{x}_1) - F(\mathbf{x}^*))}{T} + \frac{\sigma}{\sqrt{T}} \left(\frac{2(F(\mathbf{x}_1) - F(\mathbf{x}^*))}{\gamma_0} + L \gamma_0 \right).
\end{aligned}$$

■

In the convex case, we can show that SGD finds an ϵ -optimal solution with $\mathcal{O}(1/\epsilon^2)$ sample complexity.

Theorem 87 (Convex, weighted averaging). Suppose F is convex and $\mathbb{E} \|\nabla f(\mathbf{x}, \xi)\|^2 \leq B^2, \forall \mathbf{x}$. Then, SGD satisfies

$$\mathbb{E}[F(\hat{\mathbf{x}}_T) - F(\mathbf{x}^*)] \leq \frac{R^2 + B^2 \sum_{t=0}^T \gamma_t^2}{2 \sum_{t=0}^T \gamma_t},$$

where

$$\hat{\mathbf{x}}_T \doteq \frac{\sum_{t=0}^T \gamma_t \mathbf{x}_t}{\sum_{t=0}^T \gamma_t}, \quad \|\mathbf{x}_0 - \mathbf{x}^*\| \leq R.$$

Proof. First, we have

$$\begin{aligned}
\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_t - \gamma_t \nabla f(\mathbf{x}_t, \xi_t) - \mathbf{x}^*\|^2 \\
&= \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \gamma_t^2 \|\nabla f(\mathbf{x}_t, \xi_t)\|^2 - 2 \nabla f(\mathbf{x}_t, \xi_t)^\top (\mathbf{x}_t - \mathbf{x}^*).
\end{aligned}$$

Cosine theorem.

Furthermore, by the law of total expectation ($\mathbb{E}[X] = \mathbb{E}_Y[\mathbb{E}_X[X | Y]]$),

$$\begin{aligned}
\mathbb{E}_{\xi_{1:t}} \left[\nabla f(\mathbf{x}_t, \xi_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \right] &= \mathbb{E}_{\xi_{1:t-1}} \left[\mathbb{E}_{\xi_t} \left[\nabla f(\mathbf{x}_t, \xi_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \mid \mathbf{x}_t \right] \right] \\
&= \mathbb{E}_{\xi_{1:t-1}} \left[\mathbb{E}_{\xi_t} \left[\nabla f(\mathbf{x}_t, \xi_t) \mid \mathbf{x}_t \right]^\top (\mathbf{x}_t - \mathbf{x}^*) \right] && \mathbf{x}_t \text{ can be computed from } \xi_{1:t-1}. \\
&= \mathbb{E}_{\xi_{1:t-1}} \left[\nabla F(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \right] \\
&\geq \mathbb{E}[F(\mathbf{x}_t) - F(\mathbf{x}^*)]. && \text{Convexity of } F.
\end{aligned}$$

This gives us the following recursion,

$$\gamma_t \mathbb{E}[F(\mathbf{x}_t) - F(\mathbf{x}^*)] \leq \frac{1}{2} \mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{2} \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{1}{2} \gamma_t^2 B^2,$$

and the result follows by telescoping the sum from $t = 1$ to T . ■

In the strongly convex case, we can show that SGD finds an ϵ -optimal solution with $\mathcal{O}(1/\epsilon)$ complexity.

Theorem 88 (Strongly convex, diminishing stepsize). Suppose F is μ -strongly convex and $\mathbb{E}\|\nabla f(x, \xi)\|^2 \leq B^2, \forall x$, then SGD with $\gamma_t = \frac{\gamma}{t}$ and $\gamma > 1/2\mu$ satisfies

$$\mathbb{E}\|x_t - x^*\|^2 \leq \frac{C(\gamma)}{t},$$

where

$$C(\gamma) = \max\left\{\frac{\gamma^2 B^2}{2\mu\gamma - 1}, \|x_0 - x^*\|^2\right\}.$$

Proof. Like in the proof of the previous case, we have

$$\|x_{t+1} - x^*\|^2 = \|x_t - x^*\|^2 + \gamma_t^2 \|\nabla f(x_t, \xi_t)\|^2 - 2\nabla f(x_t, \xi_t)^\top (x_t - x^*).$$

Also like in the previous proof and further using strong convexity of F , we have

$$\mathbb{E}\left[\nabla f(x_t, \xi_t)^\top (x_t - x^*)\right] = \mathbb{E}\left[\nabla F(x_t)^\top (x_t - x^*)\right] \geq \mu \mathbb{E}\|x_t - x^*\|^2.$$

This gives the following recursion,

$$\mathbb{E}\|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{2\mu\gamma}{t}\right) \mathbb{E}\|x_t - x^*\|^2 + \frac{\gamma^2 B^2}{t^2}.$$

The results follows by induction. ■

Thus, in theory, we see that a diminishing stepsize is necessary for SGD to converge to an optimal solution. However, in practice, constant stepsizes are often used with great success.

13.2 Adaptive methods

Often we do not know whether the problem is convex, L -smooth, or μ -strongly convex. Thus, we want the stepsize to adapt to the landscape of the function. The generic adaptive scheme looks like the following,

$$\begin{aligned} g_t &= \nabla f(x_t, \xi_t) \\ m_t &= \phi_t(g_1, \dots, g_t) \\ V_t &= \psi_t(g_1, \dots, g_t) \\ \hat{x}_t &= x_t - \alpha_t V_t^{-1/2} m_t \\ x_{t+1} &= \operatorname{argmin}_{x \in X} \left\{ (x - \hat{x}_t)^\top V_t^{1/2} (x - \hat{x}_t) \right\}. \end{aligned}$$

Momentum.

The most popular stochastic gradient descent methods are special cases of this scheme,

- Stochastic gradient descent:

$$\mathbf{m}_t = \mathbf{g}_t, \quad \mathbf{V}_t = \mathbf{I}.$$

- AdaGrad:

$$\mathbf{m}_t = \mathbf{g}_t, \quad \mathbf{V}_t = \frac{\text{diag}(\sum_{\tau=1}^t \mathbf{g}_\tau^2)}{t}.$$

- Adam:

$$\mathbf{m}_t = (1 - \alpha) \sum_{\tau=1}^t \alpha^{t-\tau} \mathbf{g}_\tau, \quad \mathbf{V}_t = (1 - \beta) \text{diag} \left(\sum_{\tau=1}^t \beta^{t-\tau} \mathbf{g}_\tau^2 \right).$$

Or, recursively:

$$\mathbf{m}_t = \alpha \mathbf{m}_{t-1} + (1 - \alpha) \mathbf{g}_t, \quad \mathbf{V}_t = \beta \mathbf{V}_{t-1} + (1 - \beta) \text{diag}(\mathbf{g}_t^2).$$

13.3 Variance reduction

Despite having a cheaper iteration cost than gradient descent,¹⁴ SGD requires more iterations,¹⁵ due to high variance. Stochastic variance-reduced (VR) methods try to achieve the best of both worlds by reducing the variance of SGD.¹⁶ We will present VR methods in the context of finite-sum optimization, which is a special case of stochastic optimization,

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \doteq \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}).$$

In the context of deep learning, we can see n as the number of data points and f_i the function w.r.t. the i -th data point, where we wish to minimize the objective function w.r.t. every data point with equal weight.

Suppose we want to estimate $\theta = \mathbb{E}[X]$, where X is a random variable. Let Y be another random variable. We can estimate θ as $\mathbb{E}[X - Y]$ if and only if $\mathbb{E}[Y] = 0$. Furthermore, $\mathbb{V}[X - Y] \leq \mathbb{V}[X]$ if Y is highly positively correlated with X . Specifically, if $\text{Cov}(X, Y) > \frac{1}{2} \mathbb{V}[Y]$, the variance will be reduced.¹⁷

Let $\alpha \in [0, 1]$. Using the following point estimator introduces a trade-off between variance and biasedness,

$$\hat{\theta}_\alpha = \alpha(X - Y) + \mathbb{E}[Y].$$

We then have the following expected value and variance,

$$\begin{aligned} \mathbb{E}[\hat{\theta}_\alpha] &= \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y] \\ \mathbb{V}[\hat{\theta}_\alpha] &= \alpha^2 (\mathbb{V}[X] + \mathbb{V}[Y] - 2\text{Cov}(X, Y)). \end{aligned}$$

Note that the estimator is unbiased if $\alpha = 1$, but the variance decreases when α decreases.

¹⁴ $\mathcal{O}(1)$ for SGD vs. $\mathcal{O}(n)$ for GD.

¹⁵ $\mathcal{O}(\kappa/\epsilon)$ for SGD vs. $\mathcal{O}(\kappa \log 1/\epsilon)$ for GD, where $\kappa = L/\mu$.

¹⁶ Classically, one can reduce variance by mini-batching, which reduces variance by $\mathcal{O}(1/|B_t|)$, where B_t is the batch, but computational complexity increases by $\mathcal{O}(|B_t|)$. Variance can also be reduced by introducing momentum to the gradient step. However, this requires access to past stochastic gradients, which can be expensive in memory.

We will consider a more modern approach.

¹⁷ This is because $\mathbb{V}[X - Y] = \mathbb{V}[X] + \mathbb{V}[Y] - 2\text{Cov}(X, Y)$.

While SGD estimates the full gradient by $\nabla f_{i_t}(\mathbf{x}_t)$, VR methods estimate $\nabla F(\mathbf{x}_t)$ by

$$\mathbf{g}_t \doteq \alpha(\nabla f_{i_t}(\mathbf{x}_t) - Y) + \mathbb{E}[Y],$$

such that

$$\mathbb{E}\|\mathbf{g}_t - \nabla F(\mathbf{x}_t)\|^2 \xrightarrow{t \rightarrow \infty} 0.$$

VR property.

The key idea is that if \mathbf{x}_t is not too far away from previous iterates $\mathbf{x}_{1:t-1}$, then we can leverage previous gradient information to construct positively correlated control variates Y . The question is thus how to design Y , given previous gradient information, such that it has low computational and space complexity.

Stochastic average gradient. The idea behind stochastic average gradient (SAG) is to keep track of the latest gradients for all points $i \in [n]$. Then, we estimate the full gradient by the average of these recent gradients,

$$\mathbf{g}_t = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t \approx \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t) = \nabla F(\mathbf{x}_t).$$

Thus, we update the past gradients as

$$\mathbf{v}_i^t = \begin{cases} \nabla f_{i_t}(\mathbf{x}_t) & i = i_t \\ \mathbf{v}_i^{t-1} & i \neq i_t. \end{cases}$$

Equivalently, we thus have the following update rule,

$$\begin{aligned} \mathbf{g}_t &= \mathbf{g}_{t-1} - \frac{1}{n} \mathbf{v}_{i_t}^{t-1} + \frac{1}{n} \nabla f_{i_t}(\mathbf{x}_t) \\ &= \frac{1}{n} (\nabla f_{i_t}(\mathbf{x}_t) - \mathbf{v}_{i_t}^{t-1}) + \mathbf{g}_{t-1}. \end{aligned}$$

Thus, we have $\alpha = 1/n$ and $Y = \mathbf{v}_{i_t}^{t-1}$ with $\mathbb{E}[Y] = \mathbf{g}_{t-1}$.

The downside of this approach is that it has a biased gradient, a large $\mathcal{O}(nd)$ memory cost, and it is hard to analyze. But, we gain a total complexity of $\mathcal{O}((n + \kappa_{\max}) \log^{1/\epsilon})$, where $\kappa_{\max} = \max_{i \in [n]} L_i / \mu$, where L_i is the smoothness parameter of f_i .

SAGA is an unbiased version of SAG, because it sets $\alpha = 1$, but still enjoys the same benefits as SAG with a much simpler proof. However, we still have a higher memory cost than SGD, which we would like to get rid of.

Stochastic variance reduced gradient. The key idea behind stochastic variance reduced gradient (SVRG) is to build covariates based on a fixed reference point $\tilde{\mathbf{x}}$. We then need to balance the frequency of updating this reference point and variance reduction.¹⁸ The intuition behind this method is the closer $\tilde{\mathbf{x}}$ is to \mathbf{x}_t , the smaller the variance is of the gradient estimator.

¹⁸ More updates cause lower variance, but increased complexity.

Algorithm	Iterations	Iteration cost
Gradient descent	$\mathcal{O}(\kappa \log 1/\epsilon)$	$\mathcal{O}(n)$
Stochastic gradient descent	$\mathcal{O}(\kappa_{\max}/\epsilon)$	$\mathcal{O}(1)$
Variance-reduced method	$\mathcal{O}((n + \kappa_{\max}) \log 1/\epsilon)$	$\mathcal{O}(1)$

Table 1. Complexity of μ -strongly convex and L -smooth finite-sum optimization, where n is the number of functions, $\kappa = L/\mu$, and $\kappa_{\max} = \max_{i \in [n]} L_i/\mu$, where L_i is the smoothness parameter of f_i .

The algorithm works by updating $\tilde{\mathbf{x}}$ every m -th update and estimates the gradient by

$$\nabla f_{i_t}(\mathbf{x}_t) - \nabla f_{i_t}(\tilde{\mathbf{x}}) + \nabla F(\tilde{\mathbf{x}}).$$

Thus, we have $\alpha = 1$ and $Y = \nabla f_{i_t}(\tilde{\mathbf{x}})$ with $\mathbb{E}[Y] = \nabla F(\tilde{\mathbf{x}})$.

While we gain the low memory cost of $\mathcal{O}(d)$, we now need to do $\mathcal{O}(n + 2m)$ gradient evaluations per epoch, where the n comes from computing $\mathbb{E}[Y]$ and $2m$ comes from computing $\nabla f_{i_t}(\mathbf{x}_t)$ and $\nabla f_{i_t}(\tilde{\mathbf{x}})$. This method has the same complexity as SAG and SAGA.

References

- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018.
- Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24(109):23–26, 1970.
- John Greenstadt. Variations on variable-metric methods. *Mathematics of Computation*, 24(109):1–22, 1970.
- Yu Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- Julie Nutini, Mark Schmidt, Issam Laradji, Michael Friedlander, and Hoyt Koepke. Coordinate descent converges faster with the gauss-southwell rule than random selection. In *International Conference on Machine Learning*, pages 1632–1641. PMLR, 2015.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.