# 1 Optimization and coercive functions

## 1.1 Review

As we've seen before in this class, closed and bounded subsets of $\mathbb{R}^n$ are useful in minimization because of the following two results:

**Theorem 1.1** (Extreme value theorem). *If $D \subseteq \mathbb{R}^n$ is a closed and bounded set, and $f : D \to \mathbb{R}$ is a continuous function, then $f$ has a global minimizer on $D$.*

**Theorem 1.2** (Bolzano–Weierstrass theorem). *If $D \subseteq \mathbb{R}^n$ is a closed and bounded set, and $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots$ is a sequence of elements of $D$, then there is a subsequence $\mathbf{x}^{(i_1)}, \mathbf{x}^{(i_2)}, \ldots$ which converges to some $\mathbf{x}^* \in D$.*

We say that a continuous function $f : \mathbb{R}^n \to \mathbb{R}$ is *coercive* if, for all $c \in \mathbb{R}$, the sublevel set $\{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \leq c\}$ is bounded. (This set is always guaranteed to be closed.) This definition is motivated by letting us apply the results above, and in particular, we have the following corollary of the extreme value theorem:

**Corollary 1.1.** *If $f$ is coercive, then $f$ has a global minimizer on $\mathbb{R}^n$.*

*Proof.* Pick any point $\mathbf{y} \in \mathbb{R}^n$. Let $S = \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \leq f(\mathbf{y})\}$.

$S$ is closed (by an argument involving the continuity of $f$), and by the definition of coercive, $S$ is also bounded. By the extreme value theorem, $f$ has a global minimizer $\mathbf{x}^*$ on $S$.

This $\mathbf{x}^*$ is actually a global minimizer on $\mathbb{R}^n$: if $\mathbf{x} \in S$, then by definition of the global minimizer, $f(\mathbf{x}^*) \leq f(\mathbf{x})$, but if $\mathbf{x} \notin S$, then $f(\mathbf{x}^*) \leq f(\mathbf{y}) < f(\mathbf{x})$, so we still have $f(\mathbf{x}^*) \leq f(\mathbf{x})$. $\qquad\square$

## 1.2 Some more results about coercive functions

We already saw one way to get examples of coercive functions earlier in this course. It's from dealing with one-dimensional functions:

1. For functions $f : \mathbb{R} \to \mathbb{R}$, $f$ is coercive if $\lim_{x \to \infty} f(x) = \lim_{x \to -\infty} f(x) = +\infty$.

2. If $f_1, f_2, \ldots, f_n$ are coercive functions $\mathbb{R} \to \mathbb{R}$, then the function $f(\mathbf{x}) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$ is a coercive function $\mathbb{R}^n \to \mathbb{R}$.

So, for example, $f(x, y) = (x^2 - 100x) + (y^4 + y + 1)$ is a coercive function.

Another nice property is easy to show from the definition.

---

**Lemma 1.1.** *If $f : \mathbb{R}^n \to \mathbb{R}$ is coercive, and $g : \mathbb{R}^n \to \mathbb{R}$ is continuous and bounded below, then $f + g$ is coercive.*

*Proof.* Let $L \in \mathbb{R}$ be a lower bound on $g$: some number such that $g(\mathbf{x}) \leq L$ for all $\mathbf{x} \in \mathbb{R}^n$.

Then the sublevel set $S = \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) + g(\mathbf{x}) \leq c\}$ is contained inside the set $S' = \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) + L \leq c\}$.

(If $\mathbf{x} \in S$, then $f(\mathbf{x}) + g(\mathbf{x}) \leq c$, and $f(\mathbf{x}) + L \leq f(\mathbf{x}) + g(\mathbf{x})$ no matter what $\mathbf{x}$ is, so $f(\mathbf{x}) + L \leq c$ as well, proving that $\mathbf{x} \in S'$.)

But $S'$ is just $\{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \leq c - L\}$, which is a sublevel set of $f$. So $S'$ is bounded because $f$ is coercive; therefore $S$ is bounded. □

Some special cases of this lemma:

- Since coercive functions have global minimizers, they are always bounded below, so in particular, the sum of two coercive functions is coercive.

- If $f$ is coercive and $h$ is a continuous function such that $f(\mathbf{x}) \leq h(\mathbf{x})$ for all $\mathbf{x}$, then $h = f + g$, where $g = f - h$, and $g$ is bounded below (by 0), so $h$ is also coercive.

This is a simple property, but we'll later see some clever ways to apply it.

## 2 Coercive functions and the penalty method

In the penalty method, we convert the problem

$$(P) \qquad \begin{cases} \underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \end{cases}$$

into an unconstrained optimization problem: minimizing the modified objective function

$$F_k(\mathbf{x}) = f(\mathbf{x}) + k\left[(g_1^+(\mathbf{x}))^2 + (g_2^+(\mathbf{x}))^2 + \cdots + (g_m^+(\mathbf{x}))^2\right]$$

over all $\mathbf{x} \in \mathbb{R}^n$, where $k$ is some large number.

In theory, as $k \to \infty$, the global minimizer of $F_k$ should approach an optimal solution of $P$. We showed last time that *if* that convergence happens, then the point we converge to is an optimal solution (under some hypotheses which are all satisfied in the theorem below).

**Theorem 2.1.** *Suppose that $P$ is feasible (there exists so $\mathbf{y} \in \mathbb{R}^n$ satisfying $g(\mathbf{y}) \leq \mathbf{0}$), $g_1, g_2, \ldots, g_m$ are continuous, and $f$ is coercive.*

*Then there is some sequence $k_1, k_2, k_3, \ldots$ of real numbers such that $\lim_{i \to \infty} k_i = \infty$, the global minimizers $\mathbf{x}^*(k_i)$ all exist, and as $i \to \infty$, the points $\mathbf{x}^*(k_i)$ converge to some $\mathbf{x}^* \in \mathbb{R}^n$.*

*(In which case, we already know that $\mathbf{x}^*$ is an optimal solution of $P$.)*

*Proof.* First, we show that $\mathbf{x}^*(k)$ exists for all $k > 0$. For all $k > 0$, the penalty term is a nonnegative continuous function of $\mathbf{x}$. So $F_k(\mathbf{x})$ is a sum of a coercive function and a function that's bounded below; therefore $F_k(\mathbf{x})$ is a coercive function, and has a global minimizer.

Let $\mathbf{y}$ be a point satisfying $\mathbf{g}(\mathbf{y}) \leq \mathbf{0}$. Then all the global minimizers $\mathbf{x}^*(k)$ are contained in the sublevel set $S = \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{y})\}$, because points outside $S$ will have a value of $F_k$ bigger than $\mathbf{y}$ does.

Because $f$ is continuous, $S$ is closed; because $f$ is coercive, $S$ is bounded. So the sequence of $\mathbf{x}^*(k)$ for $k = 1, 2, 3, \ldots$ has a convergent subsequence by the Bolzano–Weierstrass theorem. $\square$

By the way, the reason that we need to invoke Bolzano–Weierstrass here is the choice problem: we can't guarantee that $\mathbf{x}^*(k)$ converges as $k \to \infty$, which would be the expected behavior, because maybe there are multiple global minimizers to choose from for each $k$, and by picking a specific one to be $\mathbf{x}^*(k)$, we are making bad, non-convergent choices.

# 3 When are polynomials coercive?

The more we know about coercive functions, the more useful results about them become. So here are some conditions to help us classify polynomials as coercive.

## 3.1 Quadratic forms

This is a cute result that's also an example of the extreme value theorem in action.

**Theorem 3.1.** *If $A$ is an $n \times n$ positive definite matrix, then the quadratic form $f(\mathbf{x}) = \mathbf{x}^\mathsf{T} A \mathbf{x}$ is coercive.*

*Proof.* Let $S$ be the set $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| = 1\}$. $S$ is closed and bounded, so $f(\mathbf{x})$ has a global minimizer $\mathbf{x}^*$ on $S$. Let $\alpha = f(\mathbf{x}^*)$.

Fun fact: actually $\alpha$ is the smallest eigenvalue of $A$. But we don't need to know that. All we need to know is that $\alpha = \mathbf{x}^{*\mathsf{T}} A \mathbf{x}^* > 0$, because $A$ is positive definite, and that for all $\mathbf{x}$ with $\|\mathbf{x}\| = 1$, $f(\mathbf{x}) \geq \alpha$.

Now take an arbitrary $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x} \neq \mathbf{0}$. We have

$$f(\mathbf{x}) = \mathbf{x}^\mathsf{T} A \mathbf{x} = \left( \frac{\mathbf{x}}{\|\mathbf{x}\|} \right)^\mathsf{T} A \left( \frac{\mathbf{x}}{\|\mathbf{x}\|} \right) \cdot \|\mathbf{x}\|^2 \geq \alpha \cdot \|\mathbf{x}\|^2.$$

We can show that $\alpha\|\mathbf{x}\|^2$ is a coercive function, because the sublevel set $\{\mathbf{x} \in \mathbb{R}^n : \alpha\|\mathbf{x}\|^2 \leq c\}$ is the disk around $\mathbf{0}$ of radius $\sqrt{\frac{c}{\alpha}}$, which is bounded.

Since $f(\mathbf{x}) \geq \alpha\|\mathbf{x}\|^2$ for all $\mathbf{x}$ (including $\mathbf{0}$, because both functions are 0 there), we know that $f$ is a coercive function as well. $\square$

In fact, this condition goes both ways: if $A$ is a symmetric matrix, then $\mathbf{x}^\mathsf{T} A \mathbf{x}$ is only coercive when $A$ is positive definite. If not, then we can find a nonzero $\mathbf{y}$ for which $\mathbf{y}^\mathsf{T} A \mathbf{y} \leq 0$; for any scalar $t$, we'll have $(t\mathbf{y})^\mathsf{T} A(t\mathbf{y}) \leq 0$, so the sublevel set $\{\mathbf{x} : \mathbf{x}^\mathsf{T} A \mathbf{x}\}$ will contain an entire unbounded line $\{t\mathbf{y} : t \in \mathbb{R}\}$.

## 3.2 Higher-degree polynomials

For functions of one variable, things are relatively straightforward. The leading term is the one that will affect the behavior as $x \to \pm\infty$. If the leading term has odd degree, then the polynomial will be positive in one direction and negative in the other; for example, a cubic polynomial $ax^3 + bx^2 + cx + d$ can never be coercive. If the leading term has even degree, then the polynomial is coercive if and only if the coefficient on the leading term is positive.

We'd like to apply the same logic to polynomials in several variables. For example, we want to say that if $f(x, y) = x^4 + y^4 + xy$, the $xy$ term is dominated by the $x^4$ and $y^4$ terms, so their behavior will determine the "coerciveness" of $f$. But it's not obvious how to tell when this actually works, and when it doesn't.

Here is a test; it's stated for polynomials in 3 variables, but the idea generalizes to any number. (I just didn't want to write a mess of indices.)

**Theorem 3.2.** *A polynomial $f(x, y, z)$ is coercive if both of the following hold:*

1. *It contains $x^A, y^B, z^C$ terms with positive coefficients, where $A, B, C$ are some **even** integers.*

2. *For every other term $x^a y^b z^c$ (with any coefficient) that could potentially be negative, we have $\frac{a}{A} + \frac{b}{B} + \frac{c}{C} < 1$.*

*Proof.* If $\frac{a}{A} + \frac{b}{B} + \frac{c}{C} < 1$, then we can choose positive real numbers $A' < A$, $B' < B$, and $C' < C$ such that $\frac{a}{A'} + \frac{b}{B'} + \frac{c}{C'} = 1$.

Now apply the AM-GM inequality: $\dfrac{a}{A'} x^{A'} + \dfrac{b}{B'} y^{B'} + \dfrac{c}{C'} z^{C'} \geq x^a y^b z^c$.

This holds only for $x, y, z \geq 0$, but we can extend it to all $x, y, z$ by replacing them with $|x|, |y|, |z|$ on the left-hand side. We really want the negation of this, though:

$$-\frac{a}{A'}|x|^{A'} - \frac{b}{B'}|y|^{B'} - \frac{c}{C'}|z|^{C'} \leq -|x^a y^b z^c| \leq x^a y^b z^c.$$

So we can replace the $x^a y^b z^c$ term with $|x|^{A'}$, $|y|^{B'}$, and $|z|^{C'}$ terms (with negative coefficients, but we don't care about coefficients here). This only makes the function smaller. Therefore, if the resulting function is coercive, so was the original.

After we do replace all such mixed terms (and drop any mixed terms that are always nonnegative, such as $x^2 y^2 z^2$), we have a sum of functions of $x$, $y$, and $z$. These are all coercive and therefore so is their sum. For example, the function of $x$ has an $x^A$ term, and some lower-order terms with $|x|^{A'}$ for some $A' < A$, so $x^A$ (which is coercive) determines its behavior. $\square$

This condition is sufficient, but not necessary. If we allow terms $x^a y^b z^c$ with $\frac{a}{A} + \frac{b}{B} + \frac{c}{C} = 1$, sometimes the function we get is still coercive. But then, the coefficients start to matter: for example, by the theorem about quadratic forms, $x^2 + xy + y^2$ is coercive but $x^2 + 3xy + y^2$ is not. This can get tricky.