

Some Basic Optimization, Convex Geometry, and Linear Algebra

Rasmus Kyng

Lecture 2 — Wednesday, February 26

1 Overview

In these lecture notes we will

1. Start with an overview (i.e. this list).
2. Learn some basic terminology and facts about optimization.
3. Recall our definition of convex functions and see how convex functions can also be understood in terms of a characterization based on first derivatives.
4. See how the first derivatives of a convex function can certify that we are at a global minimum.
5. Review some standard linear algebra that we will need in later lectures.

2 Optimization Problems

Focusing for now on optimization over $\mathbf{x} \in \mathbb{R}^n$, we usually write optimization problems as:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & (\text{or } \max) \quad f(\mathbf{x}) \\ \text{s.t.} \quad & g_1(\mathbf{x}) \leq b_1 \\ & \cdot \\ & \cdot \\ & \cdot \\ & g_m(\mathbf{x}) \leq b_m \end{aligned}$$

where $\{g_i(\mathbf{x})\}_{i=1}^m$ encode the constraints. For example, in the following optimization problem from the previous lecture

$$\begin{aligned} \min_{\mathbf{f} \in \mathbb{R}^E} \quad & \sum_e r(e) \mathbf{f}(e)^2 \\ \text{s.t.} \quad & \mathbf{B}\mathbf{f} = \mathbf{d} \end{aligned}$$

we have the constraint $\mathbf{B}\mathbf{f} = \mathbf{d}$. Notice that we can rewrite this constraint as $\mathbf{B}\mathbf{f} \leq \mathbf{d}$ and $-\mathbf{B}\mathbf{f} \leq -\mathbf{d}$ to match the above setting. The set of points which respect the constraints is called the *feasible set*.

Definition 2.1. For a given optimization problem the set $\mathcal{F} = \{\mathbf{x} \in \mathbb{R}^n : g_i(\mathbf{x}) \leq b_i, \forall i \in [m]\}$ is called the **feasible set**. A point $\mathbf{x} \in \mathcal{F}$ is called a **feasible point**, and a point $\mathbf{x}' \notin \mathcal{F}$ is called an **infeasible point**.

Ideally, we would like to find optimal solutions for the optimization problems we consider. Let's define what we mean exactly.

Definition 2.2. For a *maximization* problem \mathbf{x}^* is called an **optimal solution** if $f(\mathbf{x}^*) \geq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{F}$. Similarly, for a *minimization* problem \mathbf{x}^* is an optimal solution if $f(\mathbf{x}^*) \leq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{F}$.

What happens if there are *no feasible points*? In this case, an optimal solution cannot exist, and we say the problem is infeasible.

Definition 2.3. If $\mathcal{F} = \emptyset$ we say that the optimization problem is **infeasible**. If $\mathcal{F} \neq \emptyset$ we say the optimization problem is **feasible**.

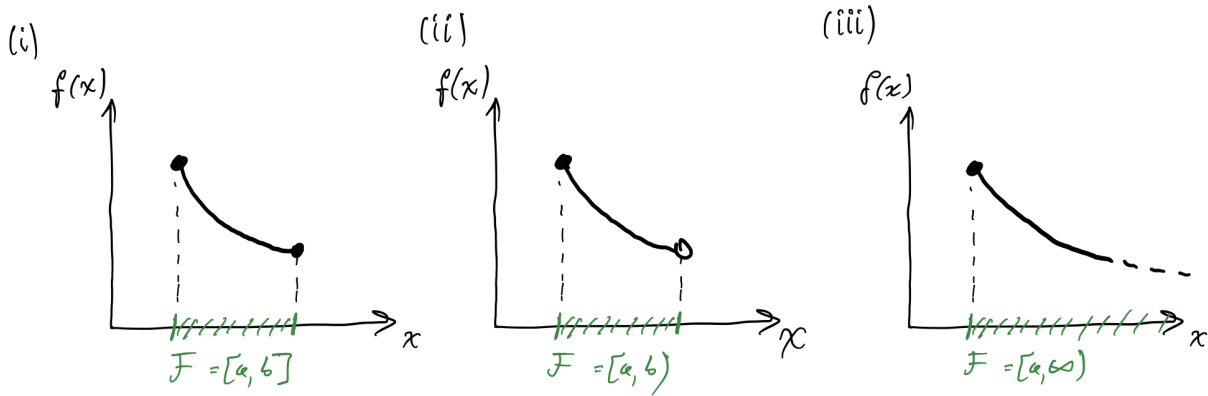


Figure 1

Consider three examples depicted in Figure 1:

- (i) $\mathcal{F} = [a, b]$
- (ii) $\mathcal{F} = [a, b)$
- (iii) $\mathcal{F} = [a, \infty)$

In the first example, the minimum of the function is attained at b . In the second case the region is open and therefore there is no minimum function value, since for every point we will choose, there will always be another point with a smaller function value. Lastly, in the third example, the region is unbounded and the function decreasing, thus again there will always be another point with a smaller function value.

Sufficient Condition for Optimality. The following theorem, which is a fundamental theorem in real analysis, gives us a sufficient (though not necessary) condition for optimality.

Theorem (Extreme Value Theorem). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function and $\mathcal{F} \subseteq \mathbb{R}^n$ be nonempty, bounded, and closed. Then, the optimization problem $\min f(\mathbf{x}) : \mathbf{x} \in \mathcal{F}$ has an optimal solution.

3 A Characterization of Convex Functions

Recall the definitions we introduced in the first lecture of convex sets and convex functions:

Definition 3.1. A set $S \subseteq \mathbb{R}^n$ is called a **convex set** if any two points in S contain their line, i.e. for any $\mathbf{x}, \mathbf{y} \in S$ we have that $\theta\mathbf{x} + (1 - \theta)\mathbf{y} \in S$ for any $\theta \in [0, 1]$.

Definition 3.2. For a convex set $S \subseteq \mathbb{R}^n$, we say that a function $f : S \rightarrow \mathbb{R}$ is **convex on S** if for any two points $\mathbf{x}, \mathbf{y} \in S$ and any $\theta \in [0, 1]$ we have that:

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

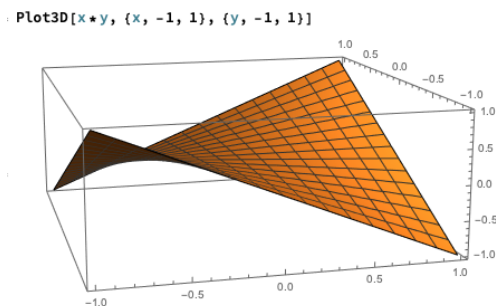


Figure 2: This plot shows the function $f(x, y) = xy$. For any fixed y_0 , the function $h(x) = f(x, y_0) = xy_0$ is linear in x , and so is a convex function in x . But is f convex?

We will first give an important characterization of convex function. To do so, we need to characterize multivariate functions via their Taylor expansion.

Notation for this lecture. In this lecture, we frequently consider a multivariate functions f whose domain is a set $S \subseteq \mathbb{R}^n$, which we will require to be open. When we additionally require that S is convex, we will specify this. Note that $S = \mathbb{R}^n$ is both open and convex and it suffices to keep this case in mind. Things sometimes get more complicated if S is not open, e.g. when the domain of f has a boundary. We will leave those complications for another time.

3.1 First-order Taylor approximation

Definition 3.3. The **gradient** of a function $f : S \rightarrow \mathbb{R}$ at point $\mathbf{x} \in S$ is denoted $\nabla f(\mathbf{x})$ is:

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f(\mathbf{x})}{\partial x(1)}, \dots, \frac{\partial f(\mathbf{x})}{\partial x(n)} \right]^\top$$

First-order Taylor expansion. For a function $f : \mathbb{R} \rightarrow \mathbb{R}$ of a single variable, differentiable at $x \in \mathbb{R}$

$$f(x + \delta) = f(x) + f'(x)\delta + o(|\delta|)$$

where by definition:

$$\lim_{\delta \rightarrow 0} \frac{o(|\delta|)}{|\delta|} = 0.$$

Similarly, a multivariate function $f : S \rightarrow \mathbb{R}$ is said to be (*Fréchet*) *differentiable* at $\mathbf{x} \in S$ when there exists $\nabla f(\mathbf{x}) \in \mathbb{R}^n$ s.t.

$$\lim_{\delta \rightarrow \mathbf{0}} \frac{\|f(\mathbf{x} + \delta) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top \delta\|_2}{\|\delta\|_2} = 0.$$

Note that this is equivalent to saying that $f(\mathbf{x} + \delta) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \delta + o(\|\delta\|_2)$.

We say that f is *continuously differentiable* on a set $S \subseteq \mathbb{R}^n$ if it is differentiable and in addition the gradient is continuous on S . A differentiable convex function whose domain is an open convex set $S \subseteq \mathbb{R}^n$ is always continuously differentiable¹.

Remark. In this course, we will generally err on the side of being informal about functional analysis when we can afford to, and we will not worry too much about the details of different notions of differentiability (e.g. Fréchet and Gateaux differentiability), except when it turns out to be important.

Theorem 3.4 (Taylor's Theorem, multivariate first-order remainder form). *If $f : S \rightarrow \mathbb{R}$ is continuously differentiable over $[\mathbf{x}, \mathbf{y}]$, then for some $\mathbf{z} \in [\mathbf{x}, \mathbf{y}]$,*

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{x}).$$

This theorem is useful for showing that the function f can be approximated by the affine function $\mathbf{y} \rightarrow f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ when \mathbf{y} is “close to” \mathbf{x} in some sense.

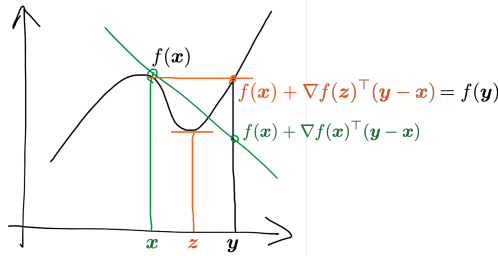


Figure 3: The convex function $f(\mathbf{y})$ sits above the linear function in \mathbf{y} given by $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$.

¹See p. 248, Corollary 25.5.1 in *Convex Analysis* by Rockafellar (my version is the Second print, 1972). Rockafellar's corollary concerns finite convex functions, because he otherwise allows convex functions that may take on the values $\pm\infty$.

3.2 Directional derivatives

Definition 3.5. Let $f : S \rightarrow \mathbb{R}$ be a function differentiable at $\mathbf{x} \in S$ and let us consider $\mathbf{d} \in \mathbb{R}^n$. We define the **derivative of f at \mathbf{x} in direction \mathbf{d}** as:

$$Df(\mathbf{x})[\mathbf{d}] = \lim_{\lambda \rightarrow 0} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda}$$

Proposition 3.6. $Df(\mathbf{x})[\mathbf{d}] = \nabla f(\mathbf{x})^\top \mathbf{d}$.

Proof. Using the first order expansion of f at \mathbf{x} :

$$f(\mathbf{x} + \lambda \mathbf{d}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\lambda \mathbf{d}) + o(\|\lambda \mathbf{d}\|_2)$$

hence, dividing by λ (and noticing that $\|\lambda \mathbf{d}\| = \lambda \|\mathbf{d}\|_2$):

$$\frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda} = \nabla f(\mathbf{x})^\top \mathbf{d} + o(\lambda \|\mathbf{d}\|_2)$$

letting λ go to 0 concludes the proof. \square

3.3 Lower bounding convex functions with affine functions

In order to prove the characterization of convex functions in the next section we will need the following lemma. This lemma says that any differentiable convex function can be lower bounded by an affine function.

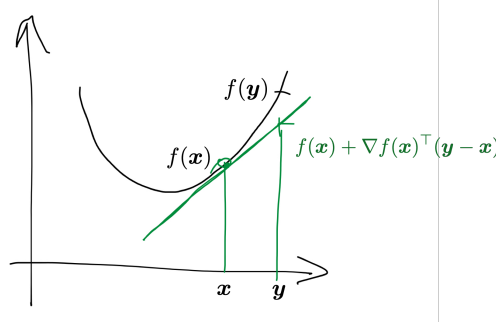


Figure 4: The convex function $f(\mathbf{y})$ sits above the linear function in \mathbf{y} given by $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$.

Theorem 3.7. Let S be an open convex subset of \mathbb{R}^n , and let $f : S \rightarrow \mathbb{R}$ be a differentiable function. Then, f is convex if and only if for any $\mathbf{x}, \mathbf{y} \in S$ we have that $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$.

Proof. $[\implies]$ Assume f is convex, then for all $\mathbf{x}, \mathbf{y} \in S$ and $\theta \in [0, 1]$, if we let $\mathbf{z} = \theta \mathbf{y} + (1 - \theta) \mathbf{x}$, we have that

$$f(\mathbf{z}) = f((1 - \theta) \mathbf{x} + \theta \mathbf{y}) \leq (1 - \theta) f(\mathbf{x}) + \theta f(\mathbf{y})$$

and therefore by subtracting $f(\mathbf{x})$ from both sides we get:

$$\begin{aligned} f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x}) &\leq \theta f(\mathbf{y}) + (1 - \theta)f(\mathbf{x}) - f(\mathbf{x}) \\ &= \theta f(\mathbf{y}) - \theta f(\mathbf{x}). \end{aligned}$$

Thus we get that (for $\theta > 0$):

$$\frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta} \leq f(\mathbf{y}) - f(\mathbf{x})$$

Applying Proposition 3.6 with $\mathbf{d} = \mathbf{x} - \mathbf{y}$ we have that:

$$\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) = \lim_{\theta \rightarrow 0^+} \frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta} \leq f(\mathbf{y}) - f(\mathbf{x}).$$

[\Leftarrow] Assume that $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in S$ and show that f is convex. Let $\mathbf{x}, \mathbf{y} \in S$ and $\mathbf{z} = \theta \mathbf{y} + (1 - \theta)\mathbf{x}$. By our assumption we have that:

$$f(\mathbf{y}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \quad (1)$$

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) \quad (2)$$

Observe that $\mathbf{y} - \mathbf{z} = (1 - \theta)(\mathbf{y} - \mathbf{x})$ and $\mathbf{x} - \mathbf{z} = \theta(\mathbf{y} - \mathbf{x})$. Thus adding θ times (1) to $(1 - \theta)$ times (2) gives cancellation of the vectors multiplying the gradient, yielding

$$\begin{aligned} \theta f(\mathbf{y}) + (1 - \theta)f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top \mathbf{0} \\ &= f(\theta \mathbf{y} + (1 - \theta)\mathbf{x}) \end{aligned}$$

This is exactly the definition of convexity. □

4 Conditions for optimality

We now want to find necessary and sufficient conditions for local optimality.

Definition 4.1. Consider a differentiable function $f : S \rightarrow \mathbb{R}$. A point $\mathbf{x} \in S$ at which $\nabla f(\mathbf{x}) = \mathbf{0}$ is called a **stationary point**.

Proposition 4.2. If \mathbf{x} is a local extremum of a differentiable function $f : S \rightarrow \mathbb{R}$ then $\nabla f(\mathbf{x}) = \mathbf{0}$.

Proof. Let us assume that \mathbf{x} is a local minimum for f . Then for all $\mathbf{d} \in \mathbb{R}^n$, $f(\mathbf{x}) \leq f(\mathbf{x} + \lambda \mathbf{d})$ for λ small enough. Hence:

$$0 \leq f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x}) = \lambda \nabla f(\mathbf{x})^\top \mathbf{d} + o(\|\lambda \mathbf{d}\|)$$

dividing by $\lambda > 0$ and letting $\lambda \rightarrow 0^+$, we obtain $0 \leq \nabla f(\mathbf{x})^\top \mathbf{d}$. But, taking $\mathbf{d} = -\nabla f(\mathbf{x})$, we get $0 \leq -\|\nabla f(\mathbf{x})\|_2^2$. This implies that $\nabla f(\mathbf{x}) = \mathbf{0}$.

The case where \mathbf{x} is a local maximum can be dealt with similarly. □

Remark 4.3. For this proposition to hold, it is important that S is open.

For convex functions however it turns out that a stationary point necessarily implies that the function is at its minimum. Together with the proposition above, this says that for a convex function on \mathbb{R}^n a point is optimal if and only if it is stationary.

Proposition 4.4. *Let $S \subseteq \mathbb{R}^n$ be an open convex set and let $f : S \rightarrow \mathbb{R}$ be a differentiable and convex function. If \mathbf{x} is a stationary point then \mathbf{x} is a global minimum.*

Proof. From Theorem 3.7 we know that for all $\mathbf{x}, \mathbf{y} \in S$: $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})(\mathbf{y} - \mathbf{x})$. Since $\nabla f(\mathbf{x}) = \mathbf{0}$ this implies that $f(\mathbf{y}) \geq f(\mathbf{x})$. As this holds for any $\mathbf{y} \in S$, \mathbf{x} is a global minimum. □

5 Linear Algebra Refresher

Semi-definiteness of a matrix. The following classification of symmetric matrices will be useful.

Definition 5.1. Let \mathbf{A} be a symmetric matrix in $\mathbb{R}^{n \times n}$. We say that \mathbf{A} is:

1. *positive definite* iff $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ for all $\mathbf{x} \in \mathbb{R}^n \setminus \{0\}$;
2. *positive semidefinite* iff $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$;
3. If neither \mathbf{A} nor $-\mathbf{A}$ is positive semi-definite, we say that \mathbf{A} is *indefinite*.

Example: indefinite matrix. Consider the following matrix \mathbf{A} :

$$\mathbf{A} := \begin{bmatrix} +4 & -1 \\ -1 & -2 \end{bmatrix}$$

For $\mathbf{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, we have $\mathbf{x}^\top \mathbf{A} \mathbf{x} = 4 > 0$. For $\mathbf{x} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ we have $\mathbf{x}^\top \mathbf{A} \mathbf{x} = -2 < 0$. \mathbf{A} is therefore indefinite.

The following theorem gives a useful characterization of (semi)definite matrices.

Theorem 5.2. *Let \mathbf{A} be a symmetric matrix in $\mathbb{R}^{n \times n}$.*

1. *\mathbf{A} is positive definite iff all its eigenvalues are positive;*
2. *\mathbf{A} is positive semidefinite iff all its eigenvalues are non-negative;*

In order to prove this theorem, let us first recall the Spectral Theorem for symmetric matrices.

Theorem 5.3 (The Spectral Theorem for Symmetric Matrices). *For all symmetric $\mathbf{A} \in \mathbb{R}^{n \times n}$ there exist $\mathbf{V} \in \mathbb{R}^{n \times n}$ and a diagonal matrix $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$ s.t.*

1. $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$.
2. $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$ (the $n \times n$ identity matrix). I.e. the columns of \mathbf{V} form an orthonormal basis. Furthermore, \mathbf{v}_i is an eigenvector of $\lambda_i(\mathbf{A})$, the i th eigenvalue of \mathbf{A} .

3. $\mathbf{A}_{ii} = \lambda_i(\mathbf{A})$.

Using the Spectral Theorem, we can show the following result:

Theorem 5.4 (The Courant-Fischer Theorem). *Let \mathbf{A} be a symmetric matrix in $\mathbb{R}^{n \times n}$, with eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. Then*

1.

$$\lambda_i = \min_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W)=i}} \max_{\mathbf{x} \in W, \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

2.

$$\lambda_i = \max_{\substack{\text{subspace } W \subseteq \mathbb{R}^n \\ \dim(W)=n+1-i}} \min_{\mathbf{x} \in W, \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

Theorem 5.2 is an immediate corollary of Theorem 5.4, since we can see that minimum value of the quadratic form $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ over $\mathbf{x} \in W = \mathbb{R}^n$ is $\lambda_1(\mathbf{A}) \|\mathbf{x}\|_2^2$.

Proof of Theorem 5.4. We start by showing Part 1.

Consider letting $W = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_i\}$, and normalize $\mathbf{x} \in W$ so that $\|\mathbf{x}\|_2 = 1$. Then $\mathbf{x} = \sum_{j=1}^i c(j) \mathbf{v}_j$ for some vector $\mathbf{c} \in \mathbb{R}^i$ with $\|\mathbf{c}\|_2 = 1$.

Using the decomposition from Theorem 5.3 $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$ where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues of \mathbf{A} , which we take to be sorted in increasing order. Then $\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top \mathbf{V}^\top \mathbf{\Lambda} \mathbf{V} \mathbf{x} = (\mathbf{V} \mathbf{x})^\top \mathbf{\Lambda} (\mathbf{V} \mathbf{x}) = \sum_{j=1}^i \lambda_j c(j)^2 \leq \lambda_i \|\mathbf{c}\|_2^2 = \lambda_i$. So this choice of W ensures the maximizer cannot achieve a value above λ_i .

But is it possible that the “minimizer” can do better by choosing a different W ? Let $T = \text{span}\{\mathbf{v}_i, \dots, \mathbf{v}_n\}$. As $\dim(T) = n + 1 - i$ and $\dim(W) = i$, we must have $\dim(W \cap T) \geq 1$, by a standard property of subspaces. Hence for any W of this dimension,

$$\begin{aligned} \max_{\mathbf{x} \in W, \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} &\geq \max_{\mathbf{x} \in W \cap T, \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \\ &\geq \min_{\substack{\text{subspace } V \subseteq T \\ \dim(V)=1}} \max_{\mathbf{x} \in V, \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_i, \end{aligned}$$

where the last equality follows from a similar calculation to our first one. Thus, λ_i can always be achieved by the “maximizer” for all W of this dimension.

Part 2 can be dealt with similarly.

□

Example: a positive semidefinite matrix. Consider the following matrix \mathbf{A} :

$$\mathbf{A} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

For $\mathbf{x} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, we have $\mathbf{Ax} = \mathbf{0}$, so $\lambda = 0$ is an eigenvalue of \mathbf{A} . For $\mathbf{x} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$, we have $\mathbf{Ax} = \begin{pmatrix} 2 \\ -2 \end{pmatrix} = 2\mathbf{x}$, so $\lambda = 2$ is the other eigenvalue of \mathbf{A} . As both are non-negative, by the theorem above, \mathbf{A} is positive semidefinite.

Since we are learning about symmetric matrices, there is one more fact that everyone should know about them. We'll use $\lambda_{\max}(\mathbf{A})$ denote maximum eigenvalue of a matrix \mathbf{A} , and $\lambda_{\min}(\mathbf{A})$ the minimum.

Claim 5.5. *For a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\|\mathbf{A}\| = \max(|\lambda_{\max}(\mathbf{A})|, |\lambda_{\min}(\mathbf{A})|)$.*