

Appendix B

Convex Sets and Functions

The notion of convexity arises when we try to formalize the property that “good local decisions lead to globally optimal solutions.” Consider a generic unconstrained optimization problem

$$\text{minimize } f(\mathbf{x}). \tag{B.0.1}$$

Here $\mathbf{x} \in \mathbb{R}^n$ is the variable of optimization, and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the objective function, which we are trying to make as small as possible using a numerical algorithm. Figure B.1 displays two objective functions f . The one on the right has many peaks and valleys – it may be very difficult to find the lowest valley, corresponding to the global optimum \mathbf{x}_* . Moreover, for the function f on the right, local information around a point \mathbf{x} is not particularly helpful for determining what direction to move to reach the global optimum. In contrast, the bowl-shaped function on the left is much more amenable to global optimization – a “gradient descent” type algorithm, that simply determined which direction to move by considering the slope of the graph of the function, would easily “ski” down to the global minimum.

The notion of *convexity* formalizes this property. Convexity is a geometric property. It is convenient to first introduce the notion of a convex set, and then extend this definition to functions.

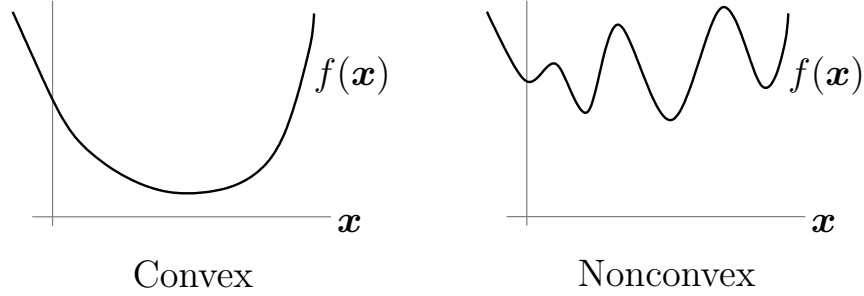


Figure B.1. **Two optimization problems** $\min f(\mathbf{x})$. The objective f at left appears to be amenable to global optimization, while the one at right appears to be more challenging.

B.1 Convex Sets

A set C is said to be *closed* if it contains its boundary. More precisely, for any converging sequence of points $\{\mathbf{x}_k\}$ in C , we must have:

$$\mathbf{x}_k \rightarrow \bar{\mathbf{x}} \Rightarrow \bar{\mathbf{x}} \in C.$$

A set $C \subseteq \mathbb{R}^n$ is *convex* if for every pair of points $\mathbf{x}, \mathbf{x}' \in C$, the line segment obtained by joining the two points also lies entirely in C :

Definition B.1.1 (Convex set). $C \subseteq \mathbb{R}^n$ is convex if

$$\forall \mathbf{x}, \mathbf{x}' \in C, \quad \alpha \in [0, 1], \quad \alpha \mathbf{x} + (1 - \alpha) \mathbf{x}' \in C. \quad (\text{B.1.1})$$

Figure B.2 gives an example of two sets, one of which is convex and one of which is not.

Example B.1.2 (Convex sets). *Show that the following are convex:*

- Every affine subspace.
- Every norm ball $\mathfrak{B}_{\|\cdot\|} = \{\mathbf{x} \mid \|\mathbf{x}\| \leq 1\}$.
- The empty set.
- Any intersection $C = C_1 \cap C_2$ of two convex sets C_1, C_2 .

Proposition B.1.3. 1. The intersection of a collection of convex sets $\bigcap_i C_i$ is convex.

2. The image of a convex set under an affine transformation is convex.

Definition B.1.4 (Convex hull). The convex hull of any given set S is the minimal convex set containing S , denoted as $\text{conv}(S)$. If S contains a finite number of $S = \{\mathbf{x}_i\}_{i=1}^n$ points, we have

$$\text{conv}(S) = \left\{ \sum_{i=1}^n \alpha_i \mathbf{x}_i \mid \forall \alpha_i \geq 0 \text{ with } \sum_{i=1}^n \alpha_i = 1. \right\}. \quad (\text{B.1.2})$$

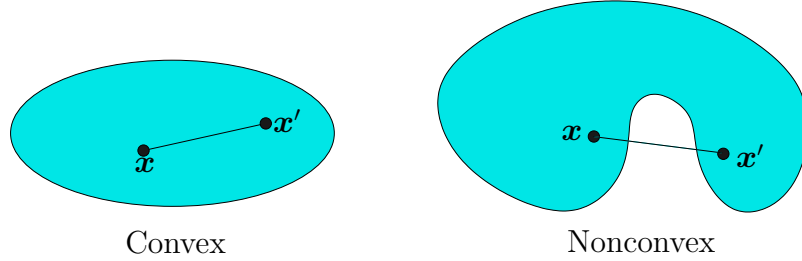


Figure B.2. **Convex and nonconvex sets.** A set is convex if we can select any pair of points \mathbf{x} , \mathbf{x}' in the set, and the line segment joining them lies entirely within the set. The set to the left has this property, while the set to the right does not.

B.2 Convex Functions

For a function $f : \mathcal{D} \rightarrow \mathbb{R}$ defined on a (convex) domain $\mathcal{D} \subseteq \mathbb{R}^n$, its *graph* is the set of pairs $(\mathbf{x}, f(\mathbf{x}))$ that can be generated by evaluating the function f at every point:

$$\text{graph}(f) = \{(\mathbf{x}, f(\mathbf{x})) \mid \mathbf{x} \in \mathcal{D}, f(\mathbf{x}) < +\infty\} \subseteq \mathbb{R}^{n+1}. \quad (\text{B.2.1})$$

We give another name to everything that lies above the graph: the *epigraph*:

$$\text{epi}(f) = \{(\mathbf{x}, t) \mid \mathbf{x} \in \mathcal{D}, t \in \mathbb{R}, f(\mathbf{x}) \leq t\} \subseteq \mathbb{R}^{n+1}. \quad (\text{B.2.2})$$

We say that f is a *convex function* if its epigraph is a convex set. Figure B.3 (right) illustrates this property. Figure B.3 (left) suggests an equivalent definition, which is sometimes easier to work with: f is convex if for any pair of points \mathbf{x} and \mathbf{x}' , the line segment joining $(\mathbf{x}, f(\mathbf{x}))$ and $(\mathbf{x}', f(\mathbf{x}'))$ lies entirely above the graph of f :

Definition B.2.1 (Convex function). *A function $f : \mathcal{D} \rightarrow \mathbb{R}$ is convex if for all $\mathbf{x}, \mathbf{x}' \in \mathcal{D}$ and $\alpha \in [0, 1]$,*

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}') \geq f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{x}'). \quad (\text{B.2.3})$$

Notice that above definitions do not require f is differentiable. If f is differentiable, the notion of convex can be characterized in terms of its derivatives. Since the epigraph is convex, then the tangent plane at each point of the graph should lie beneath the graph. The following statement makes this precise:

Proposition B.2.2 (First-Order Conditions). *Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be differentiable. Then f is convex if and only if it satisfies the condition:*

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}' - \mathbf{x})$$

for all $\mathbf{x}, \mathbf{x}' \in \mathcal{D}$.

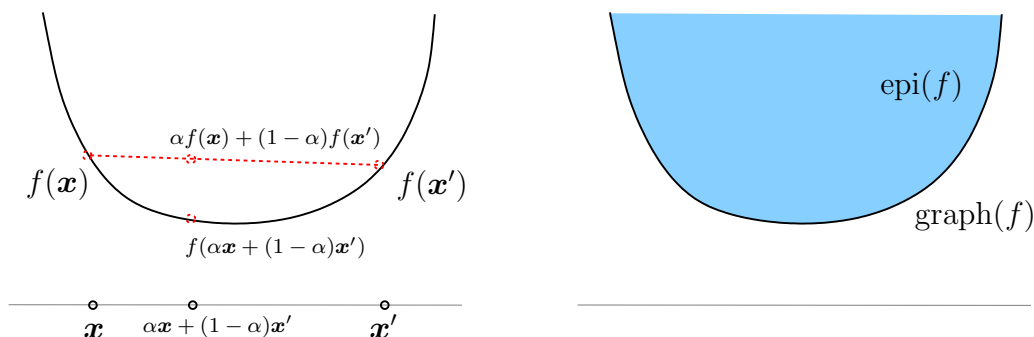


Figure B.3. **Convexity of functions:** a function f is convex if its epigraph $\text{epi}(f) = \{(\mathbf{x}, t) \mid t \geq f(\mathbf{x})\}$ is a convex set (right). This is true if and only if for every pair of points \mathbf{x} , \mathbf{x}' and scalar $\alpha \in [0, 1]$, $f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{x}') \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}')$. The picture at right illustrates this inequality: the segment joining $(\mathbf{x}, f(\mathbf{x}))$ and $(\mathbf{x}', f(\mathbf{x}'))$ lies above the graph of f .

This is precisely the geometry of the “nice” function in Figure B.1 (left). From this picture, it is clear that convexity is very favorable for global optimization.¹

You may also notice that in Figure B.3, the function $f(\mathbf{x})$ “curves upward”: its second derivative is nonnegative at every point of the domain. For twice differentiable functions, this leads to a simpler condition for con-

¹Once you’ve internalized the definition a bit, you may begin to wonder to what extent the implication “convexity \implies easy-to-optimize” is actually true. The convex functions that we encounter in this book will all possess special structure that makes them very amenable to efficient algorithms. However, this is not true of all convex functions – there exist convex functions that are NP-hard to optimize.

There also exist nonconvex functions that are easy to optimize – Chapter 9 provides a brief introduction to this emerging literature. However, if we want to talk about a class of functions, rather than a particular one, then there is a very beautiful motivation for studying convex functions. To appreciate this motivation, we need to first observe a useful fact: if $f(\mathbf{x})$ and $g(\mathbf{x})$ are convex functions, then for any $\alpha, \beta \geq 0$, $h(\mathbf{x}) = \alpha f(\mathbf{x}) + \beta g(\mathbf{x})$ is also convex. If we let \mathcal{F} be the largest class of continuously differentiable functions that satisfy the following three demands:

- Every linear function $\phi(\mathbf{x}) = \mathbf{a}^* \mathbf{x} + b$ is in \mathcal{F} ,
- Every nonnegative combination $\alpha f_1(\mathbf{x}) + \beta f_2(\mathbf{x})$ of $f_1, f_2 \in \mathcal{F}$ is in \mathcal{F} ,
- For every $f \in \mathcal{F}$, the stationarity condition $\nabla f(\mathbf{x}_*) = 0$ implies that \mathbf{x}_* is a global optimizer of f ,

it turns out that the \mathcal{F} is precisely the class of **convex**, continuously differentiable functions. You can interpret this as suggesting that for *global* solutions, convex functions really are the right general class of functions to study. For more details, see the book of Nesterov [Nesterov, 2003a].

vexity: the function is convex if and only if its second derivative at any point, and in any direction is positive. The following makes this precise:

Proposition B.2.3 (Second-Order Conditions). *Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be twice differentiable. Then f is convex if and only if its Hessian is positive semidefinite:*

$$\nabla^2 f(\mathbf{x}) \succeq 0$$

for all $\mathbf{x} \in \mathcal{D}$.

The class of convex functions includes important examples such as linear functions and norms:

Example B.2.4 (Convex functions). *Show that the following are convex functions:*

- Every affine function $f(\mathbf{x}) = \mathbf{a}^* \mathbf{x} + b$.
- Every norm $f(\mathbf{x}) = \|\mathbf{x}\|$.
- Every semidefinite quadratic $f(\mathbf{x}) = \mathbf{x}^* \mathbf{P} \mathbf{x}$, with $\mathbf{P} \succeq \mathbf{0}$.
- The maximum eigenvalue $\lambda_{\max}(\mathbf{X})$ is a convex function over the symmetric matrices \mathbf{X} .

Before continuing, we note one nice property of convex functions which will be useful for deriving an appropriate tractable replacement for the ℓ^0 norm.

Definition B.2.5 (Convex combination). *A convex combination of points $\mathbf{x}_1, \dots, \mathbf{x}_k$ is an expression of the form $\lambda_1 \mathbf{x}_1 + \dots + \lambda_k \mathbf{x}_k$, with $\lambda_i \geq 0$ for each i and $\sum_i \lambda_i = 1$.*

Lemma B.2.6 (Jensen's inequality). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. For any k , $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$, $\lambda_1, \dots, \lambda_k \in \mathbb{R}_+$, with $\sum_i \lambda_i = 1$,*

$$f\left(\sum_i \lambda_i \mathbf{x}_i\right) \leq \sum_i \lambda_i f(\mathbf{x}_i). \quad (\text{B.2.4})$$

Proof. The proof is by induction on k . For $k = 1$, there is nothing to show. Now suppose the claim is true for $1, \dots, k-1$. Then

$$f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i\right) \leq \left(\sum_{i=1}^{k-1} \lambda_i\right) f\left(\frac{\sum_{i=1}^{k-1} \lambda_i \mathbf{x}_i}{\sum_{i=1}^{k-1} \lambda_i}\right) + \lambda_k f(\mathbf{x}_k) \quad (\text{B.2.5})$$

$$\leq \sum_{i=1}^k \lambda_i f(\mathbf{x}_i) \quad (\text{B.2.6})$$

as desired. Above, the first step uses the definition of convexity, and the second uses the inductive hypothesis. \square

With this lemma, it is easy to show that any α -sublevel set of a convex function $f : \mathcal{D} \rightarrow \mathbb{R}$:

$$C_\alpha = \{\mathbf{x} \in \mathcal{D} \mid f(\mathbf{x}) \leq \alpha\} \quad (\text{B.2.7})$$

is a convex set. However, a function with all its sublevel sets being convex is not necessarily a convex function!² A function is said to be a *closed* function, if each sublevel set is a closed set. We typically only consider closed convex functions, unless otherwise stated.

Proposition B.2.7. *We can use convex functions to generate other associated convex functions:*

1. *A function is convex if and only if it is convex when restricted to any line that intersects its domain.*
2. *A weighted sum of convex functions with nonnegative weights is convex.*
3. *If f, g are convex functions and g is non-decreasing in its univariate domain, the $h(\mathbf{x}) = g(f(\mathbf{x}))$ is convex.*
4. *Given a collection of convex functions $f_\alpha : \mathcal{D} \rightarrow \mathbb{R}$, $\alpha \in \mathbf{A}$, their point-wise supremum*

$$f(\mathbf{x}) \doteq \sup_{\alpha \in \mathbf{A}} f_\alpha(\mathbf{x})$$

is also convex.

Example B.2.8. *The maximal eigenvalue of a symmetric matrix is a (closed) convex function.*

Proof. To see that, the maximal eigenvalue function can be written as

$$\sigma_{\max}(\mathbf{X}) = \sup\{\mathbf{y}^T \mathbf{X} \mathbf{y}\}, \quad \|\mathbf{y}\|_2 = 1.$$

Since the function is the point-wise supremum of a set of linear functions with respect to \mathbf{X} , it is a convex function. \square

Convex Envelop and Conjugate.

For any non-convex (closed) function $g : \mathcal{D} \rightarrow \mathbb{R}$ defined on a convex domain \mathcal{D} , it has a naturally associated convex function that bounds it from below:

Definition B.2.9 (Convex envelop). *The convex envelop of a closed function g is defined as*

$$\text{conv}g(\mathbf{x}) = \sup\{h(\mathbf{x}) \mid h(\mathbf{x}) \text{ convex \& } h(\mathbf{x}) \leq g(\mathbf{x}) \ \forall \mathbf{x} \in \mathcal{D}\}. \quad (\text{B.2.8})$$

²Such functions are called *quasi-convex*. Please find an example for yourself.

Let us define the (Fenchel) *conjugate* of a function $g(\mathbf{x})$ (not necessarily convex) as:

$$g^*(\boldsymbol{\lambda}) = \max_{\mathbf{x}} \boldsymbol{\lambda}^T \mathbf{x} - g(\mathbf{x}). \quad (\text{B.2.9})$$

The conjugate of a function g is essentially the negated dual function of g that we often see in the method of Lagrangian multipliers (see Section C.3).

Proposition B.2.10. *Assuming the conjugate is well-defined, we have the following:*

1. *The conjugate $g^*(\boldsymbol{\lambda})$ is always a convex function.*
2. $g^{**}(\mathbf{x}) = \text{conv}g(\mathbf{x})$.

Strong Convexity.

In this book, we sometimes are interested in stronger notion of convexity.

Definition B.2.11 (Strongly convex function). *A function $f : \mathcal{D} \rightarrow \mathbb{R}$ is strongly convex if f is convex and for all $\mathbf{x}, \mathbf{x}' \in \mathcal{D}$ and $\alpha \in [0, 1]$,*

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}') \geq f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{x}') + \mu \frac{\lambda(1 - \lambda)}{2} \|\mathbf{x} - \mathbf{x}'\|^2 \quad (\text{B.2.10})$$

for some $\mu > 0$.

Notice that the above definition does not require f to be differentiable. If f is first or second order differentiable, we have the following sufficient conditions for f being strongly convex.

Proposition B.2.12. *For a differentiable convex function f over \mathcal{D} , we have f is strongly convex if either of the following conditions hold:*

1. $f(\mathbf{x}') \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}' - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x}' - \mathbf{x}\|^2, \forall \mathbf{x}, \mathbf{x}' \in \mathcal{D};$
2. $\nabla^2 f(\mathbf{x}) \succeq \frac{\mu}{2} I, \forall \mathbf{x} \in \mathcal{D};$

for some $\mu > 0$.

However, as we see in Section 3.3.3, we are interested in strong convexity in a restricted sense.

B.3 Subdifferentials of Nonsmooth Convex Functions

For smooth, convex functions f , the local information encoded in the gradient ∇f and Hessian $\nabla^2 f$ characterize both the local and global behavior of f , allowing us to give optimality conditions and construct minimization algorithms. Familiar, classical algorithms such as gradient ascent, Newton's

Figure B.4. Subgradients and subdifferentials of convex functions.

method, and their variants, are all constructed using differential information. Moreover, as we saw in the previous section, these quantities play a critical role in characterizing convexity for smooth functions f .

It is a curious fact, then, that many of the most useful convex objective functions arising in high-dimensional data analysis are nondifferentiable: *their gradients and Hessians do not exist*. For example, the ℓ^1 norm $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ is nondifferentiable at any point $\mathbf{x} \in \mathbb{R}^n$ with fewer than n nonzero entries. These are precisely the points that we care about for sparse estimation! This nonsmooth behavior is actually desirable from the statistical perspective. However, it forces us to make recourse to analytical tools that are general enough to handle nondifferentiable functions. Fortunately, for convex functions, the nondifferential theory rests on simple, geometrically intuitive ideas, which we describe in this section. For accessible introductions to the general theory of convexity, we recommend [Nemirovski, 1995, Nemirovski, 2007, Nesterov, 2003a, Boyd and Vandenberghe, 2004].

The most important notion is that of a *subgradient* of a convex function, which provides a very satisfactory replacement for the gradient, when the function is not differentiable. Recall from Proposition B.2.2 that for convex, *differentiable* f ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle, \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D}. \quad (\text{B.3.1})$$

This inequality has a simple geometric interpretation, which we visualize in Figure B.4. We visualize the graph of the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$. The graph is the collection of points of the form $(\mathbf{x}, f(\mathbf{x})) \in \mathbb{R}^{n+1}$. The graph of

$$h(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

is a hyperplane, which is tangent to the graph of f at $(\mathbf{x}, f(\mathbf{x}))$. The inequality (B.3.1) says that at all points \mathbf{y} in the domain of the function f this tangent hyperplane lies below (or more precisely, not above) the graph of f .

Figure B.4 (right) visualizes the graph of another convex function f , which is not differentiable at point \mathbf{x} . The gradient of f *does not* exist at \mathbf{x} . Nevertheless, we can still define a nonvertical hyperplane $\mathcal{H} \subseteq \mathbb{R}^{n+1}$ that passes through $(\mathbf{x}, f(\mathbf{x}))$, and lies below the graph of f . This hyperplane has normal vector $(\mathbf{v}, -1)$, and can be expressed in notation as

$$\mathcal{H} = \{(\mathbf{y}, t) \mid t = f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{y} - \mathbf{x} \rangle\}. \quad (\text{B.3.2})$$

We say that $\mathbf{v} \in \mathbb{R}^n$ is a *subgradient* of f at \mathbf{x} if it defines a hyperplane that supports the graph of f at \mathbf{x} , and lies below the graph everywhere:

Definition B.3.1. Let $f : \mathcal{D} \subset \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. A vector \mathbf{v} is a subgradient of f at $\mathbf{x} \in \mathcal{D}$ if for all $\mathbf{y} \in \mathcal{D}$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{y} - \mathbf{x} \rangle. \quad (\text{B.3.3})$$

When f is differentiable, from Proposition B.2.2 it is clear that $\mathbf{v} = \nabla f(\mathbf{x})$ satisfies (B.3.3). When f is *nondifferentiable*, at a given point \mathbf{x} there can be multiple distinct hyperplanes that support the graph of f , and hence, there can be multiple subgradients \mathbf{v} (see Figure B.4). The collection of all subgradients is called the *subdifferential* of f at \mathbf{x} , and is denoted $\partial f(\mathbf{x})$. Formally:

Definition B.3.2. Let $f : \mathcal{D} \subseteq \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. The subdifferential $\partial f(\mathbf{x})$ is the collection of all subgradients of f at \mathbf{x} :

$$\partial f(\mathbf{x}) = \{\mathbf{v} \mid f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{y} - \mathbf{x} \rangle, \forall \mathbf{y} \in \mathcal{D}\}. \quad (\text{B.3.4})$$

Notice that if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable at \mathbf{x} , its subdifferential at \mathbf{x} is a singleton: $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$. This coincides with the classical definition of differentials.

A number of functions of interest have relatively simple subdifferentials.

Example B.3.3. As good exercises, the reader may try to verify what the subdifferentials for the following functions:

1. The subdifferential for $f(\mathbf{x}) = \|\mathbf{x}\|_1$ with $\mathbf{x} \in \mathbb{R}^n$.
2. The subdifferential for $f(\mathbf{x}) = \|\mathbf{x}\|_\infty$ with $\mathbf{x} \in \mathbb{R}^n$.
3. The subdifferential for $f(\mathbf{X}) = \sum_{j=1}^n \|\mathbf{X}\mathbf{e}_j\|_2$ with \mathbf{X} a matrix in $\mathbb{R}^{n \times n}$.
4. The subdifferential for $f(\mathbf{x}) = \|\mathbf{X}\|_*$ with \mathbf{X} a matrix in $\mathbb{R}^{n \times n}$.

Below are some basic properties of subdifferentials.

Lemma B.3.4 (Monotonicity Property). Given a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ and any $\mathbf{x}, \mathbf{x}', \mathbf{v}, \mathbf{v}' \in \mathbb{R}^n$ such that $\mathbf{v} \in \partial f(\mathbf{x})$ and $\mathbf{v}' \in \partial f(\mathbf{x}')$, we have

$$\langle \mathbf{x} - \mathbf{x}', \mathbf{v} - \mathbf{v}' \rangle \geq 0. \quad (\text{B.3.5})$$

Proof. From the definition of subgradient (B.3.2), we have

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{x}' - \mathbf{x} \rangle, \quad f(\mathbf{x}) \geq f(\mathbf{x}') + \langle \mathbf{v}', \mathbf{x} - \mathbf{x}' \rangle. \quad (\text{B.3.6})$$

Adding these two inequalities together we obtain:

$$f(\mathbf{x}) + f(\mathbf{x}') \geq f(\mathbf{x}) + f(\mathbf{x}') + \langle \mathbf{v} - \mathbf{v}', \mathbf{x}' - \mathbf{x} \rangle. \quad (\text{B.3.7})$$

Canceling $f(\mathbf{x}) + f(\mathbf{x}')$ from both sides obtains the desired result. \square

Lemma B.3.5. *If a convex function $f(\mathbf{x})$ has Lipschitz continuous gradient with constant L , then for any \mathbf{x}_1 and \mathbf{x}_2 , we have:*

$$\langle \nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2 \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|^2 \geq 0. \quad (\text{B.3.8})$$

Proof. Let us define a function $h(\mathbf{z}) \doteq f(\mathbf{z}) - \mathbf{z}^* \nabla f(\mathbf{x})$. Then $h(\mathbf{z})$ is convex and is minimized at $\mathbf{z} = \mathbf{x}$ (as $\nabla h(\mathbf{x}) = 0$). Hence for any \mathbf{z} , we have

$$h(\mathbf{x}) \leq h\left(\mathbf{z} - \frac{1}{L} \nabla h(\mathbf{z})\right) \leq h(\mathbf{z}) + \langle \nabla h(\mathbf{z}), -\frac{1}{L} \nabla h(\mathbf{z}) \rangle + \frac{L}{2} \left\| \frac{1}{L} \nabla h(\mathbf{z}) \right\|^2.$$

The last inequality comes from the fact that the function $f(\mathbf{x})$ (and hence $h(\mathbf{z})$) has Lipschitz continuous gradient with constant L . This gives

$$h(\mathbf{x}) \leq h(\mathbf{z}) - \frac{1}{2L} \|\nabla h(\mathbf{z})\|^2. \quad (\text{B.3.9})$$

Now applying the inequality to $\mathbf{x} = \mathbf{x}_1, \mathbf{z} = \mathbf{x}_2$ as well as the reverse case $\mathbf{x} = \mathbf{x}_2, \mathbf{z} = \mathbf{x}_1$, we get

$$\begin{aligned} f(\mathbf{x}_1) - \mathbf{x}_1^* \nabla f(\mathbf{x}_1) &\leq f(\mathbf{x}_2) - \mathbf{x}_2^* \nabla f(\mathbf{x}_1) - \frac{1}{2L} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|^2, \\ f(\mathbf{x}_2) - \mathbf{x}_2^* \nabla f(\mathbf{x}_2) &\leq f(\mathbf{x}_1) - \mathbf{x}_1^* \nabla f(\mathbf{x}_2) - \frac{1}{2L} \|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|^2. \end{aligned}$$

Adding these two together gives the desired bound (B.3.8). \square