

Lecture 1

Mathematically Preliminaries

Optimization in Data Science

Learning = Representation + Evaluation + Optimization¹

Linear Regression Given data points $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ drawn from an unknown distribution $\mathcal{P}(\mathcal{X}, \mathcal{Y})$ over the support $\mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^d \times \mathbb{R}$. Linear regression aims to find a linear function $f : \mathcal{X} \rightarrow \mathcal{Y}$ such that for any $(\mathbf{x}, y) \sim \mathcal{P}(\mathcal{X}, \mathcal{Y})$ we have $f(\mathbf{x}) = y$. Suppose the linear function has the form $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, the parameters $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ can be found by solving the following optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i + b - y_i)^2. \quad (1.1)$$

We denote the objective function as $\mathcal{L}(\mathbf{w}, b; \mathcal{D})$ known as the mean-squared error (MSE) loss in machine learning.

There are multiple common forms for linear regression:

$$\min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b; \mathcal{D}) \quad (\text{Least Squares Regression})$$

$$\min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b; \mathcal{D}) + \underbrace{\lambda \sum_{k=1}^d w_k^2}_{\text{Smooth penalty}} \quad (\text{Ridge Regression})$$

$$\min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b; \mathcal{D}) + \underbrace{\lambda \sum_{k=1}^d |w_k|}_{\text{Non-smooth penalty}} \quad (\text{Lasso Regression})$$

$$\min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b; \mathcal{D}) + \lambda_1 \sum_{k=1}^d |w_k| + \lambda_2 \sum_{k=1}^d w_k^2 \quad (\text{Elastic-Net})$$

¹Pedro Domingos. 2012. A few useful things to know about machine learning. Commun. ACM 55, 10 (October 2012), 78–87. <https://doi.org/10.1145/2347736.2347755>

Question 1.0.1

- Why the intercept b is not included in the penalty?
- How does the linear parameter \mathbf{w} behave differently in different forms of regression? Alternatively, why should we add such penalties in regression?

One-Time Graph Cut Given a graph $\mathcal{G}(V, E)$ and its adjacent matrix $\mathbf{W} = (w_{ij}) \in \{0, 1\}^{n \times n}$, the one-time graph cut partitions vertices of \mathcal{G} into two disjoint subsets A and \bar{A} such that they have minimum number of edges in between.

Definition 1.0.2: Cut

Given two disjoint subsets of vertices A and B , the cut between A and B is defined as

$$\text{cut}(A, B) = \sum_{u \in A, v \in B} w_{uv}. \quad (\text{Graph Cut})$$

By definition of graph cut, the one-time graph cut problem becomes

Find a subset $A \subset V$ such that $\text{cut}(A, \bar{A})$ attains its minimum.

We define the a vector $\mathbf{f} = (f_i)_{i=1}^n$ as

$$f_i = \begin{cases} 1 & , \text{if } i \in A \\ -1 & , \text{if } i \in \bar{A} \end{cases} \quad (1.2)$$

Define a diagonal matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ with its element as $d_{ii} = \sum_{j=1}^n w_{ij}$. Hence

$$\# \text{ Edges within } V = \sum_{i,j=1}^n w_{ij} = \sum_{i=1}^n d_{ii} = \sum_{i=1}^n f_i^2 d_{ii} = \mathbf{f}^T \mathbf{D} \mathbf{f}.$$

For any vertex i , the number of edges between vertex i and its counterpart subset \bar{i} (i.e. $i \notin \bar{i}$) can be expressed

$$\frac{1}{2} (d_{ii} - \sum_{j=1}^n f_i f_j w_{ij}).$$

So we can write the cut between A and \bar{A} as

$$\begin{aligned} \text{cut}(A, \bar{A}) &= \sum_{u \in A, v \in \bar{A}} w_{uv} = \frac{1}{2} \left(\sum_{u \in A} \sum_{v \in \bar{A}} w_{uv} + \sum_{u \in \bar{A}} \sum_{v \in A} w_{uv} \right) \\ &= \frac{1}{2} \sum_{u \in V} \sum_{v \in \bar{u}} w_{uv} = \frac{1}{4} \sum_{i=1}^n \left(d_{ii} - \sum_{j=1}^n f_i f_j w_{ij} \right) \\ &= \frac{1}{4} \left(\mathbf{f}^T \mathbf{D} \mathbf{f} - \sum_{i=1}^n \sum_{j=1}^n f_i f_j w_{ij} \right) = \frac{1}{4} \left(\mathbf{f}^T \mathbf{D} \mathbf{f} - \mathbf{f}^T \mathbf{W} \mathbf{f} \right) \\ &= \frac{1}{4} \mathbf{f}^T (\mathbf{D} - \mathbf{W}) \mathbf{f} \end{aligned}$$

Let $\mathbf{L} := \mathbf{D} - \mathbf{W}$, the one-time graph cut problem can be formulated as

$$\begin{aligned} \min_{\mathbf{f}} \quad & \mathbf{f}^T \mathbf{L} \mathbf{f} \\ \text{s.t.} \quad & \mathbf{f} \in \{-1, 1\}^n \end{aligned} \quad (\text{One-Time Graph Cut})$$

The matrix \mathbf{L} is known as the Laplacian matrix of the graph \mathcal{G} .

Question 1.0.3

- How can we generalize it to multiple cuts? Check out ratio cuts and normalized cuts.
- How can we connect the graph cuts with spectral clustering?

Optimal Transport Consider a market with m sellers and n buyers. Let the supply vector of sellers be denoted by $\mathbf{a} \in \mathbb{R}^m$, where each element represents the supply of a particular seller, and the demand vector of buyers be denoted by $\mathbf{b} \in \mathbb{R}^n$, where each element corresponds to the demand of a particular buyer. The transportation cost between seller i and buyer j is given by c_{ij} . Suppose the quantity transported from seller i to buyer j is given by p_{ij} , optimal transport seeks to determine a transport plan $\mathbf{P} = (p_{ij}) \in \mathbb{R}_+^{m \times n}$ that minimizes the total transportation cost while satisfying the supply-demand constraints.

The total cost is $\sum_{i=1}^m \sum_{j=1}^n c_{ij} p_{ij}$. The optimal transport problem can be expressed as in matrix form

$$\begin{aligned} \min_{\mathbf{P} \in \mathbb{R}_+^{m \times n}} \quad & \langle \mathbf{P}, \mathbf{C} \rangle \\ \text{s.t.} \quad & \mathbf{P} \mathbf{1}_n = \mathbf{a}, \\ & \mathbf{P}^T \mathbf{1}_m = \mathbf{b}. \end{aligned} \quad (\text{Discrete Optimal Transport})$$

where $\mathbf{1}_n$ and $\mathbf{1}_m$ denote all-one vectors of length n and m , respectively. The inner product between matrices is defined as $\langle \mathbf{P}, \mathbf{C} \rangle = \text{tr}(\mathbf{P}^T \mathbf{C}) = \sum_{i=1}^m \sum_{j=1}^n c_{ij} p_{ij}$.

Question 1.0.4

- Find you find some applications of optimal transport in economics?
- Optimal transport is a valuable tool in non i.i.d machine learning applications, such as transfer learning and multi-modal learning, for data alignment. Consider reading some of the papers that are most relevant to your interests.

Inner Products and Norms**Definition 1.0.5: Inner Product**

An **inner product** on \mathbb{R}^d is map $\langle \cdot, \cdot \rangle : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ with the following properties:

1. (**symmetry**) $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.
2. (**additivity**) $\langle \mathbf{x}, \mathbf{y} + \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle$ for any $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^d$.
3. (**homogeneity**) $\langle \lambda \mathbf{x}, \mathbf{y} \rangle = \lambda \langle \mathbf{x}, \mathbf{y} \rangle$ for any $\lambda \in \mathbb{R}$ and $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.
4. (**positive definiteness**) $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ for any $\mathbf{x} \in \mathbb{R}^d$ and $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ iff. $\mathbf{x} = \mathbf{0}$.

The **dot-product** is defined as $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^d x_i y_i$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

The definition of inner product only provide properties that it should be satisfied, but it does not involve how we calculate an inner product. Hence, the dot product is not the only possible product.

Question 1.0.6

Given a matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, we define a map $\langle \cdot, \cdot \rangle_{\mathbf{A}} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{A}} = \mathbf{x}^T \mathbf{A} \mathbf{y}.$$

Is the map an inner product? Can you define an inner product in a similar form?

Definition 1.0.7: Norm

A **norm** on \mathbb{R}^d is a function $\| \cdot \| : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying the following

1. (**nonnegativity**) $\| \mathbf{x} \| \geq 0$ for any $\mathbf{x} \in \mathbb{R}^d$.
2. (**positive homogeneity**) $\| \lambda \mathbf{x} \| = |\lambda| \| \mathbf{x} \|$ for any $\mathbf{x} \in \mathbb{R}^d$ and $\lambda \in \mathbb{R}$.
3. (**triangle inequality**) $\| \mathbf{x} + \mathbf{y} \| \leq \| \mathbf{x} \| + \| \mathbf{y} \|$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

Example.

$$\begin{aligned} \text{Euclidean norm: } \|\mathbf{x}\|_2 &= \sqrt{\sum_{i=1}^d x_i^2} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}. \\ \ell_p\text{-norms: } \|\mathbf{x}\|_p &= \sqrt[p]{\sum_{i=1}^d |x_i|^p} \quad (p \geq 1) \\ \ell_\infty\text{-norm: } \|\mathbf{x}\|_\infty &= \max_{i=1, \dots, d} |x_i|. \end{aligned}$$

Question 1.0.8

Show that $\|\mathbf{x}\|_\infty = \lim_{p \rightarrow \infty} \|\mathbf{x}\|_p$ (Hint: Reduce from right)

The Cauchy-Schwartz inequality shows the relationship between a inner product and its induced norm.

Lemma 1.0.9: Cauchy-Schwartz inequality

For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

Equality holds iff. \mathbf{x} and \mathbf{y} are linear dependent.

Similar to vector norms, we define that matrix norms as

Definition 1.0.10:

A **norm** on $\mathbb{R}^{m \times n}$ is a function $\|\cdot\| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ satisfying the following

1. (**nonnegativity**) $\|\mathbf{A}\| \geq 0$ for any $\mathbf{A} \in \mathbb{R}^{m \times n}$.
2. (**positive homogeneity**) $\|\lambda \mathbf{A}\| = |\lambda| \|\mathbf{A}\|$ for any $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\lambda \in \mathbb{R}$.
3. (**triangle inequality**) $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ for any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$.

Matrix norms can be generated from vector norms. Specifically, given two norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on \mathbb{R}^n and \mathbb{R}^m , respectively, the induced matrix norm on $\mathbb{R}^{m \times n}$ is defined as

$$\|\mathbf{A}\|_{a,b} = \max_{\mathbf{x} \in \mathbb{R}^n} \{\|\mathbf{A}\mathbf{x}\|_b : \|\mathbf{x}\|_a \leq 1\}. \quad (\text{Induced Matrix Norm})$$

If $a = b$, we simplify the notation $\|\cdot\|_{a,a}$ as $\|\cdot\|_a$.

From (Induced Matrix Norm), we can show that

$$\|\mathbf{A}\mathbf{x}\|_b \leq \|\mathbf{A}\|_{a,b} \|\mathbf{x}\|_a.$$

Example.

1. (spectral norm)

$$\|\mathbf{A}\|_2 = \max_{\mathbf{x} \in \mathbb{R}^n} \{\|\mathbf{A}\mathbf{x}\|_2 : \|\mathbf{x}\|_2 \leq 1\} = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$$

2. (1-norm: Maximum absolute column sum norm)

$$\|\mathbf{A}\|_1 = \max_{\mathbf{x} \in \mathbb{R}^n} \{\|\mathbf{A}\mathbf{x}\|_1 : \|\mathbf{x}\|_1 \leq 1\} = \max_{j=1,2,\dots,n} \sum_{i=1}^m |A_{ij}|$$

3. (∞ -norm: Maximum absolute row sum norm)

$$\|\mathbf{A}\|_\infty = \max_{\mathbf{x} \in \mathbb{R}^n} \{\|\mathbf{A}\mathbf{x}\|_\infty : \|\mathbf{x}\|_\infty \leq 1\} = \max_{i=1,2,\dots,n} \sum_{j=1}^n |A_{ij}|$$

There are several norms that are not induced matrix norms in the literature.

Example.

1. (Frobenius norm)

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2} = \sqrt{\text{tr}(\mathbf{A}^T \mathbf{A})}$$

2. (L_1 norm)

$$\|\mathbf{A}\|_1 = \sum_{i=1}^m \sum_{j=1}^n |A_{ij}|$$

3. ($L_{2,1}$ norm)

$$\|\mathbf{A}\|_{2,1} = \sum_{j=1}^n \sqrt{\sum_{i=1}^m A_{ij}^2}$$

4. (Nuclear norm)

$$\|\mathbf{A}\|_* = \text{tr}(\mathbf{A}^T \mathbf{A}) = \sum_{i=1}^{\min\{m,n\}} \sigma_i(\mathbf{A}),$$

where $\sigma_i(\mathbf{A})$ is the singular value of \mathbf{A} .

$L_{2,1}$ -norm and nuclear norm control the structure of a matrix, hence, they are favorable in high-dimensional analysis and machine learning for structure learning, e.g. Robust PCA, matrix completion.

Question 1.0.11

For a matrix \mathbf{A} of rank at most r , show that

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \|\mathbf{A}\|_* \leq \sqrt{r} \|\mathbf{A}\|_F \leq r \|\mathbf{A}\|_2.$$

Differentiability

Definition 1.0.12: Directional derivative

Let f be a real-valued function defined on a set $\mathcal{S} \subseteq \mathbb{R}^d$. Let $\mathbf{x} \in \text{int}(\mathcal{S})$ and let \mathbf{d} be any non-zero vector. If the limit

$$\lim_{t \rightarrow 0^+} \frac{f(\mathbf{x} + t\mathbf{d}) - f(\mathbf{x})}{t}$$

exists, then it is called the **directional derivative** of f at \mathbf{x} along the direction \mathbf{d} and is denoted by $f'(\mathbf{x}; \mathbf{d})$.

If the limit

$$\lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{e}_i) - f(\mathbf{x})}{t}$$

exists, we call it the i -th **partial derivative** of f at \mathbf{x} and write it as $\frac{\partial f}{\partial x_i}(\mathbf{x})$.

If all the partial derivatives exist of f at \mathbf{x} , we call the their column vector **gradient** and denote as $\nabla f(\mathbf{x})$, i.e.,

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}(\mathbf{x}), \frac{\partial f}{\partial x_2}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_d}(\mathbf{x}) \right)^T$$

A function f defined on an open set $\mathcal{U} \subseteq \mathbb{R}^d$ is called **continuously differentiable** over \mathcal{U} if all its partial derivatives exist and are continuous on \mathcal{U} for any $\mathbf{x} \in \mathbb{R}^d$. In such case, we have

$$f'(\mathbf{x}; \mathbf{d}) = \nabla f(\mathbf{x})^T \mathbf{d},$$

from that fact that

$$f'(\mathbf{x}; \mathbf{d}) = \left. \frac{\partial f(\mathbf{x} + \tau\mathbf{d})}{\partial \tau} \right|_{\tau=0} = \sum_{i=1}^d \left. \frac{\partial f(\mathbf{x} + \tau\mathbf{d})}{\partial x_i} d_i \right|_{\tau=0} = \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle.$$

Theorem 1.0.13: Taylor's Theorem

Suppose that $f : \mathcal{U} \rightarrow \mathbb{R}$ is continuously differentiable over an open set $\mathcal{U} \subseteq \mathbb{R}^d$ and for any $\mathbf{d} \in \mathbb{R}^d$ we have

$$f(\mathbf{x} + \mathbf{d}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + t\mathbf{d})^T \mathbf{d},$$

for some $t \in (0, 1)$ and any $\mathbf{x} \in \mathcal{U}$

Proof. Let $g(\tau) = f(\mathbf{x} + \tau\mathbf{d})$. By the mean-value theorem we have

$$g(1) - g(0) = g'(t),$$

for some $t \in (0, 1)$.

Since

$$g'(t) = f'(\mathbf{x} + \tau\mathbf{d})|_{\tau=t} = \nabla f(\mathbf{x} + t\mathbf{d})^T \mathbf{d}$$

by definition of gradient we complete the proof. \square

Proposition 1.0.14

Suppose that $f : \mathcal{U} \rightarrow \mathbb{R}$ is continuously differential over an open set $\mathcal{U} \subseteq \mathbb{R}^d$. Then

$$\lim_{\mathbf{d} \rightarrow \mathbf{0}} \frac{f(\mathbf{x} + \mathbf{d}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T \mathbf{d}}{\|\mathbf{d}\|} = 0,$$

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

Proof. By assumption, for any unit vector $\mathbf{v} \in \mathbb{R}^n$ we have

$$f'(\mathbf{x}; \mathbf{v}) = \nabla f(\mathbf{x})^T \mathbf{v},$$

Hence

$$\lim_{t \rightarrow 0^+} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x}) - t\nabla f(\mathbf{x})^T \mathbf{v}}{t} = 0,$$

Let $\mathbf{d} = t\mathbf{v}$. The above equality is equivalent to

$$\lim_{\mathbf{d} \rightarrow \mathbf{0}} \frac{f(\mathbf{x} + \mathbf{d}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T \mathbf{d}}{\|\mathbf{d}\|} = 0,$$

\square

The proposition implies that

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + o(\|\mathbf{x} - \mathbf{y}\|)$$

for any $\mathbf{x}, \mathbf{y} \in \mathcal{U}$, where $o(\|\mathbf{x} - \mathbf{y}\|)$ is a higher order infinitesimal of $\|\mathbf{x} - \mathbf{y}\|$.

The (i, j) -th partial derivatives of f at $\mathbf{x} \in \mathbb{R}^d$ is defined by

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) = \frac{\partial}{\partial x_i} \left(\frac{\partial f}{\partial x_j} \right)(\mathbf{x}).$$

We call a real-valued function twice continuously differentiable if all its partial derivatives exist and are continuous. In such case, we have

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) = \frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}).$$

The Hessian matrix of f at \mathbf{x} is denoted by

$$\nabla^2 f(\mathbf{x}) = \left(\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) \right)_{i,j}.$$

Theorem 1.0.15: Taylor's Theorem (Continued)

If f is twice continuously differentiable, we have

$$\nabla f(\mathbf{x} + \mathbf{d}) = \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{x} + \tau \mathbf{d}) \mathbf{d} d\tau$$

and

$$f(\mathbf{x} + \mathbf{d}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x} + t\mathbf{d}) \mathbf{d}$$

for some $t \in (0, 1)$.

The gradient can also be defined from Fréchet differentiability.

Definition 1.0.16: Fréchet differentiability

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$, and $\mathbf{x} \in \text{int}(\mathcal{S})$. Then function is said to be differentiable at \mathbf{x} if there exists a vector \mathbf{g} such that

$$\lim_{\mathbf{d} \rightarrow \mathbf{0}} \frac{f(\mathbf{x} + \mathbf{d}) - f(\mathbf{x}) - \mathbf{g}^T \mathbf{d}}{\|\mathbf{d}\|} = 0.$$

The unique vector \mathbf{g} satisfying the equality is called gradient f at \mathbf{x} and it is denoted by $\nabla f(\mathbf{x})$.

The definition of the gradient is consistent with the one we defined earlier.

Optimality Conditions for Unconstrained Optimization**Definition 1.0.17: Global and local minimum**

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over a set $\mathcal{S} \subseteq \mathbb{R}^d$. Then

1. We call \mathbf{x}^* a (strict) global minimum point of f over \mathcal{S} if $f(\mathbf{x}^*) \leq (<) f(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{S}$.
2. We call \mathbf{x}^* a (strict) local minimum point of f over \mathcal{S} if there exists $r > 0$ such that $f(\mathbf{x}^*) \leq (<) f(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, r)$.

The set of global minimum point of f over \mathcal{S} is represented by

$$\arg \min_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x}).$$

The open ball $\mathcal{B}(\mathbf{x}^*, r)$ denotes the set $\{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}^* - \mathbf{x}\| \leq r\}$.

Theorem 1.0.18: First order optimality condition for local minimum points

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$. Suppose that $\mathbf{x}^* \in \text{int}(\mathcal{S})$ is a local minimum point and that $\nabla f(\mathbf{x}^*)$ exists. Then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Proof. By assumption, $\frac{\partial f(\mathbf{x}^*)}{\partial x_i}$ exists for all $i = 1, \dots, d$. Define $g(t) = f(\mathbf{x}^* + te_i)$ and since $\mathbf{x}^* \in \text{int}(\mathcal{S})$, there exists $r > 0$ such that $t \in (-r, r)$. We have $t = 0$ is a local minimum point of $g(t)$ by local optima of \mathbf{x}^* . As $g'(0)$ exists, hence $g'(0) = \frac{\partial f(\mathbf{x}^*)}{\partial x_i}$. Such claim holds for all i , we have $\nabla f(\mathbf{x}^*) = \mathbf{0}$. \square

Definition 1.0.19: Stationary points

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$. Suppose that $\mathbf{x}^* \in \text{int}(\mathcal{S})$ and that f is differentiable over some neighborhood of \mathbf{x}^* . Then \mathbf{x}^* is called a **stationary point** of f if $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Theorem 1.0.20: Necessary second order optimality conditions

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$. Suppose that f is twice continuously differentiable over \mathcal{S} and \mathbf{x}^* is a stationary point. Then the following hold

- (a) If \mathbf{x}^* is a local minimum point of f over \mathcal{S} , then $\nabla^2 f(\mathbf{x}^*) \succcurlyeq \mathbf{0}$;
- (b) If \mathbf{x}^* is a local maximum point of f over \mathcal{S} , then $\nabla^2 f(\mathbf{x}^*) \preccurlyeq \mathbf{0}$;

Theorem 1.0.21: Sufficient second order optimality condition

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$. Suppose that f is twice continuously differentiable over \mathcal{S} and \mathbf{x}^* is a stationary point. Then the following hold:

- (a) If $\nabla^2 f(\mathbf{x}^*) \succ \mathbf{0}$, then \mathbf{x}^* is a strict local minimum point of f over \mathcal{S} .
- (b) If $\nabla^2 f(\mathbf{x}^*) \prec \mathbf{0}$, then \mathbf{x}^* is a strict local maximum point of f over \mathcal{S} .

Definition 1.0.22: Saddle point

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$. Suppose that f is continuously differentiable over \mathcal{S} . A stationary point \mathbf{x}^* is called a **saddle point** of f over \mathcal{S} if its not a local optimal point of f over \mathcal{S}

Theorem 1.0.23: Sufficient condition for a saddle point

Let $f : \mathcal{S} \rightarrow \mathbb{R}$ be defined over $\mathcal{S} \subseteq \mathbb{R}^d$. Suppose that f is twice continuously differentiable over \mathcal{S} and \mathbf{x}^* is a stationary point. If $\nabla^2 f(\mathbf{x}^*)$ is an indefinite matrix, then \mathbf{x}^* is a saddle point f over \mathcal{S} .

Attainable of Optimal

When we try find a global minimum/maximum point of a function, we need to be sure of its existence. Weierstrass theorem provides such guarantee for a continuous function over a compact set.

Theorem 1.0.24: Weierstrass theorem

Let f be a continuous function defined over a nonempty and compact set $\mathcal{C} \subseteq \mathbb{R}^d$. Then global optimal points over \mathcal{C} are attainable.

The compact property is not necessary if the function f satisfies coerciveness.

Definition 1.0.25: Coerciveness

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous function defined over \mathbb{R}^d . Then the function f is called *coercive* if

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} f(\mathbf{x}) = \infty.$$

A coercive function always achieves its global minimum point over any closed set.

Theorem 1.0.26: Attainment under coerciveness

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous and coercive function and let $\mathcal{S} \subseteq \mathbb{R}^d$ be a nonempty closed set. Then f has a global minimum point over \mathcal{S} .

Proof. Given $\mathbf{x}_0 \in \mathcal{S}$, there exists $M > 0$ such that

$$f(\mathbf{x}) > f(\mathbf{x}_0) \text{ for any } \mathbf{x} \text{ satisfying } \|\mathbf{x}\| > M,$$

by the coerciveness of f . Hence the level set $\mathcal{L}_{\mathbf{x}_0} = \{\mathbf{x} \in \mathbb{R}^d | f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ is bounded and also is closed by the continuity of f , i.e., $\mathcal{L}_{\mathbf{x}_0}$ is a non-empty and compact set. And so is $\mathcal{S} \cap \mathcal{L}_{\mathbf{x}_0} = \{\mathbf{x} \in \mathcal{S} | f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. From Weierstrass theorem we know that there exists a global minimum point of f over $\mathcal{S} \cap \mathcal{L}_{\mathbf{x}_0}$. This point is also a global minimum point of f over \mathcal{S} .

□