# Week 3 Tutorial 1

**(1.1)** *Separating Hyperplane Theorem*

We recall the separating hyperplane Theorem 3.10 discussed in Lecture 1 Week 2.

Construct counter examples if the set $C$

- is not convex.

- is not closed.

**(1.2)** *Convex cones*

- Let $\mathbf{S}^n$ denote the set of symmetric $n \times n$ matrices, that is

$$\mathbf{S}^n = \{\mathbf{A} \in \mathbf{R}^{n \times n} \colon \mathbf{A} = \mathbf{A}^T\}.$$

and by $\mathbf{S}_+^n$ the set of symmetric positive semi-definite matrices.
What is the dimension of $\mathbf{S}^n$. Show that $\mathbf{S}^n$ is a convex cone.

- The second order cone (also known as the ice cream cone) is the norm cone defined for the Euclidean norm, that is

$$C = \{(\mathbf{x}, t) \in \mathbb{R}^{n+1} \colon \|\mathbf{x}\|_2 \leq t\}$$

$$= \{\begin{pmatrix} \mathbf{x} \\ t \end{pmatrix} \colon \begin{pmatrix} \mathbf{x} \\ t \end{pmatrix}^T \begin{pmatrix} I & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ t \end{pmatrix} \leq 0, t \geq 0\}.$$

Here $I$ denotes the $n \times n$ identity matrix.
Show that $C$ is indeed a cone. Plot the cone in $\mathbb{R}^3$, that is the set $\{(x_1, x_2, t) \colon (x_1^2 + x_2^2)^{\frac{1}{2}} \leq t\}$.

**(1.3)** *Operations that preserve convexity of functions*

- <u>Affine mappings</u>: Let $f \colon \mathbb{R}^n \to \mathbb{R}$, $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^n$. Define $g \colon \mathbb{R}^m \to \mathbb{R}$ by

$$g(\mathbf{x}) = f(\mathbf{Ax} + \mathbf{b})$$

with dom $g = \{\mathbf{x} \in \mathbb{R}^m \colon \mathbf{Ax} + \mathbf{b} \in \text{ dom } f\}$. Show that if $f$ is convex, so is $g$.

- <u>Pointwise maximum:</u> Let $f_1$ and $f_2$ be convex functions and define their pointwise maximum as

$$f(x) = \max\{f_1(x), f_2(x)\}$$

with dom $f = \text{dom } f_1 \cap \text{dom} f_2$ also convex. Show that $f$ is convex.

- <u>Scalar composition</u>: Let $h \colon \mathbb{R} \to \mathbb{R}$ and $g \colon \mathbb{R} \to \mathbb{R}$ and $f = h \circ g \colon \mathbb{R} \to \mathbb{R}$ be defined as

$$f(x) = h(g(x)) \qquad \text{dom} f = \{x \in \text{dom } g \colon g(x) \in \text{dom } h\}.$$

Furthermore assume that $h$ and $g$ are twice differentiable.
Discuss under which conditions on $h$ and $g$ the function $f$ is convex.

**(1.4)**   *Stochastic gradient descent (SDG)*

A common situation in machine learning is that the objective function is of the form

$$\min_x \frac{1}{N} \sum_{i=1}^{N} f_i(x),$$

where $f_i$ is an individual loss function associated to the particular data point $x_i$ and $N \in \mathbb{N}$. In gradient descent a full step iterates $\mathbf{x}_k = (x_{k,1}, \ldots x_{k,n})$, $k = 1, 2, 3, \ldots$ would be updated according to

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \frac{\alpha}{N} \sum_{i=1}^{N} \nabla f_i(\mathbf{x}_{k-1})$$

where $f_i : \mathbb{R}^n \to \mathbb{R}$, $i = 1, \ldots N$.

*Question for students:*   What is the computational cost of each iterate? Why does it become computationally costly if you have $N = 10^6$ data points?

In SDG we update iterates $\mathbf{x}_k$ based on the descent in one component only, that is

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \alpha_k \nabla f_{i_k}(\mathbf{x}_{k-1}) \qquad \text{for } k = 1, 2, \ldots$$

where $i_k \in \{1, 2 \ldots N\}$ is a randomly chosen index at iteration $k$.

A common technique in SDG is mini-batching, where one chooses a random subset $I_k \subseteq \{1, 2, \ldots n\}$ with size $|I_k| = M \ll N$. Hence we have the following update rule

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \frac{\alpha_k}{M} \sum_{i \in I_k} \nabla f_{i_k}(\mathbf{x}_{k-1})$$

*Question for students* What is the computational complexity of SDB and mini-batch SDG?