# MATH 565: Lecture 4 (01/22/2026)

Today :
* illustration of gradient descent
* local optimality in d-dimensions
* convex functions

We saw gradient descent for $f(x) = x^2 \sin(x) + x$, with
$$f'(x) = 2x\sin x + x^2\cos(x) + 1,$$
and $f''(x) = (2-x^2)\sin(x) + 4x\cos(x)$

## A 2D Extension

Consider $G(x,y) = f(x) + f(y)$, which is a separable function, i.e., it has no terms involving both $x$ and $y$.

$$\Rightarrow \nabla G = \begin{bmatrix} \frac{\partial G}{\partial x} \\ \frac{\partial G}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x\sin(x) + x^2\cos(x) + 1 \\ 2y\sin(y) + y^2\cos(y) + 1 \end{bmatrix}$$

Gradient descent now implements $\bar{w} \leftarrow \bar{w} - \alpha \nabla G(\bar{w})$.
Check Python notebook on the course web page...

We consider one final example with a nonseparable function:
Let $d(x,y) = x^2 \sin(y)$, which gives
$$\nabla d = \begin{bmatrix} 2x\sin y \\ x^2\cos y \end{bmatrix}.$$

The point of these exercises was to demonstrate that optimization using gradient descent in high dimensions could become quite complex quickly...

# More on local optimality in $d$ dimensions

The sufficient conditions (first and second order) for local optimality in $d$ dimensions (given in Lemma 2 in lecture 3) can be qualified further to specify more nuanced cases for the second order condition.

Recall : $\nabla J = \bar{0}$ (first order optimality condition)

$HJ > 0$, i.e., $H$ is PD (positive definite)
$\longrightarrow$ (second order optimality condition)

(Result: a real symmetric matrix $H$ is PD $\Longleftrightarrow$ all its eigenvalues are $> 0$).

# Second order optimality conditions

Let $\nabla J(\bar{w}_0) = \bar{0}$.

1. If $HJ(\bar{w}_0) > 0$, then $\bar{w}_0$ is a local minimum (Lemma 2)
   ($H$ is PD)

2. If $HJ(\bar{w}_0) < 0$, then $\bar{w}_0$ is a local maximum.
   ($H$ is ND, negative definite)

3. If $HJ(\bar{w}_0)$ is indefinite, then $\bar{w}_0$ is a saddle point.

4. If $HJ(\bar{w}_0) \succeq 0$ or $HJ(\bar{w}_0) \preceq 0$, then the test is inconclusive. $\leftarrow$ positive/negative semidefinite
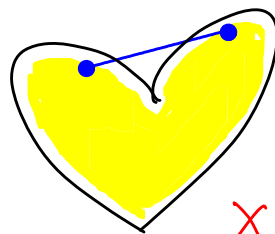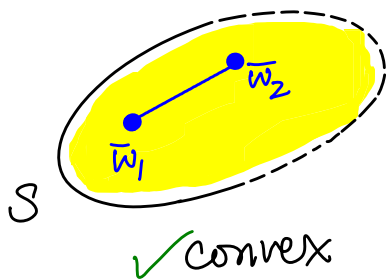
The presence of local optima can make the task of finding a global optimum hard. We now study a class of functions for which any local optima are also guaranteed to be global optima — convex functions.

# Convex Sets and Functions

We first define convex sets.

**Def** A set $S \subseteq \mathbb{R}^d$ is **convex** if $\forall \bar{w}_1, \bar{w}_2 \in S$,
$$\lambda \bar{w}_1 + (1-\lambda)\bar{w}_2 \in S \quad \forall \lambda \in [0,1].$$

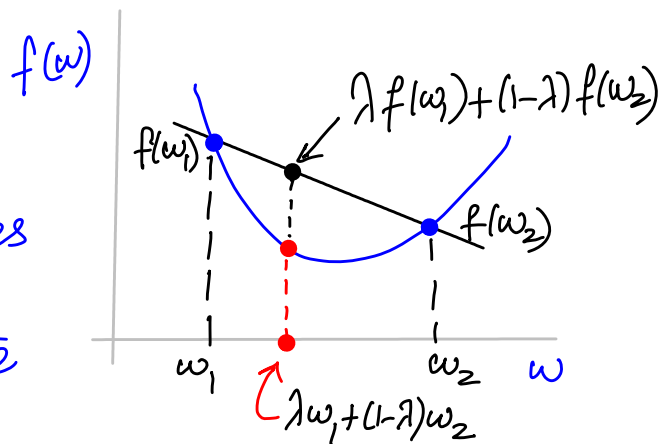In words, the line segment connecting $\bar{w}_1$ and $\bar{w}_2$ lies in $S$.



S
✓ convex

✗ not convex

**Def** A function $f: \Omega \to \mathbb{R}$ is a **convex function** for convex set $\Omega \subseteq \mathbb{R}^d$ if
$$f\left(\lambda \bar{w}_1 + (1-\lambda)\bar{w}_2\right) \leq \lambda f(\bar{w}_1) + (1-\lambda) f(\bar{w}_2) \quad\quad\quad (*)$$

$$\forall \bar{w}_1, \bar{w}_2 \in \Omega, \quad \forall \lambda \in [0,1].$$

In 1D, the value of the function at a convex combination of two points lies not above the corresponding convex combination of the function values at the two points.



$f(w)$

$f(w_1)$

$\lambda f(w_1) + (1-\lambda) f(w_2)$

$f(w_2)$

$w_1$    $w_2$   $w$

$\lambda w_1 + (1-\lambda) w_2$

**Proposition 3** If $f: \Omega \to \mathbb{R}$ is convex, then

$$f\left(\sum_{i=1}^{k} \lambda_i \overline{w}_i\right) \leq \sum_{i=1}^{k} \lambda_i f(w_i) \text{ for } w_i \in \Omega, \; 0 \leq \lambda_i \leq 1, \; \sum_{i=1}^{k} \lambda_i = 1.$$

**Proof** We use induction on $k$. → $k=1$ is trivial: $\lambda_1 = 1$ and $f(1\overline{w}_1) = 1 f(\overline{w}_1)$

<u>base case</u> $k=2$ is $(\ast)$ from the definition of a convex function. Assume result holds for $k$, and consider statement for $k+1$:

We want to show that

$$f\left(\sum_{i=1}^{k+1} \lambda_i \overline{w}_i\right) \leq \sum_{i=1}^{k+1} \lambda_i f(\overline{w}_i) \text{ holds.}$$

Assume $0 < \lambda_{k+1} < 1$ ($\lambda_{k+1} = 0$ or $1 \Rightarrow$ trivial)

note: $\sum_{i=1}^{k} \lambda_i = 1 - \lambda_{k+1}$

$$\Rightarrow f\left(\sum_{i=1}^{k} \lambda_i \overline{w}_i + \lambda_{k+1} \overline{w}_{k+1}\right) = f\left(\frac{1-\lambda_{k+1}}{(1-\lambda_{k+1})} \sum_{i=1}^{k} \lambda_i \overline{w}_i + \lambda_{k+1} \overline{w}_{k+1}\right)$$

$$= f\left((1-\lambda_{k+1})\left[\sum_{i=1}^{k} \lambda_i' \overline{w}_i\right] + \lambda_{k+1} \overline{w}_{k+1}\right) \text{ for } \lambda_i' = \frac{\lambda_i}{1-\lambda_{k+1}} \quad \to 0 \leq \lambda_i' \leq 1$$

$$\leq (1-\lambda_{k+1}) f\left(\sum_{i=1}^{k} \lambda_i' \overline{w}_i\right) + \lambda_{k+1} f(\overline{w}_{k+1}) \text{ by result for } k=2 \text{ (base case)}$$

But $\sum_{i=1}^{k} \lambda_i' = 1$, which gives that the expression is

$$\leq (1-\lambda_{k+1}) \sum_{i=1}^{k} \lambda_i' f(\overline{w}_i) + \lambda_{k+1} f(\overline{w}_{k+1}) \text{ by induction assumption}$$

$$= \sum_{i=1}^{k+1} \lambda_i f(\overline{w}_i) \quad \text{as } (1-\lambda_{k+1})\lambda_i' = \lambda_i$$

$\square$

# Useful Properties of Convex Functions

1. If $f_1, f_2$ are convex functions, then $g = f_1 + f_2$ is convex.
   (sum of convex functions is convex)

2. If $f_1, f_2$ are convex functions, then $g = \max(f_1, f_2)$ is convex.
   (max of convex functions is convex)

3. If $f: \Omega \to \mathbb{R}$ is convex and $f(\bar{w}) \geq 0$, then $g = f^2 = [f(\bar{w})]^2$ is convex.

4. If $f: \mathbb{R} \to \mathbb{R}$ is convex and $g: \mathbb{R}^d \to \mathbb{R}$ is linear, then $h: \mathbb{R}^d \to \mathbb{R}$ defined as $h = f(g(\bar{w}))$ is convex.

**Proof:** $f(\lambda w_1 + (1-\lambda) w_2) \leq \lambda f(w_1) + (1-\lambda) f(w_2)$ ————— (1)

Let $g(\bar{w}) = \bar{w}^T \bar{x} + b$ linear function.

Then, $h(\bar{w}) = f(g(\bar{w})) = f(\bar{w}^T \bar{x} + b)$.

$$\Rightarrow h(\lambda \bar{w}_1 + (1-\lambda)\bar{w}_2) = f\left(g(\lambda \bar{w}_1 + (1-\lambda)\bar{w}_2)\right)$$
$$= f\left(\bar{x}^T(\lambda \bar{w}_1 + (1-\lambda)\bar{w}_2) + b\right) \quad \to (\lambda + (1-\lambda))b$$
$$= f\left(\lambda \underbrace{(\bar{x}^T \bar{w}_1 + b)}_{w_1} + (1-\lambda)\underbrace{(\bar{x}^T \bar{w}_2 + b)}_{w_2}\right)$$
$$\leq \lambda f(\bar{x}^T \bar{w}_1 + b) + (1-\lambda) f(\bar{x}^T \bar{w}_2 + b) \quad \to \text{by (1)}$$
$$= \lambda h(\bar{w}_1) + (1-\lambda) h(\bar{w}_2)$$

□

This form of composition of functions is used in many machine learning contexts, e.g., in deep learning networks, a node may take the input from $k$ other nodes, combine them in a linear or affine format, and then apply a convex transformation to that combination.

The order in which the composition is taken in the above statement is important! For instance, if we consider $f(x)=x^2$, which is convex, and $g(x)=-x$, which is linear, and then consider

$$j(x)=g(f(x))= -x^2,$$

which is concave!