# WEEK 9: REVISION | FINAL EXAM

# CONTENTS

1. Properties of Convex Functions

2. Applications of Optimization in Machine Learning

3. Revisiting Constrained Optimization

4. Relation between Primal and Dual Problem, KKT Conditions

5. KKT conditions continued

# 1. PROPERTIES OF CONVEX FUNCTIONS

Necessary and sufficient conditions for optimality of convex functions

Goal: $\min\limits_{x} f(x)$

## Theorem

Let $f$ be a differentiable and convex function from $\mathbb{R}^d \to \mathbb{R}$, $x^* \in \mathbb{R}^d$ is a global minimum of $f$ **if and only if** $\nabla f(x^*) = 0$.

If $\exists x^*$ s.t $\nabla f(x^*) = 0$ $\Rightarrow$ $x^*$ is a global minima.

By definition of Convexity
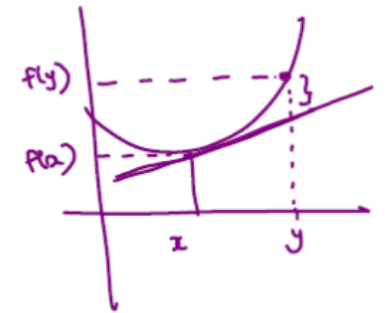
$f(y) \geq f(x) + \nabla f(x)^T (y-x)$    $\forall x, y$.

$\Rightarrow$ $f(y) \geq f(x^*) + \underbrace{\nabla f(x^*)^T (y-x)}_{0}$   $\forall y$

       $(\because \nabla f(x^*) = 0$ by assumption$)$

$\Rightarrow$    $\boxed{f(y) \geq f(x^*) \quad \forall y}$    $\Rightarrow$ $x^*$ is a global minima!

# 1. PROPERTIES OF CONVEX FUNCTIONS

If $f : \mathbb{R}^d \to \mathbb{R}, g : \mathbb{R}^d \to \mathbb{R}$ are both convex functions, then $\underline{f(x) + g(x)}$ is a convex function

Proof: Fix $\left( \lambda \in [0,1] \right)$

$$h\left( \lambda x + (1-\lambda)y \right) = f\left( \lambda x + (1-\lambda)y \right) + g\left( \lambda x + (1-\lambda)y \right) \quad [\text{by defn of } h]$$

$$\leq \lambda f(x) + (1-\lambda) f(y) + \lambda g(x) + (1-\lambda) g(y) \quad [\text{by convexity of } f \text{ and } g]$$

$$= \lambda \left( f(x) + g(x) \right) + (1-\lambda) \left( f(y) + g(y) \right)$$

$$= \lambda h(x) + (1-\lambda) h(y)$$

Sums of convex functions is convex

# 1. PROPERTIES OF CONVEX FUNCTIONS

Let $f : \mathbb{R} \to \mathbb{R}$ is a convex and non-decreasing function and $g : \mathbb{R}^d \to \mathbb{R}$ be a convex function, then their composition $h = f(g(x))$ is also a convex function.

# 1. PROPERTIES OF CONVEX FUNCTIONS

Let $f : \mathbb{R} \to \mathbb{R}$ is a convex function and $g : \mathbb{R}^d \to \mathbb{R}$ be a linear function, then their composition $h = f(g(x))$ is also a convex function.

Proof:

Fix $\lambda \in [0,1]$.

$$h\left(\lambda x + (1-\lambda)y\right) \quad = \quad f\left(g\left(\lambda x + (1-\lambda)y\right)\right)$$

$$= f\left(\lambda g(x) + (1-\lambda)g(y)\right) \quad \left[\text{by linearity of } g\right]$$

$$\leq \lambda f(g(x)) + (1-\lambda) f(g(y)) \quad = \lambda h(x) + (1-\lambda)h(y)$$

$$\left[\text{Convexity of } f\right]$$

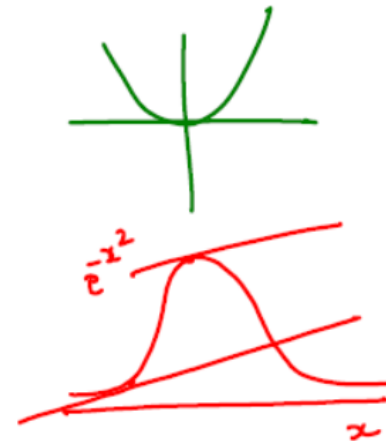# 1. PROPERTIES OF CONVEX FUNCTIONS

In general, if $f$ and $g$ are both convex functions, then $h = fog$ may not be convex function.

$$g(x) = x^2 \quad \rightarrow \quad \text{Convex}$$

$$f(x) = e^{\boxed{-x}} \quad \rightarrow \quad \text{Convex}.$$

$$fog(x) = e^{-x^2} \quad \cdot \text{ is not Convex}. \qquad e^{-x^2}$$

**Note:** $g$ is concave if and only if $f = -g$ is convex.

# 2. APPLICATIONS OF OPTIMIZATION IN ML

**Linear Regression:**

Training data $\rightarrow X_1, X_2, ..., X_n$ with corresponding outputs $y_1, y_2, ..., y_n$, where $X_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}, \forall i$.

Performance measure : Sum of Squares error.

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^{n} \left( w^T x_i - y_i \right)^2$$

$\hookrightarrow$ error of $x_i$ made by $w$.

$f(w)$

Specific Goal of linear regression

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^{n} \left( w^T x_i - y_i \right)^2$$

Gradient of the sum of squares error

$$\nabla f(w) = (X^T X) w - X^T y$$

Analytical or closed form solution of coefficients $w^*$ of a linear regression model

$$w^* = (X^T X)^{-1} X^T y$$

# 2. APPLICATIONS OF OPTIMIZATION IN ML

→ GD will reach GM if GM ⇒ LM

In linear regression, the <u>gradient descent approach</u> avoids the inverse computation by iteratively updating the weights.

$$w^{t+1} = w^t - \eta_t \nabla f(w^t))$$
$$w^{t+1} = w^t - \eta_t((X^T X)w^t - X^T y)$$

local

global

$10,000 \times 10,000$

**Stochastic gradient descent:**

- Computes approximation of gradient to make gradient computation faster (because in GD $X^T X$ will use entire dataset).
- Samples a small set of data points at random for every iteration to compute the gradient.

$$\frac{1}{T} \sum_{t=1}^{T} w_t \rightarrow w^*$$

# 3. CONSTRAINED OPTIMIZATION

Consider the constrained optimization problem as follows:

$$\min_x f(x)$$

$$\text{subject to } h(x) \leq 0$$

**Lagrangian function:**

$$L(x, \lambda) = f(x) + \lambda h(x)$$

vector   Scalar

$$\begin{array}{l} \min_x \ f(x) \\ \\ s.t \ \ h(x) \leq 0 \end{array}$$

$\equiv$

$$\min_x \left[ \max_{\lambda \geq 0} L(x, \lambda) \right]$$

*primal*

# 4. RELATION BETWEEN PRIMAL AND DUAL PROBLEM

$x^*$

$$\min_{x} \left[ \max_{\lambda \geq 0} \; \mathcal{L}(x, \lambda) \right]$$

Primal

$$\max_{\lambda \geq 0} \; \min_{x} \; \mathcal{L}(x, \lambda)$$

DUAL

$\lambda^*$

$f(x)$

$x^\#$

duality

$x^*$

$g(\lambda)$

| Weak Duality | Strong Duality |
|---|---|
| $$g(\lambda^*) \leq f(x^*)$$ | If f and g are convex functions. $$g(\lambda^*) = f(x^*)$$ |

# 5. KARUSH-KUHN-TUCKER CONDITIONS

Consider the optimization problem with multiple equality and inequality constraints as follows:

$$\min_x f(x)$$
$$\text{subject to}$$
$$h_i(x) \leq 0, \forall i = 1, ..., m$$
$$l_j(x) = 0, \forall j = 1, ..., n$$

The Lagrangian function is expressed as follows:

$$L(x, u, v) = f(x) + \sum_{i=1}^{m} u_i h_i(x) + \sum_{j=1}^{n} v_j l_j(x)$$

## Karush-Kuhn-Tucker Conditions:

Stationarity
$$\nabla f(x) + \sum_{i=1}^{n} u_i \nabla h(x) + \sum_{j=1}^{m} v_j \nabla l(x) = 0$$

Complementary slackness $\quad u_i h_i = 0 \quad \forall i$

Primal feasibility $\quad h_i(x) \leq 0 \quad \forall i$

Dual feasibility $\quad u_i \geq 0 \quad \forall i$

# SOME SOLVED PROBLEMS

# Properties of convex functions

https://www.geogebra.org/m/esqcd4he

properties of convex function



| | |
|---|---|
| 🟢 | $f(x) = 3 + x^2 + x$ |
| 🔴 | $g(x) = e^{-x}$ |
| 🔵 | $h(x) = f(g(x))$ <br> $\rightarrow 3 + (e^{-x})^2 + e^{-x}$ |
| ⚪ | $h1(x) = g(f(x))$ <br> $\rightarrow e^{-(3+x^2+x)}$ |
| 🟠 | $h2(x) = f(x)\,g(x)$ <br> $\rightarrow (3 + x^2 + x)\,e^{-x}$ |

# Given below is a set of data points and their labels.

| X | y |
|---|---|
| [1,0] | 1.5 |
| [2,1] | 2.9 |
| [3,2] | 3.4 |
| [4,2] | 3.8 |
| [5,3] | 5.3 |

How to find the optimal $w*$ using the analytical method?

Let us use Gradient descent optimization.

optimal $w^* = (X^T X)^{-1}(X^T y)$

$$X \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 3 & 2 \\ 4 & 2 \\ 5 & 3 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 55 & 31 \\ 31 & 18 \end{bmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 59.2 \\ 33.2 \end{pmatrix}$$

$$X^T y = \begin{bmatrix} 59.2 \\ 33.2 \end{bmatrix}$$

$$w^* = \begin{bmatrix} 1.255 \\ -0.317 \end{bmatrix}$$

Given $w^1 = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}$

Gradient $\nabla f(w) = (X^T X)w - X^T y$

$$\nabla f(w) = \begin{bmatrix} -50.6 \\ -28.3 \end{bmatrix}$$

update equation: $w^2 = w^1 - \eta_t \nabla f(w^1)$

$$w^2 = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} - 0.1 \begin{bmatrix} -50.6 \\ -28.3 \end{bmatrix}$$

$$w^2 = \begin{bmatrix} 5.16 \\ 2.93 \end{bmatrix}$$

$$minimize \quad 3x_1 + x_2 \rightarrow f \qquad \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix}$$

$$subject\ to$$

$$u_1 \qquad x_1 - x_2 + 4 \leq 0 \rightarrow g_1 \quad \nabla g_1$$

$$u_2 \qquad -3x_1 + 2x_2 + 10 \leq 0 \rightarrow g_2 \quad \nabla g_2$$

**Stationarity conditions** $\quad 3 + u_1 - 3u_2 = 0$ — ①

$$1 - u_1 + 2u_2 = 0$$ — ②

**Complementary slackness conditions**

$$u_1(x_1 - x_2 + 4) = 0$$ — ③

$$u_2(-3x_1 + 2x_2 + 10) = 0$$ — ④

**Primal feasibility conditions**

$$x_1 - x_2 + 4 \leq 0$$

$$-3x_1 + 2x_2 + 10 \leq 0$$

**Dual feasibility conditions**

$$u_1, u_2 \geq 0$$

---

$$\frac{u_2}{=0} \neq 0$$

$$u_1 = 0 \quad \boxed{\begin{array}{c|c} \times & \nearrow \\ \hline \times & \checkmark \end{array}}$$

$$\neq 0$$

Case (i) $\quad u_1 = 0 \ ; \ u_2 = 0$

$$3 \neq 0$$

$$1 \neq 0$$

Case (ii) $\quad u_1 = 0, \ u_2 \neq 0$

$$3 - 3u_2 = 0 \qquad u_2 = 1 \checkmark$$

$$1 + 2u_2 = 0 \qquad\qquad -3x_1 + 2x_2 + 10 = 0$$

Case (iii) $\quad u_1 \neq 0, \ u_2 = 0$

$$3 + u_1 = 0$$

$$1 - u_1 = 0 \qquad u_1 = 1 \checkmark \qquad x_1 - x_2 + 4 = 0$$

$$\text{minimize} \quad 3x_1 + x_2$$

$$\text{subject to}$$

$$x_1 - x_2 + 4 \leq 0$$

$$-3x_1 + 2x_2 + 10 \leq 0$$

**Stationarity conditions**

$$3 + u_1 - 3u_2 = 0 \quad \text{——} \quad \textcircled{1}$$

$$1 - u_1 + 2u_2 = 0 \quad \text{——} \quad \textcircled{2}$$

$$\leftharpoondown 4$$

**Complementary slackness conditions**

$$\overset{\neq 0}{u_1}(x_1 - x_2 + 4) = 0 \quad \text{——} \quad \textcircled{3}$$

$$\overset{\neq 0}{u_2}(-3x_1 + 2x_2 + 10) = 0 - \textcircled{4}$$

**Primal feasibility conditions**

$$x_1 - x_2 + 4 \leq 0$$

$$-3x_1 + 2x_2 + 10 \leq 0$$

**Dual feasibility conditions**

$$u_1, u_2 \geq 0 \quad \checkmark$$

$$u_1 \neq 0 \qquad u_2 \neq 0$$

$$4 - u_2 = 0 \qquad \therefore \boxed{u_2 = 4 \,, \qquad u_1 = 9}$$

$$x_1 - x_2 + 4 = 0 \quad \Rightarrow \quad 2x_1 - 2x_2 + 8 = 0$$

$$-3x_1 + 2x_2 + 10 = 0$$

$$-x_1 + 18 = 0 \qquad \boxed{\therefore x_1 = 18}$$

$$\boxed{\therefore x_2 = 22}$$

$$\text{minimum} \quad f(x_1, x_2) = 76$$