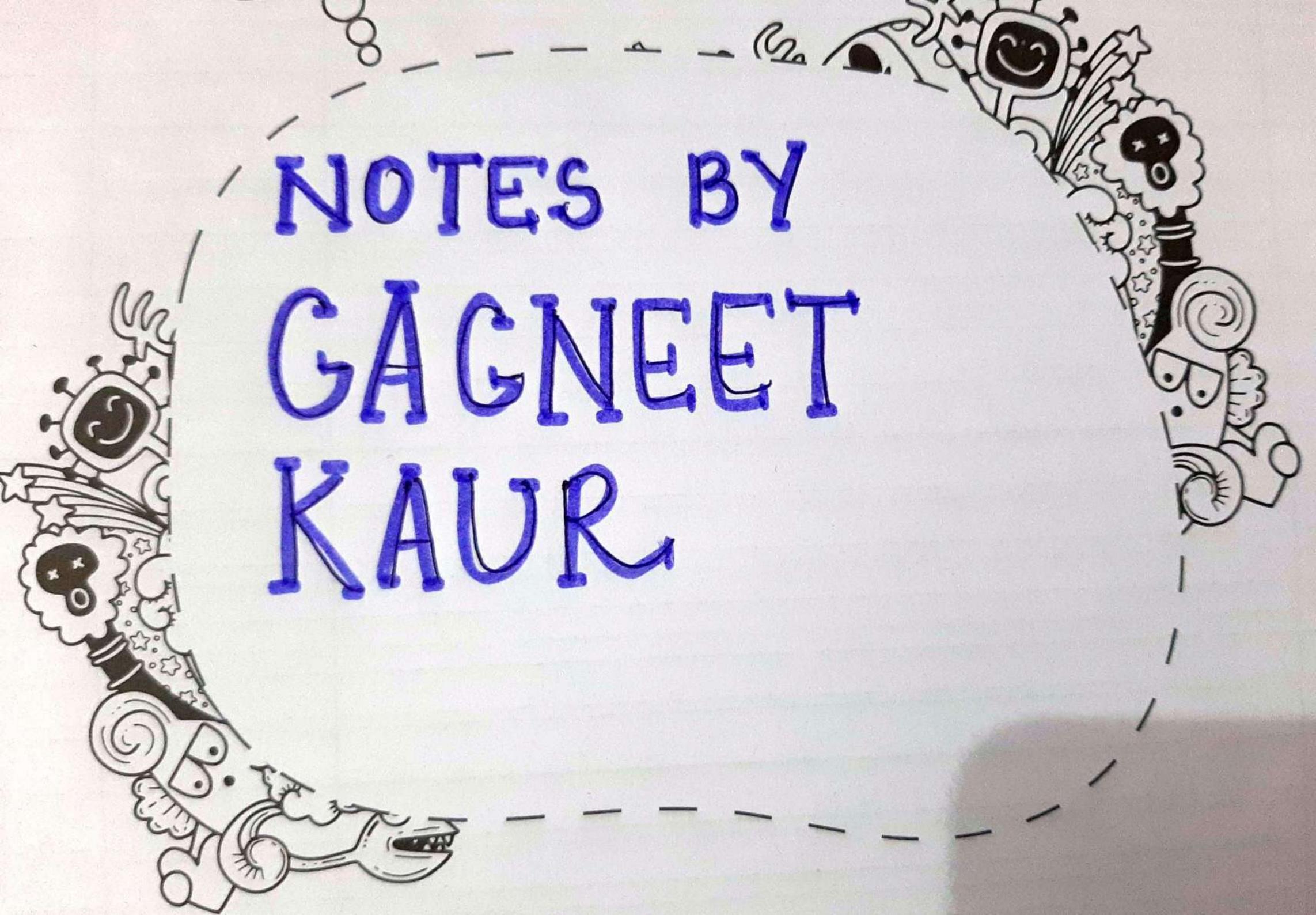


NOTES BY GAGNEET KAUR



Week 8 Optimization

Pillars of Machine Learning

- Linear Algebra : to find the structure/relationship between the data
- Probability : to model noise in the data
- Optimization : to get the best model out of all possible models.

Introduction to Optimization

- WHY OPTIMIZATION ?

→ we care about finding the "best" classifier !

"Least loss" \leftrightarrow "Maximum" reward

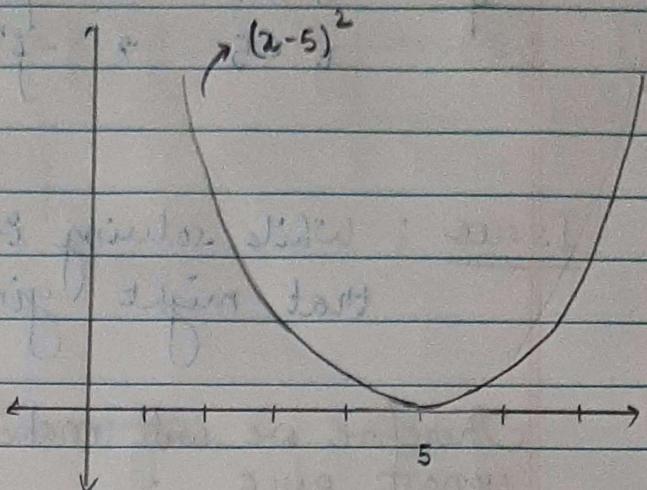
- GENERAL FORM OF OPTIMIZATION

$$\begin{array}{c}
 \text{OBJECTIVE} \\
 \min_{\mathbf{z} \in \mathbb{R}^d} f(\mathbf{z}) \quad \leftarrow \text{objective fn} \\
 \text{variable/ parameter} \rightarrow \\
 \text{CONSTRAINTS} \\
 \text{constraints} \rightarrow \begin{cases} g_i(\mathbf{z}) \leq 0 & + i=1, \dots, k \\ h_j(\mathbf{z}) = 0 & + j=1, \dots, l \end{cases} \\
 \text{Equality constraint} \qquad \qquad \qquad \text{Inequality constraints}
 \end{array}$$

#

Solving an unconstrained optimization problem

$$\min_{x \in \mathbb{R}} (x-5)^2$$

Approach

(i) Arbitrary Choice

start with $x_0 \in \mathbb{R}$

(ii) Update Rule

for $t=1, \dots, T$

update,

$$x_{t+1} = x_t + d$$

direction - "good direction to move"

(iii) end

Given x , what is a good direction?wantif $x > 5$, then more left

$d < 0$

if $x < 5$, then more right

$d > 0$

\Rightarrow direction d will depend on x , i.e., d must be a function of x .

Using derivatives, we have : if $x > 5 \Rightarrow f'(x) > 0$

If $x < 5 \Rightarrow f'(x) < 0$

Thus, we will CHOOSE $d = -f'(x)$ as it satisfies the conditions we want.

∴ We have,

$$\text{if } x > 5 \Rightarrow -f'(x) < 0$$

$$x < 5 \Rightarrow -f'(x) > 0$$

Issue : While solving it, we could encounter some x_0 that might give oscillating values.

Therefore we will make a small change in our UPDATE RULE :

$$x_{t+1} = x_t + \eta_t (-f'(x))$$

↳ step size (scalar quantity)

- How to choose STEP SIZE? (as a fn of t)

$$\boxed{\eta_t = \frac{1}{t+1}} \quad \leftarrow \text{good step size sequence}$$

$$\begin{matrix} 1, & \frac{1}{2}, & \frac{1}{3}, & \frac{1}{4}, & \dots & \text{etc.} & \frac{1}{t+1}, & \dots & \end{matrix}$$

$\uparrow \quad \uparrow \quad \uparrow \quad \dots$

$\eta_0 \quad \eta_1 \quad \dots$

Basic Algorithm for Unconstrained Optimization

GRADIENT DESCENT ALGORITHM

GOAL : $\min_{x \in \mathbb{R}} f(x)$

ALGORITHM : (1) Initialize at $x_0 \in \mathbb{R}$

(2) for $t = 1, 2, \dots$

$$x_{t+1} = x_t - n_t f'(x_t)$$

where $n_t = \frac{1}{t+1}$

(3) end

Properties of Gradient Descent

- If $n_t = \frac{1}{t+1}$, the algorithm converges.
- Gradient descent converges to "Local Minimum".

LOCAL

MINIMUM : $\exists \epsilon > 0$ such that

$$f(\hat{x}) \leq f(x) \quad \forall x \in [\hat{x} - \epsilon, \hat{x} + \epsilon]$$

GLOBAL

MINIMUM : $\forall x, f(x') \leq f(x)$

Gradient Descent and Taylor Series

Objective $\rightarrow \min_{x \in \mathbb{R}^d} f(x)$

Gradient Descent Update Rule : $x_{t+1} = x_t + n_t [f'(x)]$

what is so special about $-f'(x)$?

TAYLOR SERIES

$$f(x + \eta d) = f(x) + \eta d f'(x) + \frac{\eta^2 d^2}{2} f''(x) + \dots$$

 \downarrow x evaluations are all at x

Local Information
gives
Global Information

$$f(x + \eta d) = f(x) + \eta d f'(x) + f''(x) \frac{\eta^2 d^2}{2} + \dots$$

small
the step
size

"direction"

higher order terms

For small enough η ,

$$f(x + \eta d) \approx f(x) + \eta d f'(x)$$

$$\Rightarrow f(x + \eta d) - f(x) \approx \eta d f'(x)$$

function evaluation
at updated point
along 'direction' d

function evaluation
at the current
point

Want to choose a direction such that

$$f(x + \eta d) - f(x) < 0$$

Want ' d ' such that

$$\eta d f'(x) < 0$$

small η
constant $\Rightarrow \eta$ doesn't change
sign of this term

Want ' d ' such that

$$d f'(x) < 0$$

For the choice of $d = -f'(x)$,

$$df'(x) = -[f'(x)]^2 < 0$$

Gradient Descent For Multivariate Function

Higher Dimensions : $f(x_1, x_2) = x_1^2 + 4x_2 + 8x_2^2$

Derivative \Leftrightarrow Gradient \rightarrow (vector of partial derivatives)
(for one variable) (for multi-variables)

$$\nabla f \begin{pmatrix} [a] \\ [b] \end{pmatrix} = \left[\begin{array}{c|c} \frac{\partial f}{\partial x_1} & |_{x_1=a} \\ \hline \frac{\partial f}{\partial x_2} & |_{x_2=b} \end{array} \right]$$

Update Rule for Gradient Descent in Higher Dimensions :

$$\vec{x}_{t+1} = \vec{x}_t + \eta \left(-\nabla f(x_t) \right)$$

↑ ↑ ↑ ↑
 vector vector scalar vector

Taylor Series in Higher Dimensions

$$f(\underline{x} + \eta \underline{d}) = \underbrace{f(\underline{x})}_{\text{vector}} + \eta \underline{d}^T \nabla f(\underline{x}) + \dots$$

What is a good 'd' to choose?

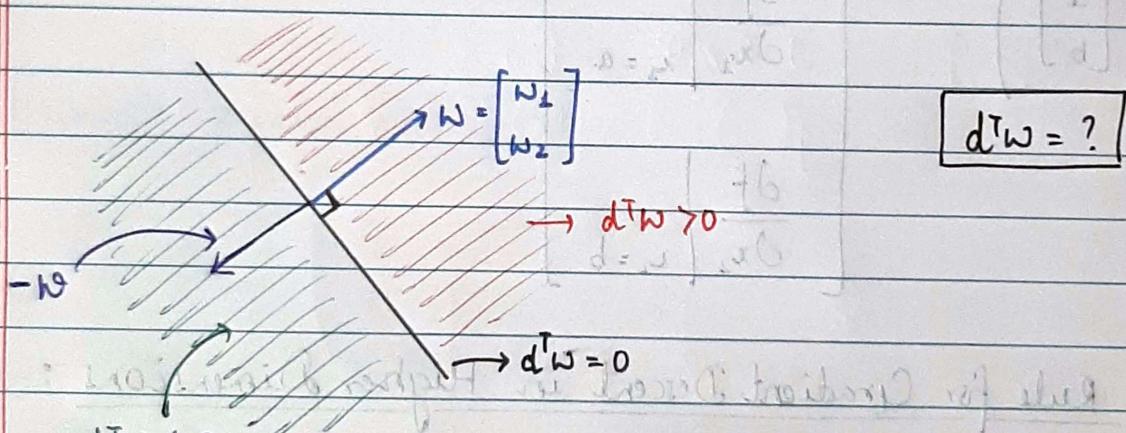
$$f(x + \eta d) - f(x) \approx \eta d^T \nabla f(x)$$

want d such that $d^T \nabla f(x) < 0$

NOTE: If $d = -\nabla f(x)$,

$$(-\nabla f(x))^T (\nabla f(x)) = -\|\nabla f(x)\|_2^2 < 0$$

- There can be other d 's in which if you move you will still reduce your function !!.



$-\nabla f(x)$ gives steepest descent

