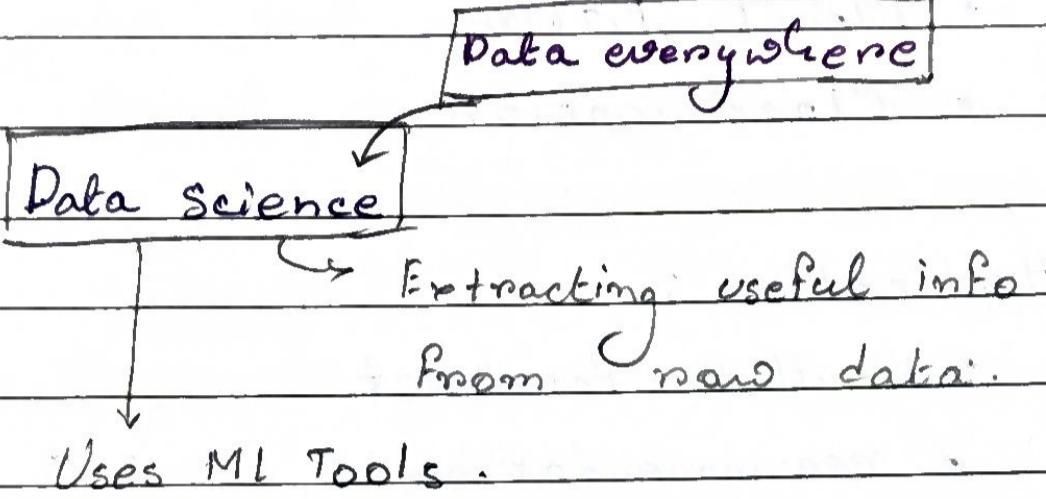


# Introduction to Numerical Optimization:

[week 1]



Data Engineering / Data Science / Data Analyst

↓                    ↓                    ↓

How to collect and store data      Making decision from that data.  
From data.      Make visualisation

## → Backbone of Datascience:

- Mathematics

- ① Linear Algebra

- ② calculus

- ③ Optimization.

- ④ Statistics.

- computer-Algo, computation, storage-

→ Broad classification of ML prob:

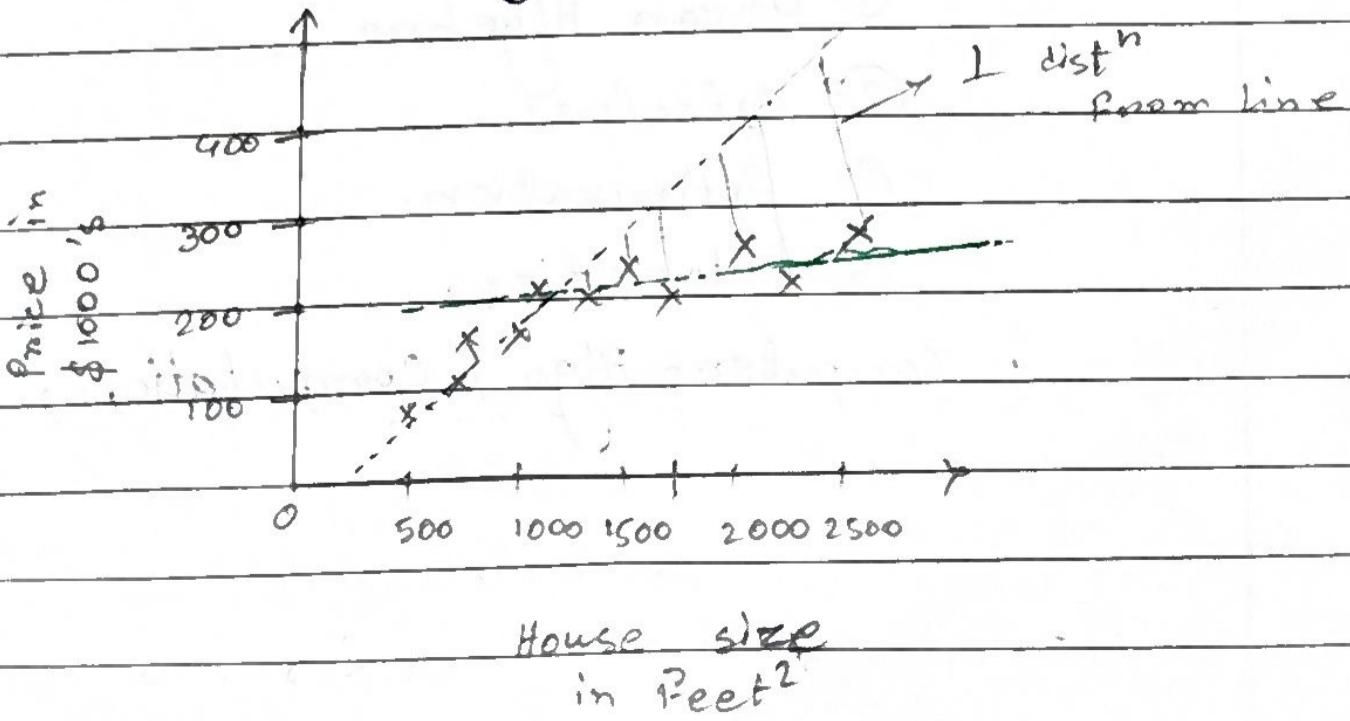
- Model Fitting
- Classification.

→ Application:

- Weather Forecast
- recommendations
- computer vision
- Pattern Recognition.

## # Why Optimization Matters?

- Regression Model
- Suppose a few data points  $\{(x_i, y_i)\}$  have been given in the plane.



we collect the data of some house of price and size and plot in graph.

Based on the price we predict the price or get an approx range of it.

- Suppose we guess there is a linear relation.

$$\text{price}(p) = \alpha \cdot \text{size}(s) + \beta$$

similar to  
straight  
line Eq.

Note: This is not the absolute model to predict price. There are more model.

Different Model give different result.

There are No unified model

There is No unified prediction

Q. Which  $\alpha$  and  $\beta$  to choose?

Bringing in optimization perspective

→ Ok, I have a model,

with that model, can I pass it

through all the data point

Typically very unlikely

of course, there will be error.

Mathematical

one

We can take 1 distri from the line and sum them up, a valid error indeed.

We can also use Misfit

$$\text{price} = \alpha(\text{size}) + \beta$$

direct correlation b/w

price  $\leftrightarrow$  size.

$$\sum (\alpha(\text{size}) + \beta - \text{price})^2$$



some of the term +ve and  
some -ve which cancel out  
so, we square.

- Typically we need/want least error in prediction.

- This is achieved by finding a straight line  $y = \alpha x + \beta$  such that.

$$f(a, b) = \sum_{i=1}^n (ax_i + b - y_i)^2 \text{ is minimized.}$$

Concept:

- Optimising problem is typically minimising or maximising  $f_n$  of some variables.

Objective (loss)  $f_n$

Decision Variables (Parameters)

Constraints (optional)

minimize  $f(x)$

subject to  $x \in S$

and

where  $f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  is a map,  $S \subseteq \mathbb{R}^n$  is a nonempty constraint set.

Eg: If we go from class to Mess and have three routes, not always the shortest route is better.

If the objective is workout then take the longest route.

## # Types of Problem:

- Unconstraint vs constraint
- Convex vs Non-convex

Unconstraint



We give a model a whole playground and tell them to figure out themselves, what the patterns are

Constraint



We give a model a whole playground with some directions/hint which make the model recognise patterns with our guidance

Why convexity Matters? ↴

We will look into it later

Convex prob → easy prob.

non-convex prob → extremely Hard prob.

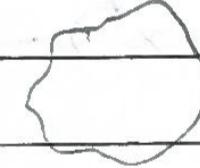
Example:

- Geometric problem: Given a Fixed perimeter, what is the largest area that can be covered?

suppose, we take rectangle  $\rightarrow L, B$   
maximize  $P(L, B) (= L \times B)$   
subject to  $2(L+B) = 10$

$$S = \{(L, B) \in \mathbb{R}^2 : 2(L+B) = 10\}$$

$$\rightarrow 2(a+b)$$



( $\Rightarrow$  We use

calculus.

For the area

## # Solution Concept:

Solution

$\hookrightarrow \text{F}(x)$

In optimization we want  
a point  $x^*$  such that

$F(x^*) \leq F(x)$  for all  $x$

when it is a minimization prob'

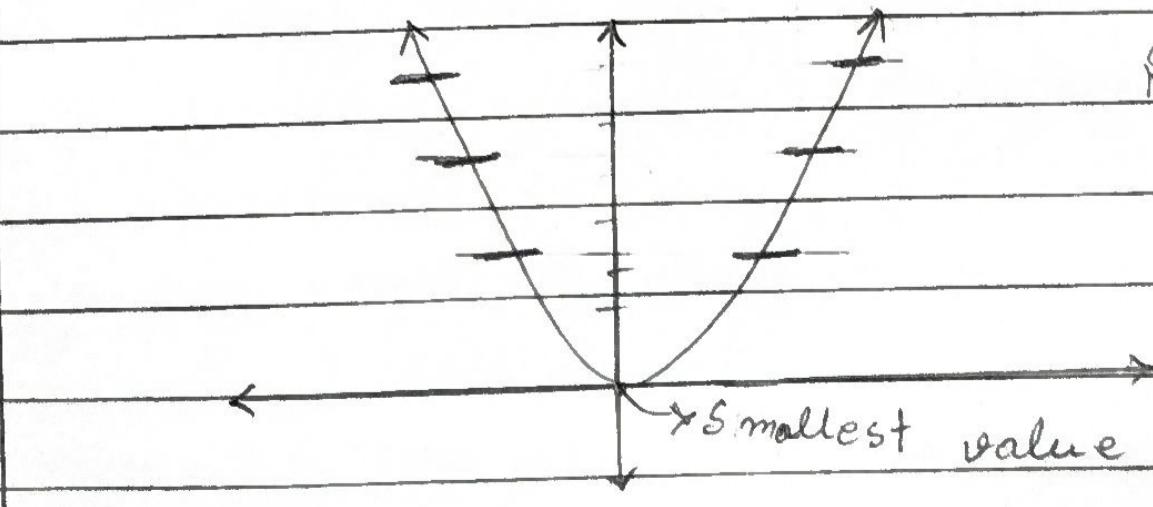
$x^*$   $\rightarrow$  Minimumpunkt / Minimum  
 $F(x^*)$   $\rightarrow$  Minimum Value

$x^* \rightarrow$  Minimizer

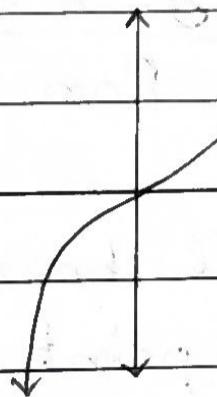
$F(x^*) \rightarrow$  Minimum Value.

E.g.:

minimize  $F(x) = x^2$  over  $S = \mathbb{R}$



if  $f(x) = x^3$



The minimum is

$-\infty$

↳ unknown

so, we say

it doesn't

have minimize

or even maximizer

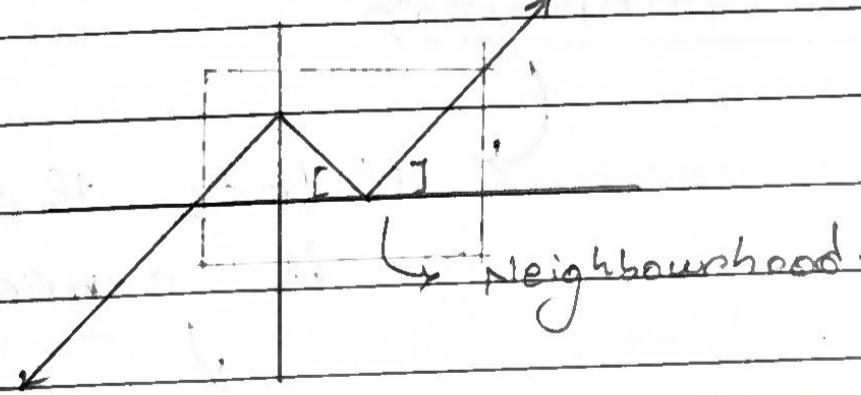
if we see other

side of the graph

$$f(x) = \begin{cases} x+2 & x \leq 0 \\ |2-x| & x > 0 \end{cases}$$

Here we use

Local Minimizer

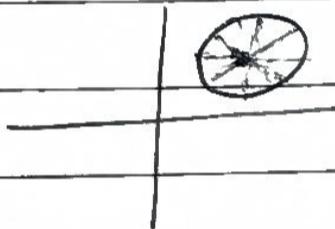


$f(x) > f(2)$  for around  
 $x \in [1, 3]$

Note: there is a break in the almost continuous function.

- A point  $x^* \in S$  is called a local minimizer if there exist  $\delta > 0$  such that  $f(x^*) \leq f(x)$  for all  $x \in B_\delta(x^*) \cap S$

$\delta \rightarrow$  that small break in the function and the point is typically



$\delta$  neighbourhood.

Why it is ~~an~~ Important?

→ Most of the numerical Algorithm we will study we have Local Minimizer.

↳ Unless the problem is convex problem

convexity will guarantee that local minimizer is global minimizer

Important thing in optimization:

- When can we say sol<sup>n</sup> exist?

→ When there is global Minimizer,  
local is considered (in the absence of)  
global minimizer and presence of  
convexity.

If a pt is sol<sup>n</sup>, how do we  
know it?

- How to Find sol<sup>n</sup>? → Main Topic  
↓ of this course

Numerically

Lecture Take away: Numerical optimization

↓  
Does Sol<sup>n</sup> exist,  
What kind of Sol<sup>n</sup> (local/global),  
How to Find it Numerically?

Data Science, its Back Bone, Application,  
Classification.

Most ML Training prob one solve via numerical optimizati-  
on of a loss fn.

# # A review of Calculus and Linear Algebra

## • Derivative

•  $f: \mathbb{R} \rightarrow \mathbb{R}$ , the derivative  $f'(x_0) = \frac{d}{dx} f(x_0)$

$$= \lim_{t \rightarrow 0} \frac{f(x_0+t) - f(x_0)}{t}$$

## • calculus of derivative.

$$\bullet (f \pm g)'(x) = f'(x) \pm g'(x)$$

$$\bullet (fg)'(x) = f'(x)g(x) + g'(x)f(x)$$

$$\bullet (\frac{1}{g})'(x) = \frac{-g'(x)}{(g(x))^2}$$

$$\bullet f''(x) = (f'(x))'$$

$$f(x) = \sin(x)$$

$$f'(x) = \cos(x)$$

$$f''(x) = -\sin(x)$$

$$f(x) = x^2$$

$$f'(x) = 2x$$

$$f''(x) = 2$$

## # Taylor's Expansion Theorem:

- If a  $f_n$  is differentiable at around point  $x_0$

$f(x) \approx$

Ist order derivative.

If  $x$  is  $x_0$  is a very good approximator of  $f(x)$

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

If  $x$  is very close to  $x_0$  this is a very good approximator of  $f(x)$

this is

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0) \frac{(x - x_0)^2}{(x - x_0)^2}$$

- If double derivative exist
- $x$  is very close to  $x_0$ .

double differentiable at around

$$+ f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2$$

Eg:  $P(x) = e^x$ ,  $x_0 = 0$

1st order:  $e^x \approx 1+x$

2nd order:  $e^x \approx 1+x+\frac{x^2}{2}$

## # Differentiability of Function of several variable:

For a function  $P: \mathbb{R}^n \rightarrow \mathbb{R}$ , the partial derivative is used,

$P$  at  $x$  with respect to  $x_i$

fix one  
of the dim

$\mathbb{R}^2 \rightarrow \mathbb{R}$

$$\frac{\partial P}{\partial x_i}(x) = \lim_{t \rightarrow 0} \frac{P(x+te_i) - P(x)}{t}$$

Eg:  $P(x, y) = x^2 + xy + y^2$

$x+te_i$  is actually a vector sum.

$$\frac{\partial P}{\partial x} = 2x + y$$

$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} + t \begin{bmatrix} 4 \\ 5 \end{bmatrix} = \begin{bmatrix} 2+4t \\ 3+5t \end{bmatrix}$$

$$\frac{\partial P}{\partial y} = x + 2y$$

• Derivative of  $F$  =  $\nabla F(x)$   
↳ Gradient

$$= \begin{bmatrix} \frac{\partial F}{\partial x_1} \\ \vdots \\ \frac{\partial F}{\partial x_n} \end{bmatrix} \rightarrow \begin{array}{l} \text{collection of all} \\ \text{the partial derivatives} \\ \text{vector.} \end{array}$$

For the above Eq :

$$\nabla F(x) = \begin{bmatrix} 2x+y \\ x+2y \end{bmatrix}$$

Let,

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

- Second-order partial derivative:

$$\frac{\partial^2 F}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_i} \left( \frac{\partial F}{\partial x_j} \right)$$

- Hessian matrix: all second-order partial derivatives.

$$H_F(x) = \nabla^2 F(x) = \begin{bmatrix} \frac{\partial^2 F}{\partial x_1^2} & \cdots & \frac{\partial^2 F}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 F}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 F}{\partial x_n^2} \end{bmatrix}$$

Eg:  $F(x, y) = x^2 + yx + y^2$

$$\frac{\partial^2 F}{\partial x \partial y} / \frac{\partial^2 F}{\partial x^2} / \frac{\partial^2 F}{\partial y^2}$$

$$H(x, y) = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

- A matrix is a rectangular arrangement of numbers, symbols or expressions written in rows and columns, usually enclosed with brackets.

## # ~~Matrices~~ Taylor's theorem of several variables:

- For a function  $F: \mathbb{R}^n \rightarrow \mathbb{R}$ , the derivative of  $F$  at  $x$ , denoted by  $\nabla F(x)$  is a vector such that

$$F(x+h) = F(x) + \langle \nabla F(x), h \rangle + O(\|h\|)$$

$$f(x) = f(x_0) + f'(x_0)(x-x_0)$$

scalar  
analog

$$y = (y - x_0) + x_0$$

$\downarrow h$

Very very  
small term

Roughly,

$$F(x+h) \approx F(x) + \langle \nabla F(x), h \rangle$$

- For a  $f_n : \mathbb{R}^n \rightarrow \mathbb{R}$ , the second derivative of  $f$  at  $x$ , denoted by  $\nabla^2 f(x)$  is a Matrix such that,

$$f(x+h) = f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2!} \left[ h, \nabla^2 f(x) h \right] + \dots$$

$\downarrow$   
dot( $\nabla f(x)$ ,  $h$ )

$O(\|h\|^2)$

Roughly,

$$f(x+h) \approx f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2!} \left[ h, \nabla^2 f(x) h \right]$$

Note : ① This is a quadratic Function.

② We need to be careful about linear  $f_n$  and quadratic  $f_n$ .  
As, they will be the function we will be frequently need data science course.

## Dot Product:

- For  $u, v \in \mathbb{R}^n : u \cdot v = \sum_{i=1}^n u_i v_i$

- $u = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, v = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$

$$u \cdot v = (1 \cdot 3) + (2 \cdot 4) = 11$$

Transpose: Flips rows and columns.

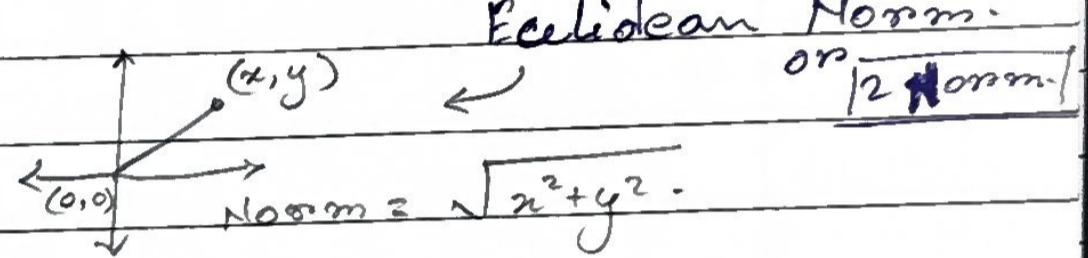
$$(A^T)_{ij} = A_{ji}$$

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$$

## \* Norm (length of a vector)

- Denoted by  $\|\vartheta\|$
- For  $\vartheta = (\vartheta_1, \vartheta_2, \dots, \vartheta_n)$

$$\|\vartheta\|_2 = \sqrt{\vartheta_1^2 + \vartheta_2^2 + \dots + \vartheta_n^2}$$



$$\|\vartheta\| = |\vartheta_1| + |\vartheta_2| + \dots + |\vartheta_n|$$

- Measures the magnitude (length) of a vector.

Dot Product :

$$\|\vartheta\|_2 = \vartheta^T \vartheta.$$

$1 \times n \leftrightarrow n \times 1$