

Special Lectures

I

- Introduction to Continuous Optimization
- Application in Machine Learning

Fundamental problem of continuous optimization :

Find \underline{x}_* such that \underline{x}_* is the minimum of the cost function $f(\underline{x})$

OP

- o Unconstrained problem
- o $\underline{x} \in \mathbb{R}^n$
- o $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous and twice differentiable.

~~Constrained Problem~~

Note :

- o Parameter space for continuous optimization is \mathbb{R}^n
- o In case of discrete optimization, the parameter space is N^n and the solution method is linear programming.

Why optimization ?

- o "Operations Research" — optimizing industrial production, logistics, supply chains...
- o Key technique in Machine Learning

Why study Mathematics behind Optimization?

- o Inbuilt optimization methods in R, Python, and Matlab.
- o Important to understand error messages \rightarrow

from these packages.

II

- Otherwise, risk of ALGO.

Fundamentals of Unconstrained Optimization:

Aim is to solve

$$\underline{x}_* = \arg \min_{\underline{x} \in \mathbb{R}^n} f(\underline{x})$$

\leftarrow Unconstrained

Terminology:

- \underline{x}_* is a global minimizer if $f(\underline{x}_*) \leq f(\underline{y}) \forall \underline{y} \in \mathbb{R}^n$.
- \underline{x}_* is a local minimizer if there exists a neighbourhood N of \underline{x}_* such that $f(\underline{x}_*) \leq f(\underline{y}) \forall \underline{y} \in N$.
- \underline{x}_* is a strict local minimizer if there exists a neighbourhood N of \underline{x}_* such that $f(\underline{x}_*) < f(\underline{y}) \forall \underline{y} \neq \underline{x}_* \in N$.

$$f(\underline{x}_*) < f(\underline{y}) \quad \forall \underline{y} \neq \underline{x}_* \in N.$$

Necessary conditions for optimality:

If \underline{x}_* solves the OP (i.e. is a local minimizer)
then $\nabla f(\underline{x}_*) = 0$. (Analogous to First Derivative Test)

Furthermore, the Hessian matrix

$$H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j} \Big|_{\underline{x}_*}$$



is positive-semi-definite (analogous to Second Derivative Test).

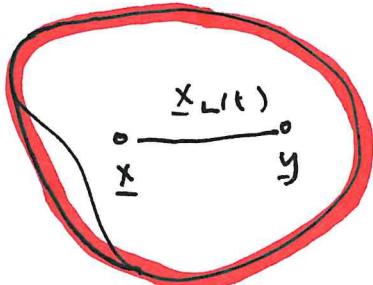
Necessary and Sufficient Conditions:

- $\nabla f(\underline{x}_*) = 0$

- H_{ij} is positive-definite:

$$\langle \underline{y}, H \underline{y} \rangle > 0 \quad \forall \underline{y} \neq 0 \in \mathbb{R}^n$$

Convexity: ~~A~~ convex set if, for each pair \underline{x} and \underline{y} in S , the line segment $\underline{x}_L(t) = \underline{x}(1-t) + t\underline{y}$ $t \in [0,1]$ is contained entirely in S .



CONVEX



NOT CONVEX

A convex function: A function

$$f: (S \subset \mathbb{R}^n) \longrightarrow \mathbb{R}$$

is convex if:

- S is a convex set

- The following inequality holds for all \underline{x} and \underline{y} in S :

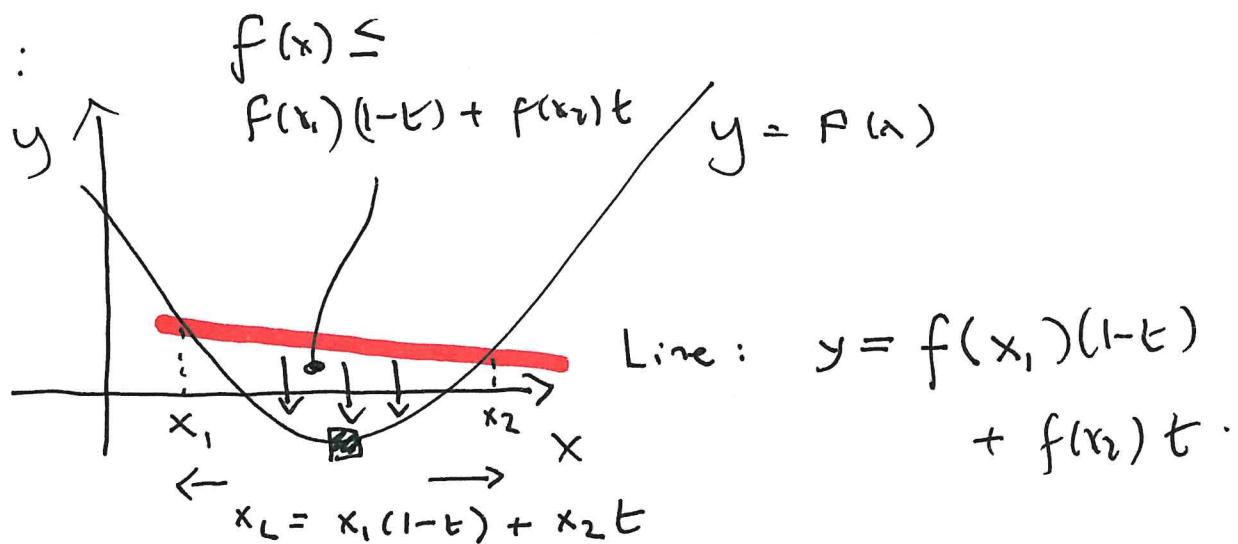
$$f(\underline{x}(1-t) + t\underline{y}) \leq (1-t)f(\underline{x}) + t f(\underline{y}) \quad \forall t \in [0,1].$$

IV Fundamental Theorem of Convex Optimization:

When f is a convex function, any local minimizer of \underline{x}_* is a global minimizer.

If, in addition f is differentiable, then any stationary point $\nabla f(\underline{x}_*) = 0$ is a global minimizer of f .

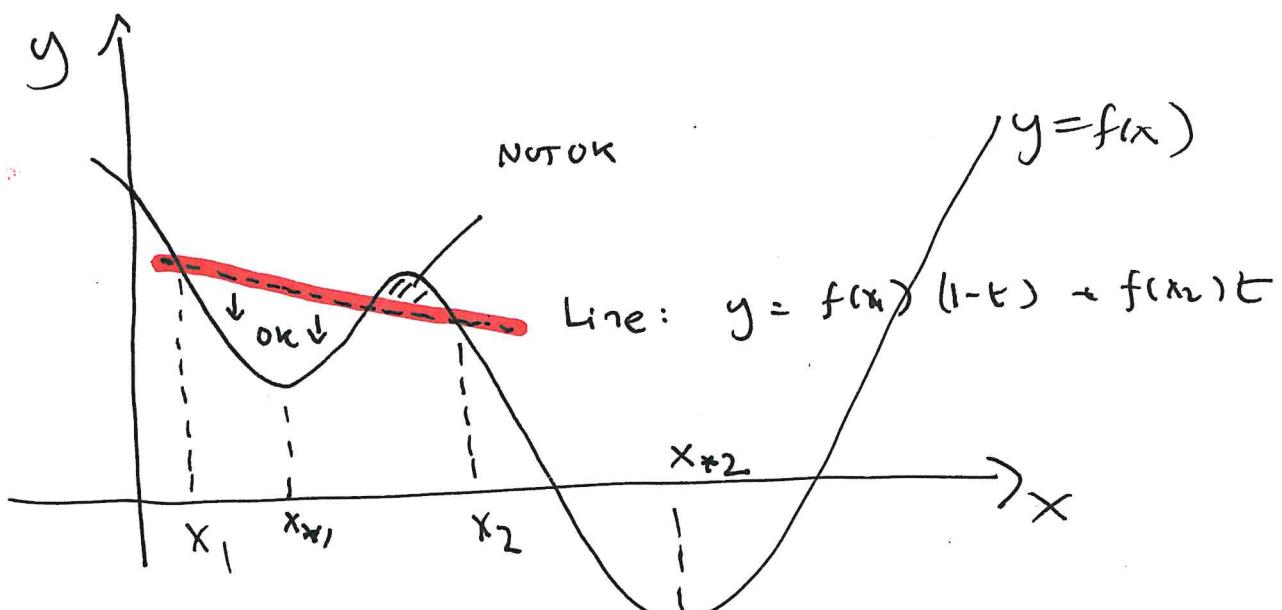
IDEA :



$$\text{Line: } y = f(x_1)(1-t) + f(x_2)t$$

$$+ f(x_2)t$$

UNIQUE GLOBAL MINIMIZER.



NOT A CONVEX FUNCTION.

- o Local min @ \underline{x}_{*1}
- o Global min @ \underline{x}_{*2}

The Model Problem :

IV

$$f(\underline{x}) = c + \langle \underline{a}, \underline{x} \rangle + \frac{1}{2} \langle \underline{x}, \underline{B} \underline{x} \rangle$$

- c is constant
- \underline{a} is a constant vector
- \underline{B} is a positive-definite matrix

$$\nabla f = \underline{a} + \underline{B} \underline{x}$$

$$\nabla f = 0 \Rightarrow \underline{x} = -\underline{B}^{-1} \underline{a}$$

$$\boxed{\underline{x}_* = -\underline{B}^{-1} \underline{a}}$$

Part 2 - Numerical Optimization

Two main methods for unconstrained optimization :

- Linesearch
- Trust Region

We look briefly at Linesearch. This is an iterative method, which aims to get closer and closer to the solution with successive guesses.

If we label the guesses as

$$\underline{x}_0, \underline{x}_1, \dots, \underline{x}_k, \underline{x}_{k+1}, \dots$$

we have :

$$\underline{x}_{k+1} = \underline{x}_k + \underline{s}_k$$

Here, the vector \underline{s}_k can be decomposed as :

$$\underline{s}_k = \alpha_k \begin{matrix} \uparrow \\ \text{STEP SIZE} \end{matrix} \begin{matrix} \nwarrow \\ \text{SEARCH DIRECTION} \end{matrix} \underline{r}_k$$

Once an appropriate search direction has been VI selected, α_k can be chosen as:

$$\alpha_k = \arg \min_{\alpha > 0} f(x_k + \alpha p_k)$$

As this is a 1D optimization problem, it is easier to solve than the original OP. Thus, we have reduced the original n -dimensional OP to a sequence of successive steps :

- o Finding the search direction
- o Solving a 1D OP.

Finding the search direction: Three main methods :

- o Steepest Descent (SD)
- o Newton / Quasi-Newton
- o Trust Region .

SD : Choose p_k to decrease cost function as much as possible. We have :

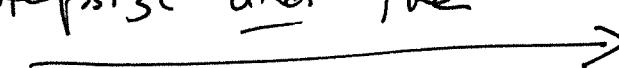
$$p_k \cdot \nabla f(x_k) = \text{Directional derivative in direction } p_k$$

Maximum negative change when $p_k \propto -\nabla f(x_k)$.

Hence :

$$p_k = \frac{-\nabla f(x_k)}{\|\nabla f(x_k)\|_2}$$

Newton : Idea is to approximate $f(x_k + s_k)$ as a quadratic function. key departure from SD :

p_k contains info about the stepsize and the search direction. 

We have :

$$\begin{aligned} \underline{x}_{n+1} &= \underline{x}_n + \underline{p}_n \\ f(\underline{x}_{n+1}) &= f(\underline{x}_n + \underline{p}_n) \\ &\simeq f(\underline{x}_n) + \underline{p}_n \cdot \nabla f(\underline{x}_n) + \frac{1}{2} \underline{p}_n \cdot \underline{p}_n \frac{\partial^2 f}{\partial \underline{x}_i \partial \underline{x}_j} \Big|_{\underline{x}_n} \\ &= m_n(\underline{p}) \end{aligned}$$

Minimize $m_n(\underline{p})$ — Quadratic model problem.

$m_n(\underline{p})$ is minimized when

$$\underline{p} = -B^{-1}\underline{a}$$

$$\text{where } \underline{a} = \nabla f(\underline{x}_n), \quad B_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j} \Big|_{\underline{x}_n}$$

Notation:

$$\underline{x}_{n+1} = \underline{x}_n + \underline{p}_n^N \leftarrow \text{for Newton}$$

$$\underline{p}_n^N = -B^{-1}\underline{a}$$

\underline{p}_n^N contains info about step size and search direction.

Convergence : For \underline{x}_0 sufficiently close to \underline{x}_* , linesearch (SD) has the property that

$$\|\underline{x}_{n+1} - \underline{x}_*\|_2 \leq C \|\underline{x}_n - \underline{x}_*\|_2$$

Constant C depends on properties of $\frac{\partial^2 f}{\partial x_i \partial x_j} \Big|_{\underline{x}_n}$, if C is close to 1 then convergence can be slow.

In contrast, Newton has the property that:

VIII

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 \leq \tilde{C} \|\mathbf{x}_k - \mathbf{x}_* \|_2^2$$

- Quadratic convergence

- Newton is much faster than Linesearch but a matrix inversion needs to be performed at each iteration ($O(n^3)$).

Secant Method: Approximation of \mathbf{P}^{-1} at each iteration. One particular approx. is the BFGS algorithm — in widespread use ($O(n^2)$).

Part 3 — Application to Machine Learning