Week 2, Lectures 2-3

Online lecture on Tuesday A.M. for the foreseeable future

Plan for today:

- Finish up looking at the quadratic model problem ( Ch. 2 )

- Line search Methods (Ch.3 )

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

When $B \in \mathbb{R}^{n \times n}$ is no longer positive-definite ?

We still assume that B is symmetric.

Theorem: $f$ attains a min. if and only if B is positive-semi-definite and $\underline{a}$ is in the range of B. If B is positive-semi-definite (PSD), then every $\underline{p}$ satisfying $B\underline{p} = -\underline{a}$ is a global minimizer of $f$.

Proof: Assume that B is P.S.D. and that $\underline{a}$ is in the range of B. Given this assumption, there exists an $\underline{x} \in \mathbb{R}^n$ such that:
$$B\underline{x} = -\underline{a} .$$

For any $\underline{w} \in \mathbb{R}^n$, consider:

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

$B$ Symmetric

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

$$= c + \langle \underline{a}, \underline{x} \rangle + \frac{1}{2} \langle \underline{x}, B\underline{x} \rangle + [\langle \underline{a}, \underline{w} \rangle + \langle \underline{x}, B\underline{w} \rangle] + \frac{1}{2} \langle \underline{w}, B\underline{w} \rangle$$

$$= f(\underline{x}) + \overset{=0}{\langle -B\underline{x}, \underline{w} \rangle + \langle \underline{x}, B\underline{w} \rangle} + \frac{1}{2} \underbrace{\langle \underline{w}, B\underline{w} \rangle}_{\geq 0}$$

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

Since $B$ is P.S.D.,

$$f(\underline{x}+\underline{w}) \geq f(\underline{x}) \quad \forall \underline{w} \in \mathbb{R}^n.$$

Hence, $\underline{x}$ is a global minimizer.

For the other way around, suppose that
$f$ ~~$B=a$~~ has a minimizer ($\underline{x}$, say).

By the first-order optimality condition,
$$\nabla f(\underline{x}) = 0 \implies B\underline{x} = -\underline{a}.$$

Hence, $\underline{a}$ is in the range of $B$.

By the second-order optimality,
$$\frac{\partial^2 f}{\partial x_i \partial x_j}\bigg|_{\underline{x}} \quad \text{is} \quad P.S.D.$$

But for the Quadratic Model Problem, this is just $B$, hence $B$ is P.S.D.

Theorem: $f$ has a unique minimizer if and only if $B$ is strictly P.D.

Proof: Assume that $B$ is P.D. Then, $B$ is invertible, so $\underline{a}$ is in the range of $B$, so let $\underline{x}$ solve $B\underline{x} = -\underline{a}$.

Consider:
$$f(\underline{x} + \underline{w}) = f(\underline{x}) + \overset{0}{\overbrace{}} + \tfrac{1}{2}\langle \underline{w}, B\underline{w} \rangle$$

But $B$ is P.D. so $\langle \underline{w}, B\underline{w} \rangle > 0 \; \forall \underline{w} \neq 0$.

Hence $f(\underline{x} + \underline{w}) > f(\underline{x}) \; \forall \underline{w} (\neq 0) \in \mathbb{R}^n$.

So $\underline{x}$ is the unique global minimizer.

For the other way around, suppose $f$ has a unique global minimizer (call it $\underline{x}$). We use a proof by contradiction: assume that $B$ is not positive-definite.

THERE IS A GAP HERE — FILLED IN LATER.

Then, we can find a non-zero vector $\underline{w}$ such that $B\underline{w} = 0$. Then,

$$f(\underline{x}+\underline{w}) = f(\underline{x}) + \cancel{\langle \underline{w}, \underline{B}\underline{w}\rangle} \cancel{\tfrac{1}{2}\langle \underline{w}, B\underline{w}\rangle}$$

$$= f(\underline{x}).$$

Hence, $\underline{x}+\underline{w}$ is also a minimizer. This contradicts uniqueness. Hence, the only way to have a unique global minimizer is if $B$ is strictly P.D.

Take—home message ⟶ in case of the model problem

Continuous optimization is a nice application of Calculus and Linear algebra. In the next chapters we will attempt to approximate a general O.P. with a quadratic problem, which can be solved using Linear Algebra.

# Chapter 3 — Line Search Methods

Notation for the O.P. :

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

Line Search Methods are iterative.
We start off with an initial guess for $\underline{x}_*$
(call it $\underline{x}_0$). We make a sequence
of improved guesses $\underline{x}_k$, such that :

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

Typically, $\underline{s}_k$ depends on $\nabla f(\underline{x}_k)$ and
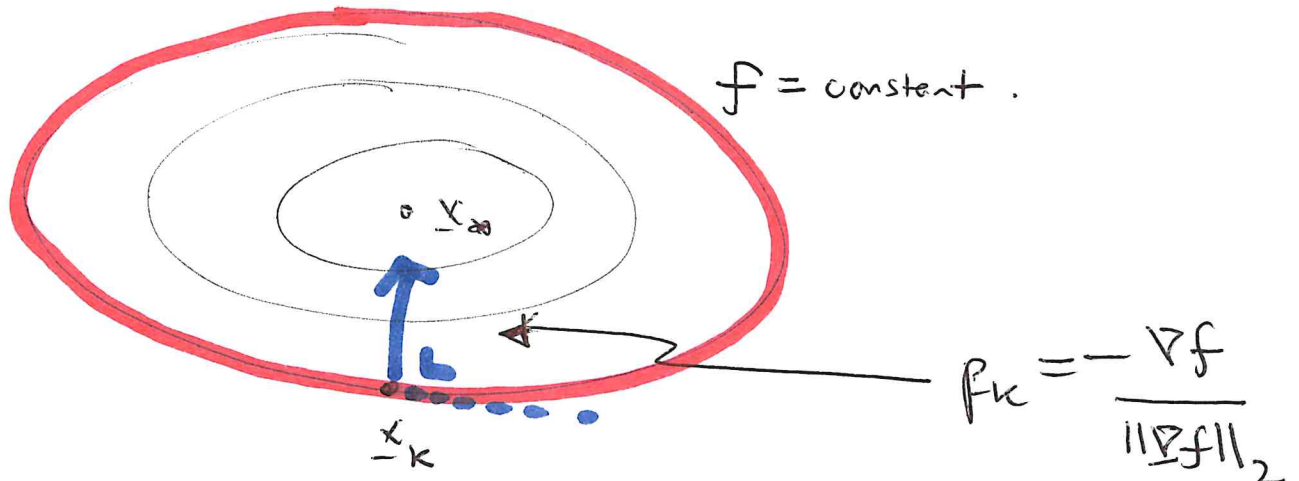is broken up into a magnitude and a
direction :

$$\underline{s}_k = \alpha_k \, \underline{p}_k \quad,$$

where typically $\underline{p}_k$ is a unit vector.
When $\underline{p}_k$ is a unit vector, $\alpha_k$ can
be found by solving a 1D OP :

$$\alpha_k = \arg\min_{\alpha > 0} f(\underline{x}_k + \alpha \, \underline{p}_k) \quad (*)$$

# § 3.2  Steepest Descent Method



$f = $ constant.

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

Look at:

$$\text{DOMINANT}$$

$$f(\underline{x}_k + \tilde{\alpha} p) = f(\underline{x}_k) + \overbrace{\alpha p_i \frac{\partial f}{\partial x_i}(\underline{x}_k)}$$

$$+ \frac{1}{2}\alpha^2 p_i p_j \frac{\partial^2 f}{\partial x_i \partial x_j}(\underline{x}_k + t p),$$

$$t \in (0, \alpha).$$

To reduce $f$ as much as possible in one iteration ($\underline{x}_k \to \underline{x}_k + \alpha p$), we need to make the dominant term as negative as possible. To do this, we simply take:

$$p = -\frac{\nabla f(\underline{x}_k)}{\|\nabla f(\underline{x}_k)\|_2}.$$

Then,

$$f(x_k + \alpha p) = f(x_k) - \alpha \left. \frac{\nabla f \cdot \nabla f}{\|\nabla f\|} \right|_{x_k} + O(\alpha^2)$$

$$= f(x_k) - \alpha \|\nabla f\|_2 \Big|_{x_k} + O(\alpha^2).$$

Thus, the ~~cha~~ reduction in $f$ is maximized.

$$p_k = - \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|_2} \qquad (1)$$

is the direction of <u>Steepest Descent</u>.

$$\boxed{\text{Pseudocode in notes (p. 23)}}$$

//

Another choice for the search direction is the Newton Method (§3.3)

Taylor approximation:

$$f(x_k + p) \simeq \underbrace{f(x_k)}_{C} + p_i \underbrace{\frac{\partial f}{\partial x_i}(x_k)}_{a_i} + \tfrac{1}{2} p_i p_j \underbrace{\frac{\partial^2 f}{\partial x_i \partial x_j}(x_k)}_{B_{ij}}$$

This is the Quadratic Model Problem:

$$m_k(p) = c + \langle a, p \rangle + \tfrac{1}{2} \langle p, Bp \rangle$$

Minimize $m_k(p)$. If $B$ is invertible, then

$$p = -B^{-1} a .$$

This gives the descent direction in the Newton method. Restoring $k$, we have:

$$p_k^N = -B^{-1}(x_k) \nabla f(x_k) \qquad (2)$$

Pseudocode in notes on p. 24

Equation (2) is the Newton descent direction

For the Newton method to yield a reduction in $f$ from $x_k$ to $x_k + p_k^N$, we require $B(x_k)$ to be positive-definite. Proof:

$$f(x_k + t p_k^N) = f(x_k) + t \sum_{i=1}^{n} (p_k^N)_i \frac{\partial f}{\partial x_i}(x_k) + O(t^2)$$

Hence,

$$f(\underline{x}_n + t\underline{p}_n^N) = f(\underline{x}_n)$$

$$= t \sum_{i=1}^{n} \sum_{j=1}^{n} \left[ (B^{-1})_{ij} \frac{\partial f}{\partial x_j} \right]_{\underline{x}_k} \left( \frac{\partial f}{\partial x_i} \right) + O(t^2)$$

Reason: $\quad \underline{p}_n^N = - B^{-1} \underline{\nabla} f$.

$$(\underline{p}_n^N)_i = - \sum_{j=1}^{n} (B^{-1})_{ij} \frac{\partial f}{\partial x_j}$$

( Matrix multiplication ).

Hence,

$$f(\underline{x}_n + t\underline{p}_n^N) = f(\underline{x}_n)$$

$$\underbrace{- t \langle \nabla f(\underline{x}_n), B^{-1} \underline{\nabla} f \rangle}_{NEG.} + O(t^2)$$

$$> 0.$$

Hence,

$$f(\underline{x}_n + t\underline{p}_n^N) < f(\underline{x}_n),$$

for $t$ suff. small.

Lecture 3

Advantages of using $P_n^N$:

- Step-length is provided, no need to solve the sub-problem (*)

- Simple criterion for method to work ($P(x_n)$ is pos.-definite)

- Quadratic convergence

We illustrate the idea of quadratic convergence here for a 1D problem where we wish to solve:

$$X_* = \arg \min_{x \in \mathbb{R}} f(x).$$

By first-order optimality, $f'(x_*) = 0$.

$$X_{k+1} = X_k - \frac{f'(x_k)}{f''(x_n)} \quad (3)$$

Error:

$$\epsilon_k = X_* - X_k \qquad \qquad X_u = x_* - \epsilon_k$$

$$\epsilon_{k+1} = X_* - X_{u+1} \qquad X_{u+1} = X_* - \epsilon_{u+1}$$

Sub in to (3):

$$X_* - \epsilon_{u+1} = X_* - \epsilon_k - \frac{f'(\overbrace{X_* - \epsilon_n}^{X_k})}{f''(X_* - \epsilon_n)}$$

$$\Rightarrow \epsilon_{k+1} = \epsilon_k + \frac{f'(x_* - \epsilon_k)}{f''(x_* - \epsilon_k)}$$

$$\Rightarrow \epsilon_{k+1} = \epsilon_k + \frac{\left[ f'(x_*) - f''(x_*)\epsilon_k + \frac{1}{2}f'''(x_*)\epsilon_k^2 + \cdots \right]}{f''(x_*) - f'''(x_*)\epsilon_k + \cdots}$$

$$\Rightarrow \epsilon_{k+1} = \epsilon_k - \frac{f''(x_*)\epsilon_k \left[ 1 - \frac{1}{2}\frac{f'''(x_*)}{f''(x_*)}\epsilon_k + \cdots \right]}{f''(x_*)\left[ 1 - \frac{f'''(x_*)}{f''(x_*)}\epsilon_k + \cdots \right]}$$

BINOMIAL THEOREM : $\quad (1+z)^p = 1 + pz + \frac{p(p-1)}{2}z^2 + \cdots$

$$\epsilon_{k+1} = \epsilon_k - \epsilon_k \left[ 1 - \frac{1}{2}\frac{f'''(x_*)}{f''(x_*)}\epsilon_k + \cdots \right]\left[ 1 + \frac{f'''(x_*)}{f''(x_*)}\epsilon_k - \cdots \right]$$

$$\Rightarrow$$

$$\epsilon_{k+1} = \epsilon_k - \epsilon_k \left[ 1 + \frac{f'''(x_*)}{f''(x_*)}\epsilon_k - \frac{1}{2}\frac{f'''(x_*)}{f''(x_*)}\epsilon_k + O(\epsilon_k^2) \right]$$

$$\Rightarrow \epsilon_{k+1} = \epsilon_k - \epsilon_k \left[ 1 + \frac{1}{2}\frac{f'''(x_*)}{f''(x_*)}\epsilon_k + O(\epsilon_k^2) \right]$$

Hence,

$$\epsilon_{k+1} = \cancel{x_k} - \cancel{x_k} - \frac{1}{2} \epsilon_k^2 \frac{f'''(x_*)}{f''(x_*)} + O(\epsilon_k^3)$$

Hence,

$$\epsilon_{k+1} = -\frac{1}{2} \epsilon_k^2 \frac{f'''(x_\circ)}{f''(x_\circ)} + O(\epsilon_k^3) \qquad \overset{x_k + \epsilon_k}{\curvearrowright}$$

OR

$$\epsilon_{k+1} = -\frac{1}{2} \epsilon_k^2 \frac{f'''(x_k)}{f''(x_*)} + O(\epsilon_k^3)$$

Drawbacks:

- Computation of Hessian at each iteration
- Requires inversion of the Hessian at each iteration ($O(n^3)$)

Amelioration — approximate the Hessian matrix using the SECANT METHOD (§3.4)

To see how the Secant Method works, we go back to the 1D problem, and we look at:

$$f'(x_k + \delta x) \simeq f'(x_n) + f''(x_n) \, \delta x.$$

Take: $\delta x = x_{n+1} - x_k$. Hence, this equation becomes:

$$\underbrace{f'(x_{n+1}) - f'(x_n)}_{y_k} \simeq f''(x_n) \underbrace{[x_{n+1} - x_n]}_{s_k}.$$

Approximate Hessian:

$$f''(x_n) \simeq y_k / s_k.$$

Equivalent $n$-dimensional analogy:

$$\underbrace{\nabla f(x_{n+1}) - \nabla f(x_n)}_{y_k} \simeq B(x_n) \underbrace{[x_{n+1} - x_n]}_{s_k}.$$

Or $\underline{y}_k = \underbrace{B(x_n) \, \underline{s}_k}_{B_{k+1}}.$

# Pseudocode:

Choose $\underline{x}_0$ sufficiently close to $\underline{x}_*$

Choose $B_0$.

for $k = 0, 1, 2, \ldots$

    Compute the descent direction $B_k \underline{p}_k = -\nabla f(\underline{x}_k)$   by solving

    Choose the stepsize $\alpha_k$.

    Write $\underline{s}_k = \alpha_k \underline{p}_k$

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

    Update $\underline{y}_k = \nabla f(\underline{x}_{k+1}) - \nabla f(\underline{x}_k)$

    Update$^{approx}$ Hessian for next iteration by

    solving $\underline{y}_k = B_{k+1} \underline{s}_k$

end for

## Chapter 4 — BFGS method.

Problem — we need to solve for $B_{k+1}$ in the equation

$$\underline{y}_k = B_{k+1} \underline{s}_k \quad (4)$$

1D: $\quad f''(x_k) \simeq y_k / s_k$.

$B_{k+1}$ — $n \times n$ matrix $\Rightarrow$ $n \times n$ unknowns.

$n$ equations in the secant approximation $(4)$

We solve for $B_{k+1}$ in an approximate sense using the **BFGS method**. We build an approximation of $B_{k+1}$ out of $\boxed{y_k}$ and $\boxed{B_k s_k}$.

We look at the outer product of $y_k$ with itself :

$$\left[ y_k \, y_k^T \right]_{ij} \qquad \left\{ \begin{array}{l} \underline{n \times 1} \\ 1 \times n \\ \underline{\phantom{}} \\ n \times n. \end{array} \right. \qquad \left( \phantom{x} \right)$$

$$= (y_k)_i (y_k)_j$$

The same outer product for $B_k s_k$.

$$B_{k+1} = B_k + \alpha \, y_k \, y_k^T + \beta \left( B_k s_k \right) \left( B_k s_k \right)^T$$

where $\alpha$ and $\beta$ are TBC.

Theorem 4.1 gives us the values for $\alpha$ and $\beta$: $\enclose{circle}{\color{red}{\text{EXAM}}}$

$$\beta = - \frac{1}{\langle s_k, B_n^T s_n \rangle}, \qquad \alpha = \frac{1}{\langle y_k, s_n \rangle}.$$

**Proof:** Our approximation of $B_{k+1}$ is:

$$B_{k+1} = B_k + \alpha \, y_n y_n^T + \beta \, (B_n S_n)(B_n S_n)^T \quad (5)$$

We require this to satisfy:

$$y_n = Q_{n+1} S_k \quad (6)$$

Sub (5) into (6):

$$y_n = \left[ B_k + \alpha \, y_n y_n^T + \beta \, (B_n S_n)(B_n S_n)^T \right] S_k$$

$$= B_n S_k + \alpha \, y_n \left( y_n^T S_k \right)$$
$$+ \beta \, B_n S_k \left( S_n^T B_n^T S_k \right)$$

$$\Rightarrow y_n = B_n S_k + \alpha \, y_n \langle y_n, S_n \rangle$$
$$+ \beta \, B_n S_n \langle S_n, B_n^T S_n \rangle$$

Re-arrange:

$$0 = y_k \left[ -1 + \alpha \langle y_u, \underline{s}_u \rangle \right]$$
$$+ B_u s_u \left[ 1 + \beta \langle \underline{s}_u, \theta_u^T \underline{s}_u \rangle \right]$$

In general, $y_u$ and $B_u s_u$ are linearly independent, so we require the square brackets to be zero:

$$-1 + \alpha \langle y_u, s_u \rangle = 0$$
$$\implies \alpha = \frac{1}{\langle y_u, s_u \rangle}$$

Also,

$$1 + \beta \langle \underline{s}_u, \theta_u^T s_u \rangle = 0$$
$$\implies \beta = \frac{-1}{\langle \underline{s}_u, \theta_u^T s_u \rangle}$$