

Soft Computing

Unit 3

Conventional Optimization

Derivative-based optimization : descent methods - method of steepest descent - classical Newton's method - step size determination - conjugate gradient methods - analysis of the quadratic case

* Derivative-based Optimization

- consists of gradient-based optimization techniques, which are capable of determining search directions according to an objective function's derivative information
- These include descent algorithms, used for solving minimization problems. These include:
 - (i) Steepest descent method
 - (ii) Newton's method.

* Descent Methods

Objective

- Minimize a real-valued objective function E defined on an n -dimensional input space $\Theta = [\theta_1, \theta_2 \dots \theta_n]^T$
- Find a minimum point $\underline{\Theta} = \underline{\Theta}^*$ that minimizes $\underline{E}(\underline{\Theta})$.

* Iterative Algorithms for computing Descent

- A given objective function E may have a non linear form with respect to an adjustable parameter Θ .
- Due to the complexity of E , an iterative algorithm is used to explore the input space efficiently.
- In iterative descent methods, the next point Θ_{next} is determined by a step down from the current point Θ_{now} , in a direction vector d :

$$\Theta_{\text{next}} = \Theta_{\text{now}} + \eta d$$

η = step-size - also called the learning rate

- An alternative formula is:

$$\Theta_{k+1} = \Theta_k + \eta_k d_k$$

k denotes the current iteration number

Θ_{now} and Θ_{next} represent two consecutive elements in a generated sequence of solution candidates $\{\Theta_k\}$. Θ_k is supposed to converge to a local minimum of Θ^* .

Procedures to compute kth step

- Iterative descent methods compute the kth step $\eta_k d_k$ through two procedures.

(i) determine direction d

(ii) calculate step size η

→ The next point Θ_{next} should satisfy the inequality:

$$E(\Theta_{\text{next}}) = E(\Theta_{\text{now}} + \eta d) < E(\Theta_{\text{now}})$$

→ The optimum step-size can be determined by renomimization.

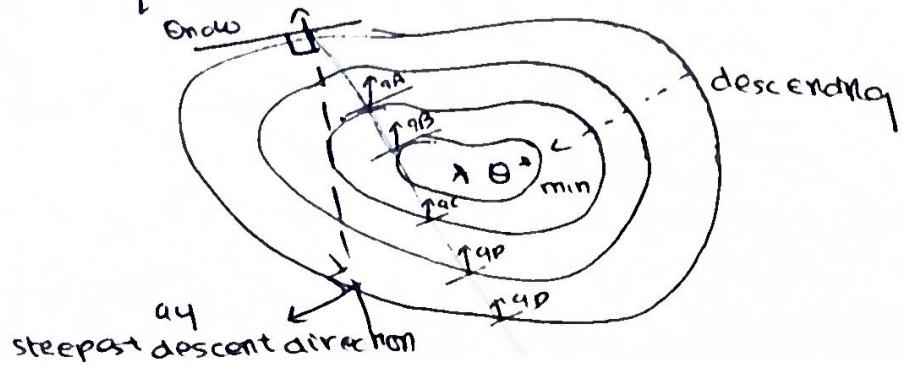
$$\eta^* = \underset{\eta > 0}{\operatorname{argmin}} \Phi(\eta)$$

$$\text{where } \Phi(\eta) = E(\Theta_{\text{now}} + \eta d)$$

* Gradient Descent Method

- Gradient descent is an optimization algorithm used to obtain the optimized network weight and bias values.
- It works by iteratively trying to minimize the cost function
- It works by calculating the gradient of the cost function and moving in the negative direction until the local/global minimum is achieved
- The gradient of a differentiable function $E: \mathbb{R}^n \rightarrow \mathbb{R}$ at Θ is the vector of the first derivatives of E , denoted by g .

$$g(\Theta) = \left[\frac{\partial E(\Theta)}{\partial \Theta_1}, \frac{\partial E(\Theta)}{\partial \Theta_2}, \dots, \frac{\partial E(\Theta)}{\partial \Theta_n} \right]^T$$



The descent is computed as:

$$\Theta_{\text{next}} = \Theta_{\text{now}} - \eta \text{ Gg}$$

Ideally, the value of O_{next} should be such that

$$g(\theta_{\text{next}}) = \frac{\partial E(\theta)}{\partial \theta} \Big|_{\theta=\theta_{\text{next}}} = 0$$

→ Stopping Criteria for Gradient Descent

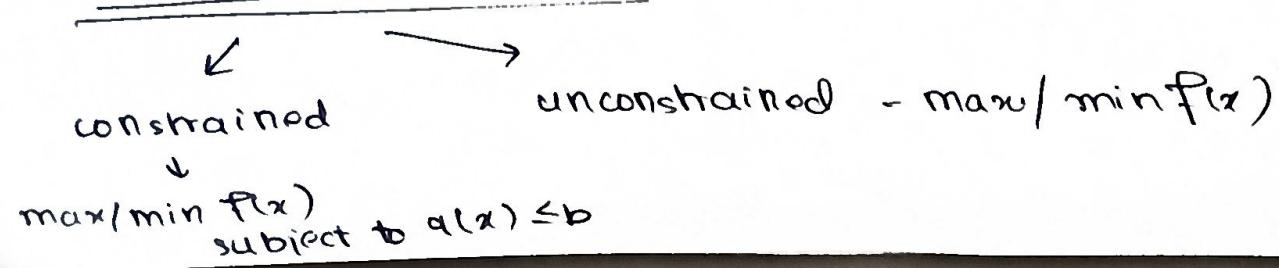
1. The objective function is sufficiently small
 2. The length of the gradient vector g is smaller than a specified value
 3. The specified computing time is exceeded.

Steps in Gradient Descent Method

1. Take the derivative of the loss function for each parameter in it
 2. Pick random values for the parameters
 3. Plug the parameter values into the derivatives (gradient or slope)
 4. Calculate the step size : $\text{step_size} = \text{slope} * \text{learning_rate}$
 5. Calculate the new parameters :

new parameters = old parameters - step-size

Optimization Problems



* Method of Steepest Descent

- It uses the negative of the gradient vector as the direction for minimization
- Start from the initial point x_i , and iteratively move along the steepest descent direction until the optimum point is reached.

Algorithm

1. Calculate s_i at x_i by
$$[s_i = -\nabla f_i]$$
2. Calculate γ_i by using
$$\boxed{\gamma_i = \frac{s_i^T s_i}{s_i^T H_i s_i}}$$
3. Calculate the new point as
$$[x_{i+1} = x_i + \gamma_i s_i]$$
4. Check the optimum of x_{i+1} by
$$[\nabla f(x_{i+1}) \approx 0]$$
5. If met, stop - otherwise implement step 1 again for the new point x_{i+1} .

* Newton Raphson's Method

- primarily used to find the roots of a real-valued function
- When applied to optimization problems, it is used to find the critical points (where the derivative is zero) of a differentiable objective function

For optimization, the method is applied iteratively to update an

initial guess of the optimal solution.

→ The function E is approximated by a quadratic form:

$$E(\theta) \approx E(\theta_{\text{now}}) + g^T (\theta - \theta_{\text{now}}) + \frac{1}{2} (\theta - \theta_{\text{now}})^T H (\theta - \theta_{\text{now}})$$

H = Hessian matrix - w/ the second partial derivatives of $E(\theta)$

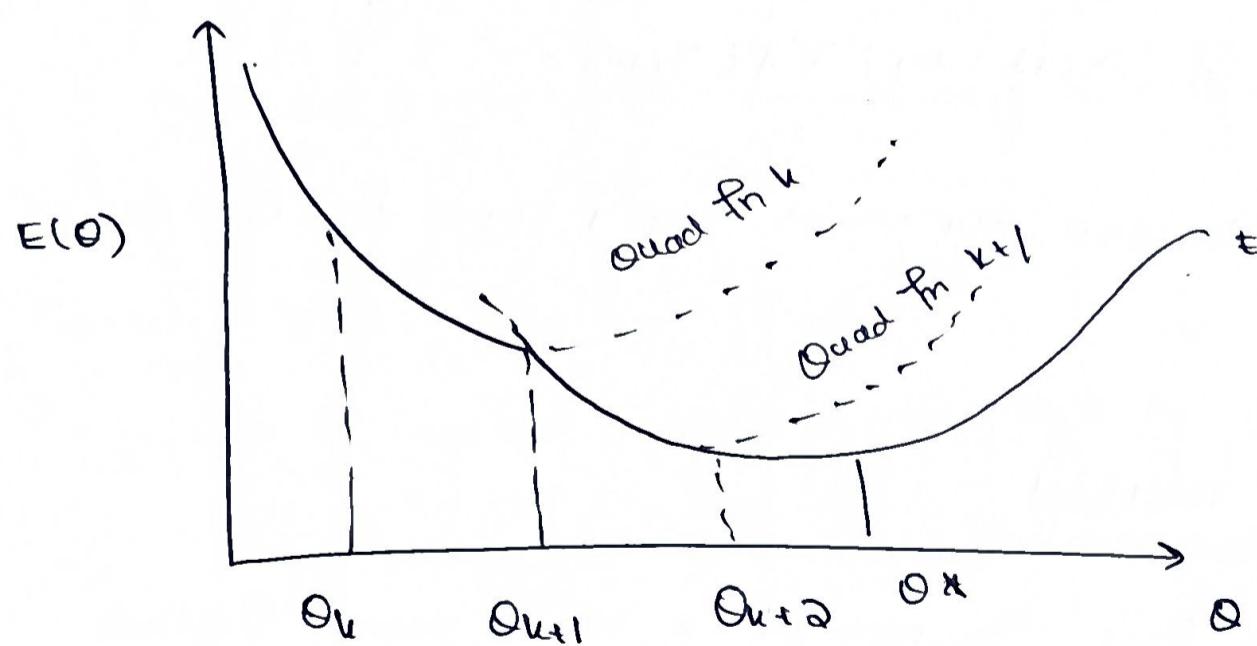
→ Differentiate and set to 0

$$0 = g + H(\hat{\theta} - \theta_{\text{now}})$$

→ rearrange

$$\hat{\theta} = \theta_{\text{now}} - H^{-1}g$$

→ Newton-Raphson method equation



* Step-Halving Procedure

→ The Newton-Raphson equation is:

$$\theta_{\text{next}} = \theta_{\text{now}} - \eta H^{-1}g$$

→ η must satisfy $E(\theta_{\text{next}}) < E(\theta_{\text{now}})$

→ To halve the step size

$$\eta_{\text{new}} = \frac{1}{2} \eta_k$$

can also be

$$O_{\text{next}} = O_{\text{now}} - \eta (1 + 2\gamma)^{-1} g$$

* Step-Size Determination

→ Given an objective function $f(x)$ and a search direction d , the goal is to find the optimal step size α , such that $f(x + \alpha d)$ is minimized

A. Line Search Algorithm

Initialization

- Start with an initial guess for the step size, denoted as α .

- Choose a small positive constant $0 < c_1 < 1$ (like the Armijo condition) to control the sufficient decrease in the objective f_n

Iterative Search

While a stopping criteria is not met:

(i) Evaluate the objective f_n at the point $P(x)$ and the proposed point $f(x + \alpha d)$

(ii) Check if the decrease is sufficient (eq. w/ the Armijo condn)

$$f(x + \alpha d) \leq f(x) + c_1 \alpha \underbrace{\nabla f(x)^T d}_{\text{directional derivative}}$$

directional derivative

(iii) If the condition is met, accept the current α and exit the loop.

(iv) If the condition is not met, update α (e.g. halve it) and repeat the iteration

Output

- The final α is the determined step size.

B. Bisection Method

Initialization

- Choose an initial interval $[\alpha_{\text{low}}, \alpha_{\text{high}}]$ where you expect the optimum step size to lie

- Ensure that $f(\alpha_{\text{low}})$ and $f(\alpha_{\text{high}})$ have opposite signs

Iteration

For each iteration

(i) Compute the midpoint: $\alpha_{\text{mid}} = \frac{\alpha_{\text{low}} + \alpha_{\text{high}}}{2}$

(ii) Evaluate the objective function at the midpoint: $f(\alpha_{\text{mid}})$

(iii) Update the interval:

$f(\alpha_{\text{mid}})$ has the same sign as $f(\alpha_{\text{low}})$

$$\Rightarrow \alpha_{\text{low}} = \alpha_{\text{mid}}$$

$f(\alpha_{\text{mid}})$ has the same sign as $f(\alpha_{\text{high}})$

$$\Rightarrow \alpha_{\text{high}} = \alpha_{\text{mid}}$$

Stopping Criterion

- Repeat the iterations until max. no of iterations

or after achieving a desired precision

Output

- The final α is the determined step size.

c. Golden Search Method



Initialization

- Choose an initial interval $[a, b]$ such that the function is unimodal within the interval.

- Define 2 interior points within the interval:

$$x_1 = b - \left(\frac{b-a}{\phi} \right)$$

ϕ = golden ratio

$$x_2 = a + \left(\frac{b-a}{\phi} \right)$$

≈ 1.618

Iteration

- Evaluate the function at the 2 interior points $f(x_1)$ and $f(x_2)$

- Compare the function values:

if $f(x_1) < f(x_2)$:

$$\text{interval} = [a, x_2]$$

if $f(x_2) < f(x_1)$

$$\text{interval} = [x_1, b]$$

- Repeat this process until the interval becomes sufficiently small or a convergence criterion is met.

Stopping Criterion & Output

- same as B.

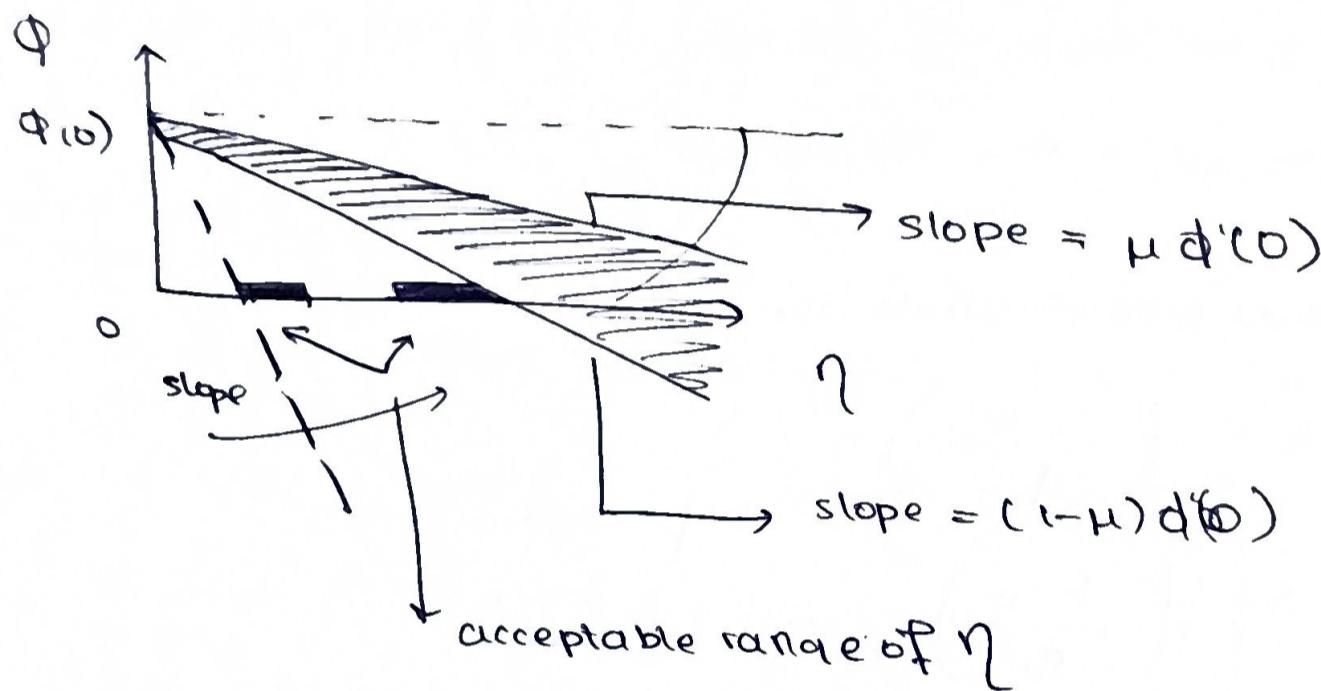
* Termination Rules

A. Goldstein Test

- used particularly with line search algorithms

The condition to be satisfied is:

$$0 \leq \mu \leq \frac{\epsilon(\theta_{\text{next}}) - \epsilon(\theta_{\text{now}})}{\eta g d} \leq 1 - \mu \leq 1$$



$\Phi(\eta)$ should lie in the shaded region.

B. Wolfe Test

→ The following condition must be satisfied so η is not too small:

$$\boxed{\Phi'(\eta) \geq (1-\mu) \Phi'(0)}$$

numericals

CAT-2 - bisection

numerical

grad. descent numerical

C. Armijo Test

→ For η to not be too small

$$\boxed{\Phi(\xi\eta) \geq \Phi(0) + \mu \Phi'(0) \xi \eta \quad 0 < \mu < 1}$$

$$\xi > 1$$

* Conjugate Gradient Method

→ used to solve constrained optimization problems such as energy minimization.

→ This method is called conjugate gradient method since the search direction is obtained by a correction of the negative gradient & conjugate to all the previous directions

Algorithm

Let f be quadratic function $f(x) = \frac{1}{2} x^T A x + b^T x + c$ which we wish to minimize.

1. Initialize

$$i=0$$

$$x_i = x_0$$

$$\text{compute } d_i = d_0 = -\nabla f(x_0)$$

2. Find best step size compute α_i to minimize the function $f(x_i + \alpha_i d_i)$ via the equation

$$\alpha_i = \frac{d_i^T (A x_i + b)}{d_i^T A d_i}$$

3. Update the current guess: $x_{i+1} = x_i + \alpha_i d_i$

4. Update the direction $d_{i+1} = -\nabla f(x_{i+1}) + \beta_i d_i$

$$\beta_i = \frac{\nabla f(x_{i+1})^T A d_i}{d_i^T A d_i}$$

5. Iterate | Repeat 2-4, until we have looked in n directions, where n is the size of the vector space.

* Analysis of the Quadratic Case

→ Consider an objective function E to have a quadratic form:

$$E(\theta) = \frac{1}{2} \theta^T A \theta + b^T \theta + c$$

The gradient would be:

$$g(\theta) = (A\theta + b)$$

To minimize $\Phi(\eta) = E(\theta_{\text{now}} + \eta d)$, set $\Phi'(\eta)$ to 0

$$\Phi'(\eta) = \frac{dE(\theta_{\text{now}} + \eta d)}{d\eta}$$

$$\boxed{\eta^* = -\frac{d^T g}{d^T A d}}$$

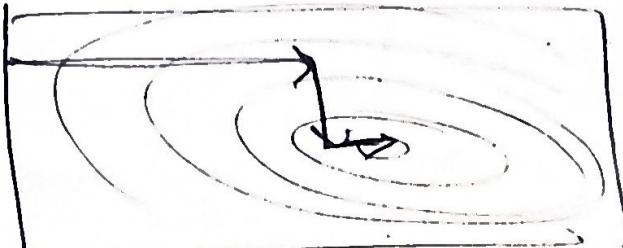
→ Consider a minimization problem of the form:

$$E(\theta) = E(x, u) = x^2 + xu + u^2 - xu$$

$$\nabla E(x, u) = [2x+u-1, x+2u+1]^T$$

Different minimization techniques that can be used are:

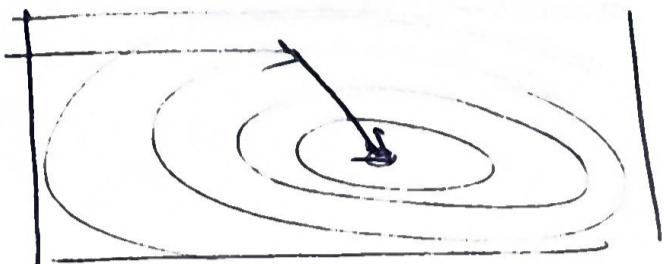
(i) steepest method



(ii) Newton's method



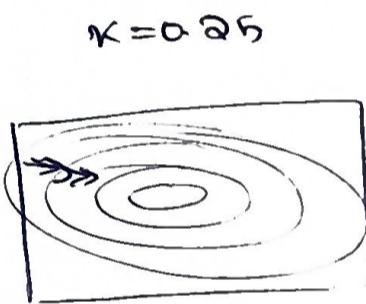
(iii) BFGS quasi-Newton method



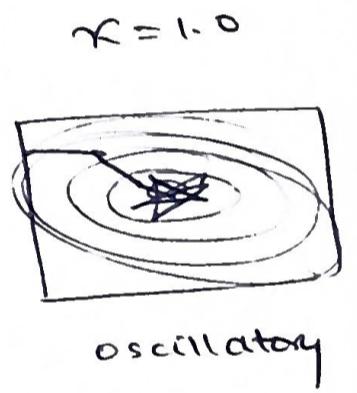
(iv) Fletcher - Reeve's conjugate method

* Steepest Descent Method without line minimization

$$\begin{aligned}\theta_{\text{next}} &= \theta_{\text{now}} - \eta g \\ &= \theta_{\text{now}} - \gamma \frac{g}{\|g\|}\end{aligned}$$



$\gamma = 0.25$



$\gamma = 1.0$

oscillatory

inefficient

update γ as follows:

- ① If the fn. undergoes m consecutive reductions, increase γ by p_1 .
- ② If the fn. undergoes n consecutive combinations of one increase & one reduction, decrease γ by q_1 .

Unit 3 Numericals

1. Compute the minimum using the method of steepest descent.

$$f(x_1, x_2) = x_1 - x_2 + 2x_1^2 + 2x_1x_2 + x_2^2$$

$$x_1 = (0, 0)$$

Step 1 Compute $\nabla f = \left\{ \begin{array}{l} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right\}$

$$\frac{\partial f}{\partial x_1} = 1 + 4x_1 + 2x_2$$

$$\frac{\partial f}{\partial x_2} = -1 + 2x_1 + 2x_2$$

$$\left. \frac{\partial f}{\partial x_1} \right|_{x=(0,0)} = 1$$

$$\left. \frac{\partial f}{\partial x_2} \right|_{x=(0,0)} = -1$$

$$\nabla f = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$s_1 = -\nabla f = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Step 2 - Find the Hessian matrix

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}$$

Step 3 Compute γ

$$\gamma = \frac{s_i^T s_i}{s_i^T H_s s_i}$$

-4+2

$$= \frac{[-1 \ 1] \begin{bmatrix} -1 \\ 1 \end{bmatrix}}{[-1 \ 1] \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix}}$$

$$= \frac{\frac{2}{2}}{[-2 \ 0] \begin{bmatrix} -1 \\ 1 \end{bmatrix}} = \frac{2}{2} = 1$$

1+2 2+2

Step 4 : Compute x_{i+1}

$$\begin{aligned} x_{i+1} &= x_i + \gamma s_i \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \end{aligned}$$

$$x_2 = \boxed{\begin{bmatrix} -1 \\ 1 \end{bmatrix}}$$

Step 5 Optimality check

$$\begin{aligned} \nabla F &= \left\{ \begin{array}{l} \frac{\partial F}{\partial x_1} \\ \frac{\partial F}{\partial x_2} \end{array} \right\} = \left\{ \begin{array}{l} 1 + 4x_1 + 2x_2 \\ -1 + 2x_1 + 2x_2 \end{array} \right\} \\ &= \left\{ \begin{array}{l} 1 - 4 + 2 \\ -1 + 2 - 2 \end{array} \right\} = \begin{bmatrix} -1 \\ -1 \end{bmatrix} \neq 0 \end{aligned}$$

Iteration 2

$$\nabla F = \boxed{\begin{bmatrix} -1 \\ -1 \end{bmatrix}}$$

$$s_2 = -\nabla f_i = \boxed{\begin{bmatrix} +1 \\ +1 \end{bmatrix}}$$

compute γ_2

$$\gamma_2 = \frac{\mathbf{s}_2^T \mathbf{s}_2}{\mathbf{s}_2^T \mathbf{H} \mathbf{s}_2} = \frac{\begin{bmatrix} 1 & +1 \end{bmatrix} \begin{bmatrix} 1 \\ +1 \end{bmatrix}}{\begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ +1 \end{bmatrix}}$$

$$= \frac{2}{\frac{6}{\cancel{4}} \cancel{0}} = \frac{2}{\frac{2}{10}} = \frac{1}{5} = 0.2$$

Q-1

Q-2

compute \mathbf{x}_3

$$\mathbf{x}_3 = \mathbf{x}_2 + \gamma_2 \mathbf{s}_2$$

$$= \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\boxed{\mathbf{x}_3 = \begin{bmatrix} -0.8 \\ 1.2 \end{bmatrix}}$$

optimality check

$$\nabla F = \left\{ \begin{array}{l} \frac{\partial F}{\partial x_1} \\ \frac{\partial F}{\partial x_2} \end{array} \right\} = \left\{ \begin{array}{l} 1 + 4x_1 + 2x_2 \\ 1 + 2x_1 + 2x_2 \end{array} \right\} = \left\{ \begin{array}{l} 0.2 \\ -0.2 \end{array} \right\} \neq 0$$

Iteration 3

$$\nabla F = \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix}$$

$$\mathbf{s}_3 = -\nabla F = \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix}$$

$$\lambda_3 = \frac{s_3^T s_3}{s_3^T H s_3}$$

$$s_3 = \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix}$$

$$\lambda_3 = \frac{\begin{bmatrix} -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix}}{\begin{bmatrix} -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix}}$$

$$\begin{bmatrix} -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix} = 0.04$$

$$\lambda_3 = \frac{-0.08}{\begin{bmatrix} -0.4 & 0 \end{bmatrix} \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix}} = \frac{-0.08}{0.08} = 1$$

-0.8 + 0.4
-0.4

$$x_4 = \begin{bmatrix} -0.8 \\ 1.0 \end{bmatrix} + 1 \begin{bmatrix} -0.2 \\ 0.2 \end{bmatrix}$$

$$x_4 = \begin{bmatrix} -0.8 & -1 \\ 1.0 & 1.0 \end{bmatrix}$$

evaluate optimality

$$\nabla f = \begin{Bmatrix} 1 + 4x_1 + 2x_2 \\ -1 + 2x_1 + 2x_2 \end{Bmatrix} = \begin{Bmatrix} 1 - 4 + 2.8 \\ -1 - 2 + 2.8 \end{Bmatrix} = \begin{bmatrix} -0.2 \\ -0.2 \end{bmatrix}$$

Iteration 4

$$\nabla f = \begin{bmatrix} -0.2 \\ -0.2 \end{bmatrix}$$

$$s_4 = \begin{bmatrix} 0.2 \\ 0.2 \end{bmatrix}$$

$$\lambda_4 = \frac{s_4^T s_4}{s_4^T H s_4}$$

$$= \frac{\begin{bmatrix} 0.2 \\ 0.2 \end{bmatrix} \begin{bmatrix} 0.2 & 0.2 \end{bmatrix}}{\begin{bmatrix} 0.2 \\ 0.2 \end{bmatrix} \begin{bmatrix} 4.2 \\ 0.2 \end{bmatrix} \begin{bmatrix} 0.2 & 0.2 \end{bmatrix}}$$

$$= \frac{0.08}{\begin{bmatrix} 1.2 & 0.8 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.2 \end{bmatrix}} = \frac{0.08}{0.40} = \boxed{\lambda_4 = 0.2}$$

$\frac{0.8 + 0.4}{0.24}$
 $\frac{0.16}{0}$

$$\lambda_5 = \cancel{0.226} + \begin{bmatrix} -1 \\ 1.4 \end{bmatrix} + 0.2 \begin{bmatrix} 0.2 \\ 0.2 \end{bmatrix}$$

$$x_5 = \boxed{\begin{bmatrix} -0.96 \\ 1.44 \end{bmatrix}}$$

evaluate optimality

$$\nabla f(x_5) = \begin{Bmatrix} 0.04 \\ -0.04 \end{Bmatrix} \approx 0$$

optimal values of $x_1 = -0.96$ 2

$x_2 = 1.44$