

Outline. Need for Cholesky Factorization, Cholesky Factorization Algorithm, Coordinate Descent Method, Coordinate Descent Method on Quadratic Functions, Conjugate Directions, Conjugate Direction - Solved Examples, Conjugate Direction Method, Basic Conjugate Direction Algorithm

1 The Cholesky Factorization

Motivation: Validating positive Definiteness:

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a symmetric matrix.

- A is positive definite if there exists a unique lower triangular matrix $\mathbf{L} \in \mathbb{R}^{n \times n}$ with positive diagonal components such that $\mathbf{A} = \mathbf{L}\mathbf{L}^T$
- This decomposition of the matrix into $\mathbf{L}\mathbf{L}^T$ is called Cholesky Decomposition or Cholesky Factorization, and it is a useful criterion to determine if a matrix is positive definite.
- An important issue that arises when employing Levenberg Marquardt algorithm is to validate whether the matrix $\nabla^2 f(\mathbf{x}_k) + \mu(\mathbf{I})$ is positive definite. This challenge can be tackled by Cholesky factorization as discussed above.

Motivation: Solution of $\mathbf{Ax} = \mathbf{b}$:

Cholesky factorization is also used to solve a system of equations.

- Let $\mathbf{A} \in \mathbb{R}^{n \times n}$, symmetric and positive definite.
- Solution of $\mathbf{Ax} = \mathbf{b}$ is $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$
- Instead use Cholesky Factorization of \mathbf{A} . $\mathbf{A} = \mathbf{L}\mathbf{L}^T$
- The given system of equations is $\mathbf{L}\mathbf{L}^T\mathbf{x} = \mathbf{b}$
- Solve the triangular system, $\mathbf{L}\mathbf{u} = \mathbf{b}$ using forward substitution to get \mathbf{u}
- Solve the triangular $\mathbf{L}^T\mathbf{x} = \mathbf{u}$ using backward substitution to get \mathbf{x}^*
- As we saw, Cholesky factorization can also be used to solve a system of equations and hence can be employed to solve the equation $\nabla^2 f(\mathbf{x}_k)\mathbf{d} = -\nabla f(\mathbf{x}_k)$ in general for Newton method.

Cholesky Factorization Algorithm:

- The computation of Cholesky factorization is done using a simple recursive approach.
- Consider the following block matrix partitioning of the matrices \mathbf{A} and \mathbf{L} .

$$\mathbf{A} = \begin{pmatrix} A_{11} & \mathbf{A}_{21} \\ \mathbf{A}_{12} & \mathbf{A}_{22} \end{pmatrix} \quad \mathbf{L} = \begin{pmatrix} L_{11} & \mathbf{0} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{pmatrix}$$

where $A_{11} \in \mathbb{R}$, $\mathbf{A}_{21} \in \mathbb{R}^{(n-1) \times 1}$, $\mathbf{A}_{22} \in \mathbb{R}^{(n-1) \times (n-1)}$, $L_{11} \in \mathbb{R}$, $\mathbf{L}_{21} \in \mathbb{R}^{(n-1) \times 1}$, $\mathbf{L}_{22} \in \mathbb{R}^{(n-1) \times (n-1)}$.

- Since $\mathbf{A} = \mathbf{L}\mathbf{L}^T$, we have

$$\mathbf{A} = \begin{pmatrix} A_{11} & \mathbf{A}_{21} \\ \mathbf{A}_{12} & \mathbf{A}_{22} \end{pmatrix} = \begin{pmatrix} L_{11}^2 & L_{11}\mathbf{L}_{21}^T \\ L_{11}\mathbf{L}_{21} & L_{11}\mathbf{L}_{21}^T + \mathbf{L}_{22}\mathbf{L}_{22}^T \end{pmatrix}$$

- Therefore, in particular $L_{11} = \sqrt{A_{11}}$, $\mathbf{L}_{21} = \frac{1}{\sqrt{A_{11}}} \mathbf{A}_{12}^T$.
- $\mathbf{L}_{22}\mathbf{L}_{22}^T = \mathbf{A}_{22} - \mathbf{L}_{21}\mathbf{L}_{21}^T = \mathbf{A}_{22} - \frac{1}{A_{11}} \mathbf{A}_{12}^T \mathbf{A}_{12}$.
- We are left with the task of Cholesky factorization of $(n-1) \times (n-1)$ matrix $\mathbf{A}_{22} - \frac{1}{A_{11}} \mathbf{A}_{12}^T \mathbf{A}_{12}$.
- We keep following the above procedure and we can get the complete Cholesky factorization.
- The algorithm for Cholesky factorization will find a solution only if all the diagonal elements l_{ii} that are computed during the process are positive, so that computing their square root is possible.
- The positiveness of these elements is equivalent to the property that the matrix to be factored is positive definite.
- Therefore, the Cholesky factorization process can be viewed as a criterion for positive definiteness.

2 Coordinate Descent Method

Motivation: This is one of the simplest techniques for minimizing a function.

Notion: Choose one coordinate direction at a time and optimize the function w.r.t that specific coordinate direction keeping other coordinates constant and the procedure is repeated for other coordinate directions till we get an optimum point.

Consider the problem,

$$\min_{\mathbf{x}} f(\mathbf{x})$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}, f \in \mathbb{C}^1$.

Notion:

- ① For every coordinate variable $x_i, i = 1, \dots, n$, minimize $f(\mathbf{x})$ w.r.t. x_i , keeping the other coordinate variable $x_j, j \neq i$ constant.
- ② Repeat the procedure in step 1 until a stopping condition is satisfied.

Algorithm 1 Coordinate Descent Method

```
(1) Initialize  $\mathbf{x}^0, \epsilon$ , set  $k := 0$ .  
(2) while  $\|\mathbf{g}^k\| > \epsilon$   
    for  $i = 1, \dots, n$   
         $x_i^{new} = \operatorname{argmin}_{x_i} f(\mathbf{x})$   
         $x_i = x_i^{new}$   
    endfor  
endwhile
```

Output : $\mathbf{x}^* = \mathbf{x}^k$, a stationary point of $f(\mathbf{x})$.

- If we assume that a search along any coordinate direction yields a unique minimum point then this method is globally convergent.

Coordinate Descent Method on Quadratic Functions:

Let us understand the discussed notion via an example.

Example 1

Consider the problem,

$$\min_{\mathbf{x}} f(\mathbf{x}) \stackrel{\text{def}}{=} 4x_1^2 + x_2^2$$

We use coordinate descent method with *exact line search* to solve this problem.

- $\mathbf{x}^0 = (-1, -1)^T$
- Let $\mathbf{d}^0 = (1, 0)^T$
- $\mathbf{x}^1 = \mathbf{x}^0 + \alpha^0 \mathbf{d}^0$ where
 $\alpha^0 = \operatorname{argmin}_{\alpha} \phi_0(\alpha) \stackrel{\text{def}}{=} f(\mathbf{x}^0 + \alpha \mathbf{d}^0)$
- $\phi_0(\alpha) = f\left(\begin{smallmatrix} x_1^0 + \alpha d_1^0 \\ x_2^0 + \alpha d_2^0 \end{smallmatrix}\right) = 4(\alpha - 1)^2 + 1$
- $\phi_0'(\alpha) = 0 \Rightarrow \alpha^0 = 1 \Rightarrow \mathbf{x}^1 = (0, -1)^T$
- $\mathbf{d}^1 = (0, 1)^T$, $\mathbf{x}^2 = \mathbf{x}^1 + \alpha^1 \mathbf{d}^1$, $\alpha^1 = \operatorname{argmin}_{\alpha} \phi_1(\alpha) \stackrel{\text{def}}{=} f\left(\begin{smallmatrix} 0 \\ \alpha - 1 \end{smallmatrix}\right) = (\alpha - 1)^2 \Rightarrow \alpha^1 = 1 \Rightarrow \mathbf{x}^2 = (0, 0)^T = \mathbf{x}^*$

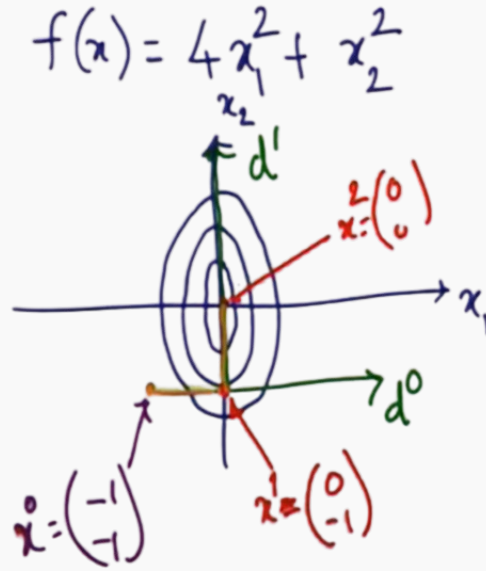


Figure 1: The above diagram depicts the contours of the given function. For the given function, *Coordinate Descent* method reach minimum in two steps - from \mathbf{x}_0 to \mathbf{x}_1 and then from \mathbf{x}_1 to \mathbf{x}_2 .

For the above problem,

- Moving along the coordinate directions and using exact line search gives the solution in at most two steps.
- The Same result is obtained even if \mathbf{d}^0 and \mathbf{d}^1 are interchanged.

Now, let us consider another case via the following example.

Example 2

Consider the problem,

$$\min_{\mathbf{x}} f(\mathbf{x}) \stackrel{\text{def}}{=} 4x_1^2 + x_2^2 - 2x_1x_2$$

We use coordinate descent method with *exact line search* to solve this problem.

- $\mathbf{x}^0 = (-1, -1)^T$
- Let $\mathbf{d}^0 = (1, 0)^T$
- $\mathbf{x}^1 = \mathbf{x}^0 + \alpha^0 \mathbf{d}^0$ where

$$\alpha^0 = \operatorname{argmin}_{\alpha} \phi_0(\alpha) \stackrel{\text{def}}{=} f(\mathbf{x}^0 + \alpha \mathbf{d}^0)$$
- $\phi_0(\alpha) = f \begin{pmatrix} x_1^0 + \alpha d_1^0 \\ x_2^0 + \alpha d_2^0 \end{pmatrix} = 4(\alpha - 1)^2 + 1 + 2(\alpha - 1)$
- $\phi'_0(\alpha) = 0 \Rightarrow \alpha^0 = \frac{3}{4} \Rightarrow \mathbf{x}^1 = (-\frac{1}{4}, -1)^T$
- $\mathbf{d}^1 = (0, 1)^T$, $\mathbf{x}^2 = \mathbf{x}^1 + \alpha^1 \mathbf{d}^1$, $\alpha^1 = \operatorname{argmin}_{\alpha} \phi_1(\alpha) \stackrel{\text{def}}{=} f \begin{pmatrix} -\frac{1}{4} \\ \alpha - 1 \end{pmatrix} = (\alpha - 1)^2 + \frac{\alpha - 1}{2} + \frac{1}{4} \Rightarrow$

$$\alpha^1 = \frac{3}{4} \Rightarrow \mathbf{x}^2 = (-\frac{1}{4}, -\frac{1}{4})^T \neq \mathbf{x}^*$$

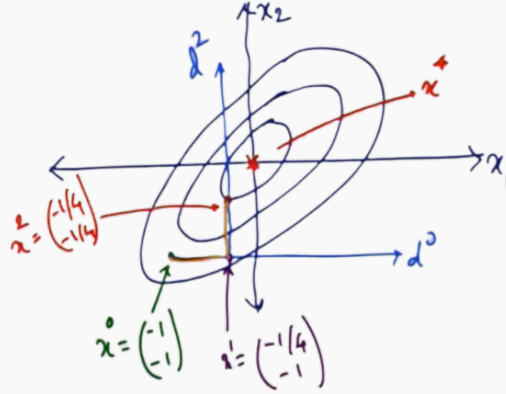


Figure 2: The above diagram depicts the contours of the given function. The path traced by *Coordinate Descent* algorithm in two steps is shown in orange highlight ($x^0 \rightarrow x^1$ and $x^1 \rightarrow x^2$ where $x^2 \neq x^*$). As we can see, the algorithm does not converge to minimum in two steps.

Convergence Criteria:

Example 2 suggests that perhaps we need more iterations to reach the solution (\mathbf{x}^*).

Let us analyze the reason behind reaching the solution in exact two steps in example 1 but not in example 2.

$\min_{\mathbf{x}} f_1(\mathbf{x}) = 4x_1^2 + x_2^2$	$\min_{\mathbf{x}} f_2(\mathbf{x}) = 4x_1^2 + x_2^2 - 2x_1x_2$
$\mathbf{H} = \nabla^2 f_1(\mathbf{x}) = \begin{pmatrix} 8 & 0 \\ 0 & 2 \end{pmatrix}$	$\mathbf{H} = \nabla^2 f_2(\mathbf{x}) = \begin{pmatrix} 8 & -2 \\ -2 & 2 \end{pmatrix}$
\mathbf{x}^* , attained in <i>at most two steps</i> using the coordinate descent method.	\mathbf{x}^* , could not be attained in two steps using the coordinate descent method (if \mathbf{x}^0 is not on one of the principal axes of the elliptical contours).
Terms involving \mathbf{x}_1 and \mathbf{x}_2 are separable.	Terms involving \mathbf{x}_1 and \mathbf{x}_2 are non-separable.

We observe that separability of the objective function in terms of its variables is an important criteria to determine the convergence of Coordinate Descent algorithm.

The separability of the objective function is also reflected in the Hessian matrix.

For $f_1(\mathbf{x})$, \mathbf{H} is a positive definite, *diagonal matrix* whereas for $f_2(\mathbf{x})$, \mathbf{H} is although positive definite but *not a diagonal matrix*. [cite lec 18]

Ques: Is there a way to transform an inseparable function such as $f_2(\mathbf{x})$ into a separable function in terms of its variables, so that the resulting Hessian matrix (\mathbf{H}) becomes a diagonal matrix. Doing so will enable us to use *Coordinate Descent Method* in this new space and get the solution **at the most n-steps for an n-dimensional problem**.

Let us consider a general Quadratic programming problem for minimization.

Conjugate Directions: Solution of General Quadratic Function:

Consider the problem,

$$\min_{\mathbf{x}} f(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{c}^T \mathbf{x} \quad (1)$$

where \mathbf{H} is a symmetric positive definite matrix.

- Let $\{\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}\}$ be a set of linearly independent directions and $\mathbf{x}^0 \in \mathbb{R}^n$

- Any $\mathbf{x} \in \mathbb{R}^n$ can be represented as

$$\mathbf{x} = \mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \quad (2)$$

- Given $\{\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}\}$ and $\mathbf{x}^0 \in \mathbb{R}^n$, substituting RHS of (2) in (1), given problem becomes the problem of minimizing $\Psi(\boldsymbol{\alpha})$, defined as,

$$\frac{1}{2} \left(\mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \right)^T \mathbf{H} \left(\mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \right) + \mathbf{c}^T \left(\mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \right) \quad (3)$$

Solving the problem in eq. (3) wrt to $\boldsymbol{\alpha}$, will result in the solution $\boldsymbol{\alpha}^*$.

Plugging $\boldsymbol{\alpha}^*$ in eq. (2), we get \mathbf{x}^* which will be the solution to the problem in eq. (1). i.e., the general quadratic equation.

Note that since \mathbf{H} is a symmetric positive definite matrix, so the function in eq.(1) is strictly convex and therefore, there exists a strict local minima and that local minima is the \mathbf{x}^* that we got by solving another problem i.e., minimizing $\Psi(\boldsymbol{\alpha})$.

- Define $\mathbf{D} = (\mathbf{d}^0 | \mathbf{d}^1 | \dots | \mathbf{d}^{n-1})$ and $\boldsymbol{\alpha} = (\alpha^0, \alpha^1, \dots, \alpha^{n-1})$.
- Re-writing $\Psi(\boldsymbol{\alpha})$ in terms of this compact notation

$$\Psi(\boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\alpha}^T \underbrace{\mathbf{D}^T \mathbf{H} \mathbf{D}}_{\mathbf{Q}} \boldsymbol{\alpha} + (\mathbf{H} \mathbf{x}^0 + \mathbf{c})^T \mathbf{D} \boldsymbol{\alpha} + \underbrace{\frac{1}{2} \mathbf{x}^{0T} \mathbf{H} \mathbf{x}^0 + \mathbf{c}^T \mathbf{x}^0}_{\text{constant}}$$

While minimizing the above function $\Psi(\boldsymbol{\alpha})$, ignore the constant term.

Furthermore, the Hessian matrix of this quadratic function is $\mathbf{D}^T \mathbf{H} \mathbf{D}$. Let us call that matrix as \mathbf{Q} .

Matrix \mathbf{Q} looks like:

$$\mathbf{Q} = \mathbf{D}^T \mathbf{H} \mathbf{D} = \begin{pmatrix} \mathbf{d}^{0T} \mathbf{H} \mathbf{d}^0 & \mathbf{d}^{0T} \mathbf{H} \mathbf{d}^1 & \dots & \mathbf{d}^{0T} \mathbf{H} \mathbf{d}^{n-1} \\ \mathbf{d}^{1T} \mathbf{H} \mathbf{d}^0 & \mathbf{d}^{1T} \mathbf{H} \mathbf{d}^1 & \dots & \mathbf{d}^{1T} \mathbf{H} \mathbf{d}^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{d}^{n-1T} \mathbf{H} \mathbf{d}^0 & \mathbf{d}^{n-1T} \mathbf{H} \mathbf{d}^1 & \dots & \mathbf{d}^{n-1T} \mathbf{H} \mathbf{d}^{n-1} \end{pmatrix}$$

From *Example1*, we recall that because the $\mathbf{H}(f_1(\mathbf{x}))$ was diagonal, we could apply coordinate descent method and get the solution at the most 2 steps.

Suppose we want to make the \mathbf{Q} hessian matrix diagonal. The way to do that is to make all the non-diagonal elements of this matrix as zero.

- \mathbf{Q} will be **diagonal** matrix if $\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^j = 0, \forall i \neq j$.
- Let $\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^j = 0, \forall i \neq j$.

$$\mathbf{Q} = \mathbf{D}^T \mathbf{H} \mathbf{D} = \begin{pmatrix} \mathbf{d}^{0T} \mathbf{H} \mathbf{d}^0 & 0 & \dots & 0 \\ 0 & \mathbf{d}^{1T} \mathbf{H} \mathbf{d}^1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{d}^{n-1T} \mathbf{H} \mathbf{d}^{n-1} \end{pmatrix}$$

- Therefore,

$$\mathbf{Q}_{ij}^{-1} = \begin{cases} \frac{1}{\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^i} & \text{if } j = i \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \Psi(\boldsymbol{\alpha}) &= \frac{1}{2} \left(\mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \right)^T \mathbf{H} \left(\mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \right) + \mathbf{c}^T \left(\mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^i \mathbf{d}^i \right) \\ &= \frac{1}{2} \sum_i \left[(\mathbf{x}^0 + \alpha^i \mathbf{d}^i)^T \mathbf{H} (\mathbf{x}^0 + \alpha^i \mathbf{d}^i) + 2\mathbf{c}^T (\mathbf{x}^0 + \alpha^i \mathbf{d}^i) \right] + \text{constant} \end{aligned}$$

- $\Psi(\boldsymbol{\alpha})$ is **separable** in terms of $\alpha^0, \alpha^1, \dots, \alpha^{n-1}$

Once we have this separability, it is easy to optimize this objective function in terms of individual alphas.

Let us see how to do that.

$$\Psi(\boldsymbol{\alpha}) = \frac{1}{2} \sum_i \left[(\mathbf{x}^0 + \alpha^i \mathbf{d}^i)^T \mathbf{H} (\mathbf{x}^0 + \alpha^i \mathbf{d}^i) + 2\mathbf{c}^T (\mathbf{x}^0 + \alpha^i \mathbf{d}^i) \right] + \text{constant}$$

$$\frac{\partial \Psi}{\partial \alpha^i} = 0 \Rightarrow \alpha^{i*} = -\frac{\mathbf{d}^{iT} (\mathbf{H} \mathbf{x}^0 + \mathbf{c})}{\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^i}$$

- Therefore,

$$\mathbf{x}^* = \mathbf{x}^0 + \sum_{i=0}^{n-1} \alpha^{i*} \mathbf{d}^i$$

So computation of \mathbf{x}^* becomes easy provided that we have $\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^j = 0$.

Definition (*Conjugate Directions*)

Let $\mathbf{H} \in \mathbb{R}^{n \times n}$ be a symmetric matrix. The vectors $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}$ are said to be ***H-conjugate*** if they are linearly independent and $\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^j = 0 \forall i \neq j$.

With this knowledge, let's try solving for Example 2 again.

Example 2 (*Conjugate Directions*)

Consider the problem,

$$\min_{\mathbf{x}} f(\mathbf{x}) \stackrel{\text{def}}{=} 4x_1^2 + x_2^2 - 2x_1x_2$$

- \mathbf{H} = The matrix $\begin{pmatrix} 8 & -2 \\ -2 & 2 \end{pmatrix}$ is invertible.
- Let $\mathbf{x}^0 = (-1, -1)^T$
- Let $\mathbf{d}^0 = (1, 0)^T$
- $\mathbf{x}^1 = \mathbf{x}^0 + \alpha^0 \mathbf{d}^0$ where

$$\alpha^0 = \underset{\alpha}{\operatorname{argmin}} \phi_0(\alpha) \stackrel{\text{def}}{=} f(\mathbf{x}^0 + \alpha \mathbf{d}^0)$$
- $\phi_0(\alpha) = f\left(\begin{smallmatrix} x_1^0 + \alpha d_1^0 \\ x_2^0 + \alpha d_2^0 \end{smallmatrix}\right) = 4(\alpha - 1)^2 + 1 + 2(\alpha - 1)$
- $\phi'_0(\alpha) = 0 \Rightarrow \alpha^0 = \frac{3}{4} \Rightarrow \mathbf{x}^1 = (-\frac{1}{4}, -1)^T$
- Choose a non-zero direction \mathbf{d}^1 such that $\mathbf{d}^{1T} \mathbf{H} \mathbf{d}^0 = 0$
- Let $\mathbf{d}^1 = (a, b)^T$. Therefore,
 $(a, b) \begin{pmatrix} 8 & -2 \\ -2 & 2 \end{pmatrix} (1, 0) = 0 \Rightarrow 8a - 2b = 0$
- Let $\mathbf{d}^1 = (1, 4)^T$,
- $\mathbf{x}^2 = \mathbf{x}^1 + \alpha^1 \mathbf{d}^1$ where

$$\alpha^1 = \underset{\alpha}{\operatorname{argmin}} \phi_1(\alpha) \stackrel{\text{def}}{=} f\left(\begin{smallmatrix} \alpha - \frac{1}{4} \\ 4\alpha - 1 \end{smallmatrix}\right) = \frac{3}{4} (4\alpha - 1)^2$$

- $\phi'_1(\alpha) = 0 \Rightarrow \alpha^1 = \frac{1}{4}$
- $\mathbf{x}^2 = \mathbf{x}^1 + \alpha^1 \mathbf{d}^1 = (0, 0)^T = \mathbf{x}^*$

A convex quadratic function can be minimized in, *at most*, n steps, provided we search along conjugate directions of the Hessian matrix.

Ques: Given \mathbf{H} , does a set of \mathbf{H} -conjugate vectors exist? If yes, how do we get a set of such vectors?.

Choosing Conjugate Directions:

Let $\mathbf{H} \in \mathbb{R}^{n \times n}$ be a symmetric matrix.

- Do there exist n conjugate directions w.r.t \mathbf{H} ?
 \mathbf{H} is symmetric $\Rightarrow \mathbf{H}$ has n mutually orthogonal eigenvectors.
Let $\boldsymbol{\nu}_1$ and $\boldsymbol{\nu}_2$ be two orthogonal eigenvectors of \mathbf{H} .
 $\therefore \boldsymbol{\nu}_1^T \boldsymbol{\nu}_2 = 0$.

$$\begin{aligned} \mathbf{H}\boldsymbol{\nu}_1 &= \lambda_1 \boldsymbol{\nu}_1 \Rightarrow \boldsymbol{\nu}_2^T \mathbf{H}\boldsymbol{\nu}_1 = \lambda_1 \boldsymbol{\nu}_2^T \boldsymbol{\nu}_1 \\ &\Rightarrow \boldsymbol{\nu}_2^T \mathbf{H}\boldsymbol{\nu}_1 = 0 \\ &\Rightarrow \boldsymbol{\nu}_1 \text{ and } \boldsymbol{\nu}_2 \text{ are } \mathbf{H}\text{-conjugate.} \end{aligned}$$

$\therefore n$ orthogonal eigenvectors of \mathbf{H} are \mathbf{H} -conjugate.

Property of Conjugate Directions:

Lemma

Let \mathbf{H} be a symmetric positive definite $n \times n$ matrix. Let directions $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^k \in \mathbb{R}^n$ ($k \leq n-1$) are \mathbf{H} -conjugate, then they are linearly independent.

- Let \mathbf{H} be a symmetric positive definite matrix and $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}$ be non-zero directions such that

$$\mathbf{d}^{iT} \mathbf{H} \mathbf{d}^j = 0, \quad i \neq j.$$

Are $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}$ linearly independent?

$$\Rightarrow \mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1} \text{ are linearly independent}$$

$$\sum_{i=0}^{n-1} \mu^i \mathbf{d}^i = 0 \Rightarrow \sum_{i=0}^{n-1} \mu^i \mathbf{d}^{jT} \mathbf{H} \mathbf{d}^i = 0 \quad \text{for every } j = 0, \dots, n-1$$

$$\Rightarrow \mu^j \mathbf{d}^{jT} \mathbf{H} \mathbf{d}^j = 0$$

$$\Rightarrow \mu^j = 0 \text{ for every } j = 0, \dots, n-1$$

$$\Rightarrow \mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1} \text{ are linearly independent}$$

Basic Conjugate Direction Algorithm:

Given starting point \mathbf{x}_0 and \mathbf{H} conjugate directions $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{n-1} \in \mathbb{R}^n$, the Conjugate Algorithm works as follows.

- For ($k = 0, 1, \dots, n-1$)
 - $\nabla f(\mathbf{x}_k) = \mathbf{H}\mathbf{x}_k + \mathbf{c}$

$$\begin{aligned}
- \alpha_k &= -\frac{\nabla f(\mathbf{x}_k)^T \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{H} \mathbf{d}_k} \\
- \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{d}_k
\end{aligned}$$

Theorem

Consider minimization problem $\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{c}^T \mathbf{x}$, where \mathbf{H} is a symmetric positive definite matrix. For any starting point \mathbf{x}_0 and \mathbf{H} conjugate directions $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{n-1} \in \mathbb{R}^n$, the **Basic Conjugate Direction Algorithm** converges to the unique \mathbf{x}^* (that solves $\mathbf{H}\mathbf{x}^* + \mathbf{c}$) in n -steps; that is $\mathbf{x}_n = \mathbf{x}^*$.

Proof:

Convergence of Basic Conjugate Direction Method:

- Consider $\mathbf{x}^* - \mathbf{x}_0 \in \mathbb{R}^n$
- Because $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{n-1}$ are \mathbf{H} -conjugate, they are linearly independent. Thus, $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{n-1}$ span \mathbb{R}^n
- So, we can write $\mathbf{x}^* - \mathbf{x}_0$ as a linear combination of $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{n-1}$
- So, there exists $\beta_0, \beta_1, \dots, \beta_{n-1}$ such that $\mathbf{x}^* - \mathbf{x}_0 = \beta_0 \mathbf{d}_0 + \beta_1 \mathbf{d}_1 + \dots + \beta_{n-1} \mathbf{d}_{n-1}$
- Pre-multiplying on both sides with $\mathbf{d}_j^T \mathbf{H}$, we get $\mathbf{d}_j^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_0) = \beta_j \mathbf{d}_j^T \mathbf{H} \mathbf{d}_j$

$$\beta_j = \frac{\mathbf{d}_j^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_0)}{\mathbf{d}_j^T \mathbf{H} \mathbf{d}_j} \quad (1)$$

- Following the conjugate direction algorithm, we get

$$\begin{aligned}
\mathbf{x}_k &= \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{d}_{k-1} \\
&= \mathbf{x}_0 + \sum_{i=0}^{k-1} \alpha_i \mathbf{d}_i \\
\Rightarrow \mathbf{x}_k - \mathbf{x}_0 &= \sum_{i=0}^{k-1} \alpha_i \mathbf{d}_i \\
\mathbf{x}^* - \mathbf{x}_0 &= \mathbf{x}^* - \mathbf{x}_k + \mathbf{x}_k - \mathbf{x}_0 \\
&= (\mathbf{x}^* - \mathbf{x}_k) + \sum_{i=0}^{k-1} \alpha_i \mathbf{d}_i
\end{aligned} \quad (2)$$

- Pre-multiplying on both sides with $\mathbf{d}_k^T \mathbf{H}$, we get

$$\begin{aligned}
\mathbf{d}_k^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_0) &= \mathbf{d}_k^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_k) + \sum_{i=0}^{k-1} \alpha_i \mathbf{d}_k^T \mathbf{H} \mathbf{d}_i \\
\mathbf{d}_k^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_0) &= \mathbf{d}_k^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_k) \quad (3)
\end{aligned}$$

- But, $\mathbf{H}\mathbf{x}^* = -\mathbf{c}$ (because \mathbf{x}^* is optimal)

$$\nabla f(\mathbf{x}_k) = \mathbf{H}\mathbf{x}_k + \mathbf{c}$$

- So, $\mathbf{H}\mathbf{x}^* - \mathbf{H}\mathbf{x}_k = -\mathbf{c} - \nabla f(\mathbf{x}_k) + \mathbf{c} = -\nabla f(\mathbf{x}_k)$
- Using this in eq. (3), we get,

$$\mathbf{d}_k^T \mathbf{H} (\mathbf{x}^* - \mathbf{x}_0) = -\mathbf{d}_k^T \nabla f(\mathbf{x}_k) \quad (4)$$

- Using eq. (4) in (1),

$$\beta_j = -\frac{\mathbf{d}_j^T \nabla f(\mathbf{x}_j)}{\mathbf{d}_j^T \mathbf{H} \mathbf{d}_j} = \alpha_j \quad (5)$$

- But we know that,

$$\begin{aligned} \mathbf{x}^* - \mathbf{x}_0 &= \sum_{i=0}^{n-1} \beta_i \mathbf{d}_i = \sum \alpha_i \mathbf{d}_i \quad (\text{using eq. (5)}) \\ \Rightarrow \mathbf{x}^* &= \mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i \end{aligned}$$

References

- [1] Introduction to Nonlinear Optimization by Amir Beck (Ch. 5 & Ch. 10)
- [2] [NPTEL Lecture 18 Conjugate Directions](#)