

Lecture-9

Conjugate Direction Algorithm
(Solution of Linear System or
Minimization of A Quadratic
Function)

Conjugate Gradient

- Linear conjugate gradient: for solving linear systems $Ax=b$ with PD matrix, A .
 - Exact solution in n steps (Hestenes & Stiefel, 1950s)
 - Approximate solution in fewer than n steps
- Non-linear conjugate gradient: for solving large-scale non-linear optimization problems.
 - Fletcher and Reeves, 1964
 - Polk-Ribiere, 1969

Conjugate Gradient

A is symmetric PD. (1)

Or minimize the following function:

(2)

$r(x)$ is the residual

$S = \{p_0, p_1, \dots, p_{n-1}\}$ The set S is conjugate wrt A if

$$p_i^T A p_j = 0 \quad \forall i \neq j$$

Linear Independence

S is linearly independent

$$\text{if } \mathbf{s}_0 p_0 + \mathbf{s}_1 p_1 + \dots + \mathbf{s}_{n-1} p_{n-1} = 0$$

$$\text{then } \mathbf{s}_0 = \mathbf{s}_1 = \mathbf{s}_2 = \dots = \mathbf{s}_{n-1} = 0$$

Conjugate set is also linearly independent.

$$p_i^T A p_j = 0 \quad \forall i \neq j$$

Conjugate Direction Method

$$x_{k+1} = x_k + \mathbf{a}_k p_k \quad \text{Line search}$$

$$p_i^T A p_j = 0 \quad \forall i \neq j$$

$$\mathbf{f}(x) = \frac{1}{2} x^T A x - b^T x \quad \text{1D minimizer of a quadratic function}$$

$$\mathbf{a}_k = -\frac{\nabla \mathbf{f}_k^T p_k}{p_k^T A p_k}$$

Convergence Rate of Steepest Descent

$$\frac{d}{d\mathbf{a}} f(x_k - \mathbf{a}_k) = \frac{d}{d\mathbf{a}} \left(\frac{1}{2} (x_k - \mathbf{a}_k)^T Q (x_k - \mathbf{a}_k) - b^T (x_k - \mathbf{a}_k) \right) = 0$$

$$= -(x_k - \mathbf{a}_k)^T Q g_k + b^T g_k = 0$$

$$-x_k^T Q g_k + \mathbf{a}_k^T Q g_k + b^T g_k = 0$$

$$\mathbf{a}_k^T Q g_k = x_k^T Q g_k - b^T g_k$$

$$\mathbf{a} = \frac{x_k^T Q g_k - b^T g_k}{g_k^T Q g_k}$$

$$\mathbf{a} = \frac{(x_k^T Q - b^T) g_k}{g_k^T Q g_k} \quad \nabla f(x) = Qx - b$$

From Lecture-5

$$\mathbf{a} = \frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k}$$

$$x_{k+1} = x_k - \mathbf{a}_k \nabla f_k$$

$$x_{k+1} = x_k - \frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k} \nabla f_k$$

Conjugate Direction Method

$$\mathbf{a} = \frac{x_k^T Q g_k - b^T g_k}{g_k^T Q g_k}$$

$$\mathbf{a} = \frac{(x_k^T A - b^T)(-p_k)}{(-p_k)^T A (-p_k)}$$

$$\mathbf{a}_k = -\frac{\nabla f_k^T p_k}{p_k^T A p_k} \quad \nabla f(x) = Ax - b = r(x)$$

$$\mathbf{a}_k = -\frac{r_k^T p_k}{p_k^T A p_k} \quad p_i^T A p_j = 0 \quad \forall i \neq j$$


Theorem 5.1

For any x^0 the sequence $\{x_k\}$ generated by the conjugate direction algorithm, converges to the solution x^* of the linear system in at most n steps.

- Sequence $\{x_k\}$
- Linearly independent vectors
- Conjugate vectors

Proof

$$x_{k+1} = x_k + \mathbf{a}_k p_k \quad \mathbf{a}_k = -\frac{r_k^T p_k}{p_k^T A p_k}$$


$$x_k = x_0 + \mathbf{a}_0 p_0 + \mathbf{a}_1 p_1 + \dots + \mathbf{a}_{k-1} p_{k-1}$$

$$x_k - x_0 = \mathbf{a}_0 p_0 + \mathbf{a}_1 p_1 + \dots + \mathbf{a}_{k-1} p_{k-1}$$

Proof

S is linearly independent

Therefore:

$$x^* - x_0 = \mathbf{s}_0 p_0 + \mathbf{s}_1 p_1 + \dots + \mathbf{s}_{n-1} p_{n-1}$$

$$p_k^T A(x^* - x_0) = p_k^T A(\mathbf{s}_0 p_0 + \mathbf{s}_1 p_1 + \dots + \mathbf{s}_{n-1} p_{n-1})$$

$$p_k^T A(x^* - x_0) = (0 + 0 + \dots + \mathbf{s}_k p_k^T A p_k + \dots + 0) \quad \text{conjugate}$$

$$\mathbf{s}_k = \frac{p_k^T A(x^* - x_0)}{p_k^T A p_k}$$

Proof

$$x_{k+1} = x_k + \mathbf{a}_k p_k \quad \mathbf{a}_k = -\frac{r_k^T p_k}{p_k^T A p_k}$$

$$x_k = x_0 + \mathbf{a}_0 p_0 + \mathbf{a}_1 p_1 + \dots + \mathbf{a}_{k-1} p_{k-1}$$

$$x_k - x_0 = \mathbf{a}_0 p_0 + \mathbf{a}_1 p_1 + \dots + \mathbf{a}_{k-1} p_{k-1}$$

$$p_k^T A(x_k - x_0) = 0$$

$$p_k^T A x_k = p_k^T A x_0$$

$$p_k^T A(x^* - x_0) = p_k^T A(x^* - x_k) = p_k^T (b - A x_k) = -p_k^T r_k$$

$$p_k^T A(x^* - x_0) = -p_k^T r_k$$

Proof

$$p_k^T A(x^* - x_0) = -p_k^T r_k$$

$$\mathbf{s}_k = \frac{p_k^T A(x^* - x_0)}{p_k^T A p_k} \quad \mathbf{a}_k = -\frac{r_k^T p_k}{p_k^T A p_k}$$

Therefore:

$$\mathbf{s}_k = \mathbf{a}_k$$

QED

Interpretation of Theorem 5.1

If A is a diagonal matrix, then we can minimize (1-D) the function along coordinate axes in n iterations.

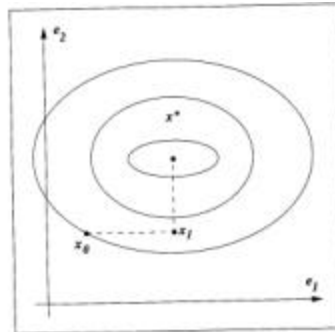


Figure 5.1 Successive minimizations along the coordinate directions find the minimum of a quadratic with a diagonal Hessian in n iterations.

Interpretation of Theorem 5.1

If A is not a diagonal matrix, then we can not minimize the function along coordinate axes in n iterations.

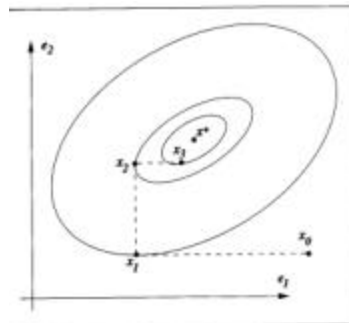


Figure 5.2 Successive minimization along coordinate axes does not find the solution in n iterations, for a general convex quadratic.

Transformed Problem

Let

$$\hat{x} = S^{-1}x \quad \text{where} \quad S = [p_0, p_1, \dots, p_{n-1}]$$

$$f(x) = \frac{1}{2}x^T Ax - b^T x$$

$$J(\hat{x}) = f(S\hat{x}) = \frac{1}{2}\hat{x}^T (S^T A S) \hat{x} - (S^T b)^T \hat{x} \quad \text{By conjugacy } S^T A S \text{ is a diagonal matrix.}$$

Now we can minimize along coordinate directions in transformed space.

However, each coordinate direction in transformed space correspond to the conjugate direction in the original space due to $\hat{x} = S^{-1}x$

Therefore, we conclude the conjugate direction algorithm converges in n steps.

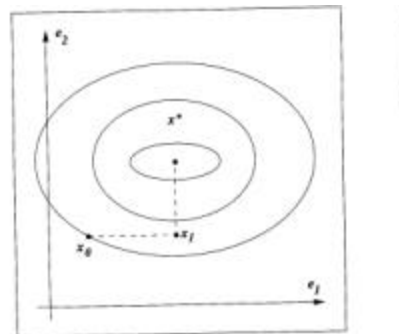


Figure 5.1 Successive minimizations along the coordinate directions of a quadratic with a diagonal Hessian in n iterations.

When Hessian is diagonal, each coordinate minimization correctly determines of the components of the solution x^* . Therefore, after k 1-D minimizations, the quadratic has been minimized on the subspace spanned by e_1, e_2, \dots, e_k .

Theorem 5.2

Let x_0 be any starting point and suppose that the sequence $\{x_k\}$ is generated by the conjugate direction algorithm. Then

$$r_k^T p_i = 0 \quad \text{for } i = 0, \dots, k-1$$

and x_k is minimizer of $f(x) = \frac{1}{2}x^T Ax - b^T x$ over the set

$$\{x \mid x = x_0 + \text{span} \{p_0, \dots, p_{k-1}\}\} \quad (3)$$

Proof

First show that a point \tilde{x} minimizes f over the set (3) if and only if

$$r(\tilde{x})^T p_i = 0 \quad \text{for } i = 0, \dots, k-1$$

$$\{x \mid x = x_0 + \text{span} \{p_0, \dots, p_{k-1}\}\}$$

Where

Let $h(\mathbf{s}) = f(x_0 + \mathbf{s}_0 p_0 + \dots + \mathbf{s}_{k-1} p_{k-1})$

$$\mathbf{s} = (\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_{k-1})$$

Since $h(\mathbf{s})$ is strictly convex quadratic, it has a unique minimizer:

$$\frac{\partial h}{\partial \mathbf{s}_i} = 0, \quad i = 0, \dots, k-1$$

$$\nabla f(x_0 + \mathbf{s}_0^* p_0 + \dots + \mathbf{s}_{k-1}^* p_{k-1})^T p_i = 0 \quad i = 0, \dots, k-1 \quad \text{Chain rule}$$

$r(x)$ is the residual

$$r(\tilde{x})^T p_i = 0 \quad i = 0, \dots, k-1$$

Proof

$$\nabla \mathbf{f}(x) = Ax - b = r(x) \quad x_{k+1} = x_k + \mathbf{a}_k p_k$$

$$r_{k+1} = r_k + \mathbf{a}_k A p_k$$

$$r_k = r_{k-1} + \mathbf{a}_{k-1} A p_{k-1} \quad (\text{A})$$

Use induction:

True for $k=1$

From (A)

$$r_1 = r_0 + \mathbf{a}_0 A p_0$$

$$r_1^T p_0 = (r_0 + \mathbf{a}_0 A p_0)^T p_0$$

$$r_1^T p_0 = r_0^T p_0 + \mathbf{a}_0 p_0^T A p_0$$

$$r_1^T p_0 = 0$$

Because

$$\mathbf{a}_k = -\frac{r_k^T p_k}{p_k^T A p_k}$$

Proof

$$r_k = r_{k-1} + \mathbf{a}_{k-1} A p_{k-1} \quad (\text{A})$$

Assume true for $k-1$

From (A)

$$p_{k-1}^T r_k = p_{k-1}^T r_{k-1} + \mathbf{a}_{k-1} p_{k-1}^T A p_{k-1} = 0$$

$$\mathbf{a}_k = -\frac{r_k^T p_k}{p_k^T A p_k}$$

Definition

And

$$p_i^T r_k = p_i^T r_{k-1} + \mathbf{a}_{k-1} p_i^T A p_{k-1} = 0 \quad i = 0, \dots, k-2$$

$$r_k^T p_i = 0 \quad \text{for } i = 0, \dots, k-2$$

Conjugacy

induction

Therefore $r_k^T p_i = 0$ for $i = 0, \dots, k-1$ QED

How do we select conjugate directions

- Eigenvalues of A are mutually orthogonal and conjugate wrt to A .
- Gram-Schmidt process to produce conjugate directions instead of orthogonal vectors.

Basic Properties of the CG

Each direction is chosen to be a linear combination of the steepest descent direction and the previous direction.

$$p_k = -\nabla f_k + \mathbf{b}_k p_{k-1}$$

$$p_k = -r_k + \mathbf{b}_k p_{k-1}$$

$$p_{k-1}^T A p_k = -r_k^T A p_{k-1} + \mathbf{b}_k p_{k-1}^T A p_{k-1}$$

$$\mathbf{b}_k = \frac{r_k^T A p_{k-1}}{p_{k-1}^T A p_{k-1}}$$

Algorithm 5.1

Given x_0 ;

set $r_0 \leftarrow Ax_0 - b$, $p_0 \leftarrow -r_0$, $k \leftarrow 0$

p_0 is steepest descent

While $r_k \neq 0$

$$\mathbf{a}_k \leftarrow -\frac{r_k^T p_k}{p_k^T A p_k};$$

$$x_{k+1} \leftarrow x_k + \mathbf{a}_k p_k;$$

$$r_{k+1} \leftarrow Ax_{k+1} - b;$$

$$\mathbf{b}_{k+1} \leftarrow \frac{r_{k+1}^T A p_k}{p_k^T A p_k};$$

$$p_{k+1} \leftarrow -r_{k+1} + \mathbf{b}_{k+1} p_k;$$

$$k \leftarrow k + 1;$$

end (while)