

**Numerical Optimization**  
**Prof. Shirish K. Shevade**  
**Department of Computer Science and Automation**  
**Indian Institute of Science, Bangalore**

**Lecture - 19**  
**Conjugate Gradient Method**

(Refer Slide Time: 00:53)

Consider the problem:

$$\min_x f(x) \stackrel{\text{def}}{=} \frac{1}{2}x^T Hx + c^T x, \quad H \text{ symmetric positive definite matrix.}$$

Let  $d^0, d^1, \dots, d^{n-1}$  be  $H$ -conjugate.  $\therefore d^0, d^1, \dots, d^{n-1}$  are linearly independent.

Let  $B^k$  denote the subspace spanned by  $d^0, d^1, \dots, d^{k-1}$ .

Clearly,  $B^k \subset B^{k+1}$ .

Let  $x^0 \in \mathbb{R}^n$  be any arbitrary point.

Let  $x^{k+1} = x^k + \alpha^k d^k$  where  $\alpha^k$  is obtained by doing exact line search:

$$\alpha^k = \arg \min_{\alpha} f(x^k + \alpha d^k)$$

Claim:

$$x^k = \arg \min_{x \in B^k} f(x)$$

s.t.  $x \in x^0 + B^k$

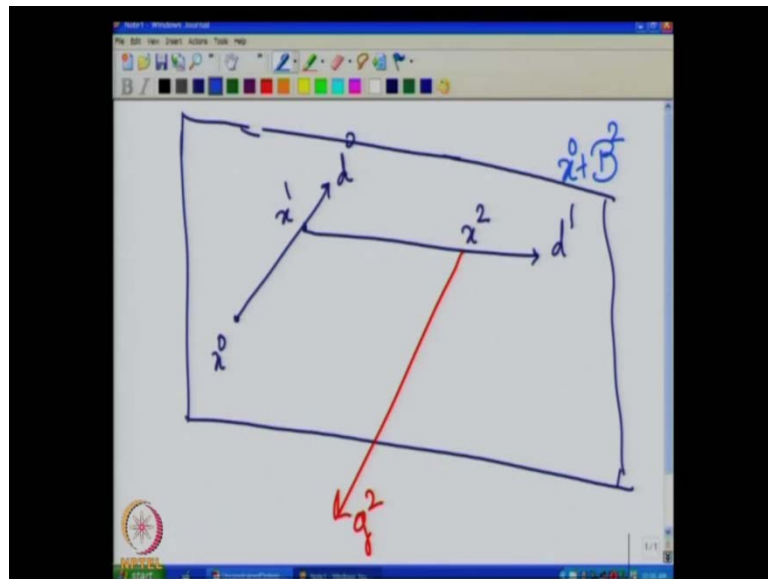
Welcome back to this series of lectures on optimization. So, in the last class, we started studying about conjugate directions and we saw some of the properties of conjugate directions. So in today's class, we will discuss about expanding subspace theorem, and the use of conjugate gradient method to quadratic as well as non-quadratic functions. So let us first consider the quadratic function, which is to be optimized. So, let the function to be optimized be half  $x$  transpose  $H$   $x$  plus  $c$  transpose  $x$ , where  $H$  is a symmetric positive definite matrix. And let us assume that we have  $n$  conjugate directions  $d^0$  to  $d^{n-1}$ , and we already saw in the last class that these directions which are  $H$  conjugate are linearly independent. So they form a basis for  $n$ -dimensional space.

Now, let us take the first  $k$  vectors  $d^0$  to  $d^{k-1}$ , and since they are independent they will form a basis of a  $k$ -dimensional space, which is a subspace of  $\mathbb{R}^n$ . So, let us denote by  $B^k$  that subspace which is spanned by  $d^0$  to  $d^{k-1}$ . Now, vector  $d^k$  is independent of this, so clearly the space spanned by the  $k$  vectors  $d^0$  to  $d^k$

minus 1 is a subset of the subspace spanned by  $k + 1$  vectors  $d_0$  to  $d_k$ , because  $d_0$  to  $d_k$  are linearly independent.

Now, let us take any arbitrary point in the  $n$ -dimensional space, and will denote by  $x_k$ . Let us denote  $x_{k+1}$  by  $x_k + \alpha_k d_k$  as usual, where  $\alpha_k$  is obtained by doing the exact line search that is  $\alpha_k$  is optimum minimizing the function  $f$  of  $x_k + \alpha d_k$  with respect to  $\alpha$ . Now,  $x_k$  and  $d_k$  are known and  $\alpha$  is a variable here and the value of  $\alpha$  that gives the minimum will be denoted by  $\alpha_k$ . Now, under these conditions that  $d_0$  to  $d_{n-1}$  are independent, now the claim is that then we find  $x_k$  as by using the conjugate direction method, the minimum value of  $f$  of  $x$  in the space, spanned by  $x_0$  plus  $d_k$ , that is nothing but  $x_k$  or in other words.

(Refer Slide Time: 03:42)



So, suppose, we have a point  $x_0$  from which we start, we get a direction which is say  $d_0$  and then we get a point which is  $x_1$  on this direction by doing the line search. And then we have a direction  $d_1$  along which we do again line search, to get a point  $x_2$ . Now, the important point is that,  $d_0$  and  $d_1$  are independent, so they will span a two-dimensional space. So this is two-dimensional space spanned by  $d_0$  and  $d_1$ , and now if we considered  $g_2$ , the gradient direction at  $x_2$ . Now,  $g_2$  will be perpendicular to the space spanned by  $d_0$  and  $d_1$  and containing  $x_0$ , so this is the space which contains  $x_0$  plus  $B_2$ .

So, the space containing  $x_0$ , and the space spanned by  $d_0$  and  $d_1$  which are linearly independent, now the expanding subspace theorem says that the direction  $g_2$  that you get at  $x_2$  is orthogonal to this space. Now, if we choose the direction  $d_2$ , which is in the which is along the or which is the combination of  $g_2$  and the previous directions. Then what happens is that when we do the minimization of  $f$  of  $x$  in the now three-dimensional space spanned by  $d_0$   $d_1$  and  $B_2$ . Then the claim is that, we will have got the minimum of  $f$   $x$  with respect to the three-dimensional space  $d_0$   $d_1$  and  $d_2$ , and if we keep on repeating this, for  $n$  steps. Then we will have found the minimum of  $f$  of  $x$  with respect to the  $n$ -dimensional space and that is what we will see now.

So, the claim is that,  $x_k$  that we get is the minimum of  $f$  of  $x$  subject to the constrained that  $x$  belongs to the space spanned by  $d_0$  to  $d_{k-1}$ , and the space containing  $x_{naught}$ . Now, remember that  $x_{naught}$  is any arbitrary point, so it really does not matter what our  $x_{naught}$  is as long as we get the conjugate directions, then we can say that the  $x_k$  that we get is the minimum of  $f$  of  $x$  subject to this space. And if we keep on repeating this then finally, what happenses that at the end of any tritions we will have found  $x_n$  which will be equal to  $x_{star}$ .

(Refer Slide Time: 07:32)

Exact line search:

$$\alpha^k = \arg \min_{\alpha \in \mathbb{R}} f(x^k + \alpha d^k)$$

Therefore,

$$\nabla f(x^k + \alpha^k d^k)^T d^k = 0 \Rightarrow g^{k+1^T} d^k = 0 \quad \forall k = 0, \dots, n-1$$

$$x^k = x^{k-1} + \alpha^{k-1} d^{k-1} = x^j + \sum_{i=j}^{k-1} \alpha^i d^i \quad (j = 0, \dots, k-1)$$

$$\therefore Hx^k + c = Hx^j + c + \sum_{i=j}^{k-1} \alpha^i H d^i$$

$$\therefore g^k = g^j + \sum_{i=j}^{k-1} \alpha^i H d^i$$

$$\therefore g^{k^T} d^{j-1} = g^{j^T} d^{j-1} + \sum_{i=j}^{k-1} \alpha^i d^{i^T} H d^{j-1} = 0$$

Therefore,  $g^{k^T} d^j = 0 \quad \forall j = 0, \dots, k-1$  or  $g^k \perp_{\mathbb{R}} B^k$

So, let us true this thing first, now we know that alpha k is obtained using exact line search. So, this is the one-dimensional optimization problem, and we know that the way to solve it is that take the gradient of this function and equated to 0, because it is one-

dimensional optimization problem. And note that  $f$  is also convex quadratic functions in this case. So at the new point, the gradient of, of the function at a new point is perpendicular to  $d_k$ .

And therefore, we can say that  $g_{k+1}^T d_k$  is equal to 0, for all  $k$  going from 0 to  $n-1$ , so or in other words  $g_{k+1}$  is perpendicular to  $d_k$ . So in our previous example  $g_2$  will be perpendicular to  $d_1$ , but the claim is that  $g_2$  will be not only perpendicular to  $d_1$ . But also,  $d_0$  or in general  $g_k$  will be perpendicular to  $d_0$  to  $d_{k-1}$  that is the claim, and let us see how to show that. So, now, let us look at a general step where  $x_k$  is equal to  $x_{k-1}$ , plus  $\alpha_{k-1} d_{k-1}$  and  $x_k$ . Therefore, can be written as some of  $x_j$  and  $\sum \alpha_i d_i$ ,  $i$  going from  $j$  to  $k-1$ . So  $x_k$  can be written in terms of  $x_j$  and the linear sum of the remaining conjugate vectors  $d_i$ . Now, now if we premultiply both sides by  $H$  and add  $c$ , so what we get is  $Hx_k + c$ , that will be equal to  $Hx_j + c + \sum_{i=j}^{k-1} \alpha_i H d_i$ .

So note that, we want to find out what happens to  $g_{k+1}^T d_j$  for any  $j$  which is less than  $k$ . We know the  $g_{k+1}^T d_k$  is 0, or  $g_{k+1}^T d_j$  is equal to 0 when  $j$  is equal to  $k$ . But what happens when  $j$  is less than  $k$ , and that is why we do this multiplication so that now we can treat this as a gradient of the objective function. Note that, we are working with the quadratic function of the type  $f$  of  $x$  is equal to  $\frac{1}{2} x^T H x + c^T x$ . So, these quantity denotes the gradient of the objective function at  $x_k$ , and this quantity denotes the gradient of the objective function at  $x_j$ , so we can write  $g_k$  is equal to  $g_j + \sum \alpha_i H d_i$ . And therefore, now we want to see what happens to  $g_k^T d_j$ , so let us multiply both sides by take the dot product of both sides with respect to  $d_j$ . So  $g_k^T d_j$  is equal to  $g_j^T d_j + \sum \alpha_i d_i^T H d_j$ .

Now,  $g_j^T d_j$  is equal to 0, because of this result. Here we have shown that  $g_{k+1}^T d_k$  is equal to 0. So similarly,  $g_j^T d_{j-1}$  is equal to 0, so this quantity is 0 and now what about this quantity? So you will see that the  $d$ s are the conjugate directions so  $d_i^T H d_{j-1}$ , where  $i$  goes from  $j$  to  $k-1$ . So certainly, none of the  $i$  is going to take the value  $g_j^T d_{j-1}$ , so all these quantities with respect to all  $i$  is  $d_i^T H d_{j-1}$  will be 0.

So together this whole quantity will be 0, and this will be true for any  $j$ . So therefore, what we have is  $g^k$  transpose  $d^j$  is equal to 0 for all  $j$  going from 0 to  $k$  minus 1. So, here we were shown that  $g^k$  transpose  $d^{j-1}$  is equal to 0 for all  $j$ , and from this we know that  $g^k$  transpose  $d^{k-1}$  also equal to 0. So, in other words, what we have shown is that  $g^k$  is perpendicular to all the  $d^j$   $j$  going from 0 to  $k$  minus 1, or in other words  $g^k$  is perpendicular to the space spanned by  $d^0$  to  $d^{k-1}$  and that space is nothing but  $B^k$ . And if we consider the, space containing  $B^k$  as space spanned by  $d^0$  to  $d^{k-1}$  that is  $B^k$  that contains  $x^0$ , then  $g^k$  will be perpendicular to that entire space. So this is an interesting observation, and therefore we can say that  $g^k$  is perpendicular to  $B^k$ .

(Refer Slide Time: 13:04)

Note that for every  $j = 0, \dots, n-1$ ,

$$\alpha^j = \arg \min_{\alpha} f(x^j + \alpha d^j)$$

$$\therefore f(x^j + \alpha^j d^j) \leq f(x^j + \mu^j d^j), \quad \mu^j \in \mathbb{R}$$

$$\therefore f(x^j) + \alpha^j g^{jT} d^j + \frac{1}{2} \alpha^{j2} d^{jT} H d^j \leq f(x^j) + \mu^j g^{jT} d^j + \frac{1}{2} \mu^{j2} d^{jT} H d^j$$

We need to show that  $f(x^k) \leq f(x) \forall x \in x^0 + B^k$  or

$$f(x^0 + \sum_{j=0}^{k-1} \alpha^j d^j) \leq f(x^0 + \sum_{j=0}^{k-1} \mu^j d^j), \quad \mu^j \in \mathbb{R} \forall j.$$

That is,

$$f(x^0) + \sum_{j=0}^{k-1} (\alpha^j g^{jT} d^j + \frac{1}{2} \alpha^{j2} d^{jT} H d^j) \leq f(x^0) + \sum_{j=0}^{k-1} (\mu^j g^{jT} d^j + \frac{1}{2} \mu^{j2} d^{jT} H d^j)$$

where  $\mu^j \in \mathbb{R} \forall j$ .

Now, again let us go back to our earlier alpha determination of alpha  $j$  using exact line search. And that it shown here where we want to find alpha  $j$  by minimizing  $f$  of  $x^j$  plus alpha  $d^j$  with respect to alpha. Now, what we want to show is that,  $f$  of  $x$  what we want to show is that, now if we minimize  $f$  of  $x$  subject to the constraint that  $x$  belongs to this space, then that is nothing but  $x^k$ . Or in other words, so  $x^2$  here is the minimum of  $f$  of  $x$ , subject to the this space which is given by  $x^0$  plus  $B^2$ . And now suppose if you find the direction  $d^3$  using  $g^2$  and other the  $d^0$  and  $d^1$ , such that  $d^3$  is  $d^2$  is  $H$  conjugate to  $d^0$  and  $H$  conjugate to  $d^1$ .

Then, when we find the point  $x_3$  that point  $x_3$  will be the minimum of  $f$  of  $x$  in the space, which contains  $x_0$  and spanned by  $d_0, d_1$  and  $d_2$ , so that is going to be our next claim. So let us see how to do that, now because of the exact line search we know that  $f$  of  $x_j$  plus  $\alpha_j d_j$  is less than or equal to  $f$  of  $x_j$  plus  $\mu_j d_j$  when  $\mu_j$  belongs to  $\mathbb{R}$ . This is because of the exact line search that along the line, the value of the function at  $x_j$  plus  $\alpha_j d_j$  will be always less than or equal to the value of the function at any point on the line  $x_j$  plus  $\mu_j d_j$  where  $\mu_j$  belongs to  $\mathbb{R}$ . And since the function  $f$  is quadratic, we can use the Taylor series to write  $f$  of  $x$  as to write the or to expand  $f$  of  $x_j$  plus  $\alpha_j d_j$  as  $f$  of  $x_j$  plus  $\alpha_j$  times the gradient of  $f$  of  $x_j$  transpose  $d_j$  plus half  $\alpha_j$  square  $d_j$  transpose  $H d_j$ . And that will be less than or equal to similar expansion of  $f$ , then using the Taylor series, so we can write this for all  $\mu_j$  belong to  $\mathbb{R}$ .

So, what we know is that, along one direction, the function value is minimum at a particular point. So, if we consider if we start from  $x_0$  and take  $d_0$  as the direction  $x_1$  is the minimum is  $x_1$  is the point at which the function attains the minimum value. So as far as this direction is concerned no other point in this direction has a function value which is greater than which is strictly less than  $f$  of  $x_1$ . Now, the same thing holds here also, that if we start from  $x_1$  and take this direction now no other point in this direction has a function value which is strictly less than  $f$  of  $x_2$ . So, this is true for one-dimensional case, but what is the guarantee that the  $x_2$  that we here is actually the minimum of  $f$  in the space spanned by  $d_0$  and  $d_1$ , and that is what we want to show. Now, what we need to show is that  $f$  of  $x_k$  is less than or equal to  $f$  of  $x$  for all  $x$  in the space spanned by  $d_0$  to  $d_{k-1}$  and that space containing  $x_{naught}$ .

So, this is what we want to prove because that is our claim, so that means what we want to show if we write  $x_k$  as  $x_0$  plus  $\sum \alpha_j d_j$  going from 0 to  $k-1$ . Remember that these  $\alpha_j$  is are obtained using the line search in the respective directions  $d_j$  and that quantity, we want to show that that quantity is less than or equal to  $f$  of  $x$ . Now, any general  $x$  in this space can be written as  $x_0$  plus  $\sum \mu_j d_j$ , where  $d_j$  is the basis of this  $k$ -dimensional space  $j$  going from 0 to  $k-1$  and  $\mu_j$  is a real number so this is what we want to show. Now, if we expand this using Taylor series, because this is a quadratic function, what we can do is that we can write this as,  $f$  of the  $x_0$  plus  $\sum \alpha_j g_0$  transpose  $d_j$  plus half  $\alpha_j$  square  $d_j$  transpose  $H d_j$ .

Now, let us compare these 2 quantities, so suppose here, we write, we combine all the  $j$  is accumulated so far. So, we take the sum as,  $\sum_{j=0}^{k-1}$  here and similarly, here and then combine these 2 quantities then we see that there is a difference of these 2 quantities, so here we have  $d_j^T d_j$  and here we have  $g_0^T d_j$ . Now, suppose we show that  $d_j^T d_j$  is nothing but  $g_0^T d_j$  then we can replace this quantity by  $g_0^T d_j$  and then cancel  $f(x_j)$  from both sides and add  $f(x_0)$  on both sides and then what we get is this quantity? So what we need to show is that  $g_0^T d_j$  should be equal to  $d_j^T d_j$ . So, in other words, the dot product of  $g_j$  and  $d_j$  should be same as dot product of  $g_0$  and  $d_j$ , and if that holds let us see what happens. So suppose  $g_j^T d_j$  is nothing but  $g_0^T d_j$  remember that we had got this inequality earlier.

(Refer Slide Time: 19:33)

For every  $j = 0, \dots, n-1$ ,

$$f(x^j) + \alpha^j g^j{}^T d^j + \frac{1}{2} \alpha^2 d^j{}^T H d^j \leq f(x^j) + \mu^j g^j{}^T d^j + \frac{1}{2} \mu^2 d^j{}^T H d^j$$

Suppose  $g^j{}^T d^j = g^0{}^T d^j \forall j$

$$\therefore \alpha^j g^j{}^T d^j + \frac{1}{2} \alpha^2 d^j{}^T H d^j \leq \mu^j g^j{}^T d^j + \frac{1}{2} \mu^2 d^j{}^T H d^j \quad \forall j$$

Therefore,

$$f(x^0) + \sum_{j=0}^{k-1} (\alpha^j g^j{}^T d^j + \frac{1}{2} \alpha^2 d^j{}^T H d^j) \leq f(x^0) + \sum_{j=0}^{k-1} (\mu^j g^j{}^T d^j + \frac{1}{2} \mu^2 d^j{}^T H d^j)$$

$$\therefore f(x^0) + \sum_{j=0}^{k-1} \alpha^j d^j \leq f(x^0) + \sum_{j=0}^{k-1} \mu^j d^j, \quad \mu^j \in \mathbb{R} \forall j$$

$$\therefore f(x^k) \leq f(x), \quad \forall x \in x^0 + B^k$$

So, this is the inequality that we had got by doing line search along the direction  $d_j$ . Now, suppose this holds, then what we can do is that we can write this quantity as so the  $f(x_j)$  gets cancelled. And we can write this as  $\alpha_j g_0^T d_j$  this quantity remains the same is less than or equal to  $\mu_j g_0^T d_j$  so the  $g_j^T d_j$  in either side is replaced by  $g_0^T d_j$ , and the reason for doing it is that then it we are more closed to the the desired thing.

Now, we just add  $f(x_0)$  on both sides, and then some over all  $j$  is. So if we add  $f(x_0)$  on both sides and some over  $j$  going from 0 to  $k$  minus, because for each  $j$  this inequality

holds. So if we sum them up, then what we get is  $f(x^0) + \sum_{j=0}^{k-1} \alpha_j \nabla f(x^j)^T d^j$  less than or equal to  $f(x^0) + \sum_{j=0}^{k-1} \alpha_j \nabla f(x^0)^T d^j$ . So, in other words, what we have shown is that at the end of  $k$  iterations, we will have reached the point  $x^k$  and that  $f(x^k)$  is less than or equal to  $f(x)$  for any  $x$  in that  $k$ -dimensional space containing  $x^0$  and spanned by  $d^0$  to  $d^{k-1}$ . So this is a very important observation, that  $f(x^k)$  is less than or equal to  $f(x)$  for all  $x^0$  for all  $x$  belong in the space  $x^0 + B^k$ . So all these was possible, because of the fact that  $\nabla f(x^j)^T d^j = \nabla f(x^0)^T d^j$ , now we want to see whether this really holds or not.

(Refer Slide Time: 21:43)

We need to show that

$$g^j{}^T d^j = g^0{}^T d^j \quad \forall j$$

Consider,  $x^j = x^0 + \sum_{i=0}^{j-1} \alpha^i d^i$ .

$$\therefore Hx^j + c = Hx^0 + c + \sum_{i=0}^{j-1} \alpha^i H d^i$$

$$\therefore g^j = g^0 + \sum_{i=0}^{j-1} \alpha^i H d^i$$

$$\therefore g^j{}^T d^j = g^0{}^T d^j + \sum_{i=0}^{j-1} \alpha^i d^i{}^T H d^i$$

So, let us now show that  $\nabla f(x^j)^T d^j$  is equal to  $\nabla f(x^0)^T d^j$ , now consider any point  $x^j$  which can be written as  $x^0 + \sum_{i=0}^{j-1} \alpha^i d^i$ , so which can be written as a point in  $j$ -dimensional space spanned by  $d^0$  to  $d^{j-1}$  and containing  $x^0$ . Now, again we want to show the relationship between  $\nabla f(x^j)^T d^j$  and  $\nabla f(x^0)^T d^j$ . So, we need a gradient information, so we pre-multiply throughout by  $H$  and add  $c$ . So, what we will get is the gradient of  $f$  at  $x^j$  on the left hand side, and this is the gradient of  $f$  at  $x^0$ , so that means  $\nabla f(x^j)$  is equal to  $\nabla f(x^0) + \sum_{i=0}^{j-1} \alpha^i H d^i$ . Now, what we want to do is that we want to show what happens  $\nabla f(x^j)^T d^j$  so  $\nabla f(x^j)^T d^j$  is equal to  $\nabla f(x^0)^T d^j + \sum_{i=0}^{j-1} \alpha^i d^i{}^T H d^i$ .



So, now you will see that, the  $i$  where is from 0 to  $j$  minus 1, so  $d_j^T H d_i$  or  $d_i^T H d_j$  will be equal to 0, because these are  $H$  conjugate directions. So, this entire quantity becomes 0 and therefore, what we get is  $g_j^T d_j$  is equal to  $g_0^T d_j$ . And because of that, we can say that if this holds then we can write  $g_j^T d_j$  to be  $g_0^T d_j$  and then this holds for all  $j$  so sum them up. Then we get this quantity then add  $f(x_0)$  on both sides the inequality does not change the direction. And therefore, we can say that  $f(x_k)$  is clearly less than or equal to  $f(x)$ , for all  $x$  in the space  $x_0 + B_k$ . Now, if you repeat this step  $n$  times, at the end of  $n$  iterations we will have got the  $n$  conjugate directions  $d_0$  to  $d_{n-1}$ . And we know that they are linearly independent so at the end of  $n$  iterations whatever  $x_n$  that we get will be, the minimum of the objective function.

(Refer Slide Time: 24:18)

**Expanding Subspace Theorem**

Consider the problem to minimize  $f(x) \stackrel{\text{def}}{=} \frac{1}{2}x^T H x + c^T x$  where  $H$  is symmetric positive definite matrix. Let  $d^0, d^1, \dots, d^{n-1}$  be  $H$ -conjugate and let  $x^0 \in \mathbb{R}^n$  be any initial point. Let

$$\alpha^k = \arg \min_{\alpha \in \mathbb{R}} f(x^k + \alpha d^k), \quad \forall k = 0, \dots, n-1$$

and  $x^{k+1} = x^k + \alpha^k d^k, \quad \forall k = 0, \dots, n-1.$

Then, for all  $k = 0, \dots, n-1,$

- 1  $g^k^T d^j = 0, \quad j = 0, \dots, k$
- 2  $g^k^T d^k = g^0^T d^k$
- 3

$$x^{k+1} = \arg \min_x f(x)$$

s.t.  $x \in x^0 + B^k$

So, we have this expanding subspace theorem, if we consider the problem to minimize  $f$  of  $x$  to be half  $x$  transpose  $H$   $x$  plus  $c$  transpose  $x$ , where  $H$  is a symmetric positive definite matrix. And suppose  $d_0$  to  $d_{n-1}$   $H$  conjugate directions, so which clearly from the basis  $\mathbb{R}^n$  and  $x_0$  is any initial point then and if  $\alpha^k$  is chosen using exact line search and  $x_{k+1}$  is said to  $x_k + \alpha^k d^k$ . Then what we have is,  $g_k^T d_j$  is equal to 0, so  $g_k$  is perpendicular to the space spanned by  $d_0$  to, this should be  $k-1$ . So,  $g_k^T d_j$  is equal to 0 for all  $j$  going from 0 to  $k-1$ .  $g_k^T d_k$  is equal to  $g_0^T d_k$ . And the important thing is that  $x_{k+1}$  is nothing but  $\arg \min$  of  $f$  of  $x$ ,  $x_0 + B^k$ .

(Refer Slide Time: 25:36)

Given a set of  $n$  directions,  $d^0, d^1, \dots, d^{n-1}$  which are  $H$ -conjugate and  $x^0 \in \mathbb{R}^n$ , it is easy to determine  $\alpha^i, \forall i = 0, \dots, n-1$ ,

$$\alpha^i = -\frac{d^{i\top}(Hx^0 + c)}{d^{i\top}Hd^i}$$

and get

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

- How do we construct the  $H$ -conjugate directions,  $d^0, d^1, \dots, d^{n-1}$ ?
- Given the  $H$ -conjugate directions,  $d^0, d^1, \dots, d^{k-1}$ , how do we determine  $\alpha^k$  where

$$\alpha^k = \arg \min_{\alpha} f(x^k + \alpha d^k)?$$

NPTEL

So, this is the important theorem, now we are given a set of  $n$  directions  $d^0$  to  $d^{n-1}$ , which are  $H$  conjugate, and  $x^0$  is in  $\mathbb{R}^n$ . Now, we have seen that it is easy to determine  $\alpha^i$  star to be, minus  $d^i$  transpose  $Hx^0$  plus  $c$  divided by  $d^i$  transpose  $Hd^i$ , and then we know that  $x^*$  can be obtained using this all this is possible provided  $d^0$  to  $d^{n-1}$  or  $H$  conjugate.

Now, the important question is that, how do we construct the  $H$  conjugate directions  $d^0$  to  $d^{n-1}$ , and then we will come up with a iterative procedure, which will construct  $d^0, d^1, d^2$  and so on. So, at the end of  $k$  iterations, how do you find out  $\alpha^k$  star, based on the point  $x^k$ , because this  $\alpha^i$  star depends on  $Hx^0$  plus  $c$  and all the previous conjugate directions. So, so how do we get  $\alpha^k$  at a current point  $x^k$ , so given the  $H$  conjugate directions  $d^0$  to  $d^{k-1}$  how do we determine  $\alpha^k$  which basically the minimizer of this so does there exist a close form expression depending  $\alpha^k$ . So, let us first look at this and then go to the other question where we want to construct the  $H$  conjugate directions.

(Refer Slide Time: 27:16)

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

$$\therefore d^k T H (x^* - x^0) = \alpha^k d^k T H d^k$$

$$\therefore \alpha^k = \frac{d^k T H (x^* - x^0)}{d^k T H d^k}$$

Suppose that after  $k$  iterative steps and obtaining  $k$   $H$ -conjugate directions,

$$x^k - x^0 = \sum_{i=0}^{k-1} \alpha^i d^i$$

$$\therefore d^k T H (x^k - x^0) = 0$$

Now, note that  $x^*$  can be written as a combination as a sum of  $x^0$  and the linear combination of  $d^i$  is where  $i$  going from 0 to  $n - 1$ . So, if we premultiply throughout by  $d^k T H$  then only 1 quantity on the right side remains the rest of the quantity become 0 and the quantity which remains is  $\alpha^k d^k T H d^k$ . And therefore,  $\alpha^k$  is equal to  $d^k T H (x^* - x^0)$  divided by  $d^k T H d^k$ . Now, this  $\alpha^k$  is still dependent on  $x^*$  and  $x^0$  and we want to get the dependence of  $\alpha^k$  only on  $x^k$ , so let us see how to do that. Now, after  $k$  iterative steps, and obtaining  $k$   $H$  conjugate directions what we have is  $x^k$  is nothing, but  $x^0$  plus the linear combination of the  $k - 1$  vectors  $d^0$  to  $d^{k-1}$ . Now, clearly, we can say that  $d^k T H (x^k - x^0)$  will be 0, because  $d^k$  will be conjugate to all  $d^i$  is going from 0 to  $k - 1$ .

(Refer Slide Time: 28:46)

Given,  $d^{kT} H(x^* - x^k) = 0,$

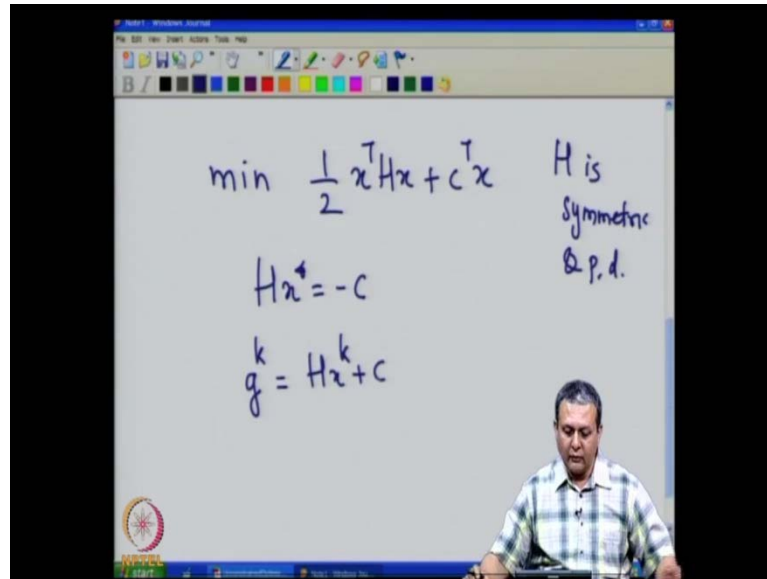
$$\begin{aligned} \therefore \alpha^k &= \frac{d^{kT} H(x^* - x^k + x^k - x^0)}{d^{kT} H d^k} \\ &= \frac{d^{kT} (Hx^* - Hx^k)}{d^{kT} H d^k} \\ &= \frac{d^{kT} (-c - Hx^k)}{d^{kT} H d^k} \\ &= \frac{g^{kT} d^k}{d^{kT} H d^k} \end{aligned}$$

Therefore,

$$\alpha^k = \frac{g^{kT} d^k}{d^{kT} H d^k}$$

Now, we let us use this fact in this formula, so given that  $d^k \text{ transpose } H \text{ into } x^k \text{ minus } x^0 \text{ equal to } 0$ . We can write  $\alpha^k$  in the previous formula the previously we had written that  $\alpha^k$  is equal to  $d^k \text{ transpose } H \text{ into } x^* \text{ minus } x^k$  by  $d^k \text{ transpose } H d^k$ , so we just add subtract and add  $x^k$  in this formula. Now, if you expand this, so take the first 2 terms and combine them with this and then take the last 2 terms, and combine them with this so  $d^k \text{ transpose } H \text{ into } x^k \text{ minus } x^0$  that is equal to 0. So, the quantity involving the the last 2 terms vanishes, so what we are left with is  $d^k \text{ transpose } H \text{ into } x^* \text{ minus } x^k$ . Now, we know that the function that we want to optimize is half  $x^T H x$  plus  $c \text{ transpose } x$ .

(Refer Slide Time: 29:53)



So, the function that we want to minimize is half  $x$  transpose  $H$   $x$  plus  $c$  transpose  $x$ ,  $H$  is symmetric and positive definite. So, taking the so we take the gradient equated to 0, so  $H$   $x$  star should be equal to minus  $c$ . So, we use this fact, so this  $H$   $x$  star is nothing but minus  $c$ . And so we replace  $H$   $x$  star by minus  $c$  and again if we look at the gradient of the function so  $g^k$  will be  $H$   $x^k$  plus  $c$ . So, if we look at the the formula here so what we have is negative of  $H$   $x^k$  plus  $c$ . So, which is nothing but negative of  $g^k$ , so therefore, we have got  $\alpha^k$  using  $g^k$  and  $d^k$  and  $H$ .

So, we have got  $\alpha^k$  using  $g^k$   $d^k$  and  $H$ , so you will see that this  $\alpha^k$  that we have got now does not depend on  $x$  naught, but it just depends on the the current gradient value. So, compare this  $\alpha^k$  with the  $\alpha^k$  that we have got earlier, that  $d^k$  transpose  $H$  into  $x$  star minus  $x$  naught or even before that we had got  $\alpha^k$  which depended on  $x$  naught so we got the formula for  $\alpha^k$  based only on  $g^k$  and  $d^k$ .

(Refer Slide Time: 31:51)

Suppose  $\{-g^0, -g^1, \dots, -g^{n-1}\}$  is a linearly independent set of vectors.

Use Gram-Schmidt procedure to determine the  $H$ -conjugate vectors,  $d^0, d^1, \dots, d^{n-1}$ .

- Let  $d^0 = -g^0$
- In general,

$$d^k = -g^k + \sum_{j=0}^{k-1} \beta^j d^j, \quad k = 1, \dots, n-1$$

But we want  $d^0, d^1, \dots, d^{n-1}$  to be  $H$ -conjugate vectors.

$$d^i T H d^k = -d^i T H g^k + \sum_{j=0}^{k-1} \beta^j d^i T H d^j, \quad i = 0, \dots, k-1$$

$$\therefore 0 = -d^i T H g^k + \beta^i d^i T H d^i, \quad i = 0, \dots, k-1$$

$$\therefore \beta^i = \frac{g^{k T} H d^i}{d^{i T} H d^i}$$

Now, the important question of finding the conjugate directions,  $H$  conjugate directions  $d_0$  to  $d_{n-1}$ . Now, let us assume that we have got  $n$  linearly independent set of vectors and those we are going to denote by  $-g_0$  to  $-g_{n-1}$ . Now, these are the negatives of the gradient vectors that we have got at every iteration, and let us assume that they are linearly independent. If they are linearly dependent then at some point of time we have the gradient is 0, so gradient is 0 for quadratic function means that we have found the solution. So, this is the case where in  $n-1$  iterations we have not found the solution. And finally, at the end of last iteration, what we have is  $-g_0$  to  $-g_{n-1}$  forming a linearly independent set of vectors. Now, we have already seen that Gram-Schmidt procedure can be used to find out, the orthonormal basis given a set of linearly independent vectors.

Now, that procedure can be extended to find out the  $H$  conjugate directions  $d_0$  to  $d_{n-1}$ , and this is the basis that we are going to use. Remember that, we are assuming that  $-g_0$  to  $-g_{n-1}$  are linearly independent later on we will show that this is indeed case. So, to begin with, we start with  $d_0$  to be the negative of the gradient at  $x_0$  and in general we will write the Gram-Schmidt procedure as  $d_k$  to be  $-g_k$  plus, linear combination of  $k$  vectors  $d_0$  to  $d_{k-1}$ . Now, this  $\beta_j$  is obtained using Gram-Schmidt procedure, and what we want is that  $d_0$  to  $d_{n-1}$  to be  $H$  conjugate. So, if we premultiply this with  $d_i^T H$ , there  $i$  varies from 0 to  $k-1$ .

So, then what we get is  $d^i$  transpose  $H d^k$  is nothing but minus  $d^i$  transpose  $H g^k$  plus this quantity now  $d^i$  transpose  $H d^k$  is 0 because  $i$  is going from 0 to  $k-1$  so this holds for every  $i$  going from 0 to  $k-1$ . So, we get 0 equal to minus  $d^i$  transpose  $H g^k$  plus, only 1 term in this expression remains and that corresponds to the vector  $d^i$ . And therefore what we get is the expression for  $\beta_i$  and which is nothing, but  $g^k$  transpose  $H d^i$  divided by  $d^i$  transpose  $H d^i$ .

So, this is the usual Gram-Schmidt procedure, at  $i$  mean to get  $d^k$  using minus  $g^k$  which is the linearly independent vector in from this set. So, at the  $k$  iteration we get 1 vector from this set, and we use the combination of all the previous vectors to get  $d^k$  now when we talk about  $H$  conjugate directions, this formula becomes much simpler and that is what we will see now. So, we will see that the  $\beta_j$  will be 0 for all  $j$  going from 0 to  $k-2$ , and only the  $k-1$  direction comes into picture.

(Refer Slide Time: 35:43)

$$\therefore d^k = -g^k + \sum_{j=0}^{k-1} \left( \frac{g^{kT} H d^j}{d^{jT} H d^j} \right) d^j$$

We now need to show that  $\{-g^0, -g^1, \dots, -g^{n-1}\}$  is a *linearly independent* set of vectors.

Note that  
 $\text{span}\{d^0, d^1, \dots, d^{k-1}\} = \text{span}\{-g^0, -g^1, \dots, -g^{k-1}\}$   
 We have already shown that

$\{d^0, d^1, \dots, d^{k-1}\}$  are  $H$ -conjugate  $\Rightarrow g^k \perp B^k$   
 $\therefore -g^k \perp \text{span}\{d^0, d^1, \dots, d^{k-1}\}$   
 $\therefore -g^k \perp \text{span}\{-g^0, -g^1, \dots, -g^{k-1}\}$

Therefore,  $\{-g^0, -g^1, \dots, -g^{n-1}\}$  is a *linearly independent* set of vectors.

So, let us see how to get that, so if we rewrite that formula  $d^k$  equal to minus  $g^k$  plus sigma  $\beta_j d^j$  where  $\beta_j$  is actually replace by this. Now, first let us show that minus  $g^0$  to minus  $g^{n-1}$ , they form a linearly independent set of vectors. Now, note that, the span of  $d^0$  to  $d^{k-1}$  is same as span of minus  $g^0$  to minus  $g^{k-1}$ , because the way  $d^k$  is formed is using minus  $g^k$  and all the previous vectors. So, if all  $g^j$  are independent, up to  $k-1$  then span of  $d^0$  to  $d^{k-1}$  is same as span of minus  $g^0$  to minus  $g^{k-1}$ .

And we have already shown the  $d^0$  to  $d^{k-1}$  are H conjugate that is  $g^k$  perpendicular to, and if they are H conjugate then  $g^k$  is perpendicular to  $B^k$ , this is what we have already shown. So, the new direction that we get is perpendicular to the space spanned by  $d^0$  to  $d^{k-1}$ , and that essentially means that  $g^k$  will be perpendicular to the space spanned by  $d^0$  to  $d^{k-1}$ . So, which means that,  $g^k$  is independent of  $d^0$  to  $d^{k-1}$ , so therefore  $d^0$  to  $d^{n-1}$  forms a linearly independent set of a vectors.

(Refer Slide Time: 37:30)

Now, consider

$$d^0 = -g^0$$

$$d^k = -g^k + \sum_{j=0}^{k-1} \underbrace{\left( \frac{g^{kT} H d^j}{d^{jT} H d^j} \right)}_{\beta_j} d^j \quad \forall k = 1, \dots, n-1$$

Note that  $x^{j+1} = x^j + \alpha^j d^j$  and  $g^{j+1} = g^j + \alpha^j H d^j$ .  
Therefore,

$$H d^j = \frac{1}{\alpha^j} (g^{j+1} - g^j)$$

Thus,

$$d^k = -g^k + \sum_{j=0}^{k-1} \left( \frac{g^{kT} (g^{j+1} - g^j)}{d^{jT} (g^{j+1} - g^j)} \right) d^j$$

$$= -g^k + \left( \frac{g^{kT} g^k}{d^{k-1T} (g^k - g^{k-1})} \right) d^{k-1}$$

Now, we have got  $d^k$  to be  $-g^k$  plus  $\sum \beta_j d^j$  where  $\beta_j$  is this quantity, and we know that  $x^{j+1}$  it should be  $x^j$  plus  $\alpha^j d^j$  and  $g^{j+1}$  is equal to  $g^j$  plus  $\alpha^j H d^j$ . So, so let us see what happens to the  $x$  quantity in the numerator here, so  $H d^j$ ,  $H d^j$  is nothing but  $H$  into  $x^k$ ,  $x^j$  plus 1 minus  $x^j$  by  $\alpha^j$ . And that is nothing but this quantity, and therefore, what we get is  $d^k$  to be  $-g^k$  plus  $\sum_j$  going from 0 to  $k-1$  into this quantity.

So, we have  $g^{j+1}$  equal to  $g^j$  plus  $\alpha^j$  into  $H d^j$ , and we use this quantity to write this and this  $H d^j$  in the numerator is now replaced by both in the numerator and denominator is replaced by this quantity. So, this  $\alpha^j$  gets canceled and what we get is  $g^{j+1} - g^j$  transpose  $g^k$  divided by  $g^{j+1} - g^j$  transpose  $d^j$ . Now, let us look at this quantity, so you will see that if you do the exact line search, then  $g^k$



transpose  $g_j^T$  goes from 0 to  $k-1$  will be 0. And what we are left with this quantity  $g_k^T d_{k-1}$ .

So, the important quantity, that is left here is this quantity related to only  $j = k-1$ , the rest of the quantities become 0, because  $g_k$  is perpendicular to all  $g_j$  that we have already shown. So,  $g_k$  is perpendicular to all  $g_j$ ,  $j$  going from 0 to  $k-1$ , so this quantity vanishes and because of the line search,  $d_j^T g_j$  going from 0 to  $k-1$  becomes 0. And we are left with only quantity corresponding to  $j = k-1$ .

So, you will see that, as compared to the usual Gram-Schmidt procedure which used all the previous directions and the minus  $g_k$  direction to find  $d_k$ , here we use only minus  $g_k$  and only the previous direction  $d_{k-1}$ . So, because of the H-conjugacy of these vectors, the formula becomes very simple, so we do not have to worry about all the previous directions to find out  $d_k$ . Or it is sufficient for us to determine  $d_k$  using minus  $g_k$  and  $d_{k-1}$ , and this is possible because  $g_k$  is perpendicular to all the previous  $g_j$ 's and we do the exact line search.

(Refer Slide Time: 41:12)

$$d^k = -g^k + \left( \frac{g^{kT} g^k}{d^{k-1T} (g^k - g^{k-1})} \right) d^{k-1}$$

Due to exact line search,  $g^{kT} d^{k-1} = 0$ .

$$d^{k-1} = -g^{k-1} + \beta^{k-2} d^{k-2}$$

$$-d^{k-1T} g^{k-1} = g^{k-1T} g^{k-1} + \beta^{k-2} g^{k-2T} d^{k-2}$$

Therefore,

$$d^k = -g^k + \frac{g^{kT} g^k}{g^{k-1T} g^{k-1}} d^{k-1}, \quad k = 1, \dots, n-1$$

**Fletcher-Reeves method**

So, due to exact line search,  $g_k^T d_{k-1}$  is equal to 0 and, if we expand this  $d_{k-1}^T g_{k-1}$ , then what we get is this quantity becomes 0 and only this quantity remains. And therefore, what we get is formula for finding the conjugate directions and which is called the Fletcher-Reeves formula or Fletcher-Reeves

method. So, the direction  $d^k$  is found using  $-g^k + (g^k)^T g^k / (g^k)^T g^k$ . So it is a combination of  $-g^k$  and  $d^{k-1}$ . Now, here we have assumed that  $d^0$  to be  $-g^0$ , so we can get all the  $n$  conjugate directions by making use of the gradient at the respective points, and then the previous direction.

(Refer Slide Time: 42:15)

**Conjugate Gradient Algorithm (Fletcher-Reeves)**

For Quadratic function,  $\frac{1}{2}x^T Hx + c^T x$ ,  $H$  symmetric positive definite

(1) Initialize  $x^0$ ,  $\epsilon$ ,  $d^0 = -g^0$ , set  $k := 0$ .

(2) while  $\|g^k\| > \epsilon$

(a)  $\alpha^k = -\frac{g^k{}^T d^k}{d^k{}^T H d^k}$

(b)  $x^{k+1} = x^k + \alpha^k d^k$

(c)  $g^{k+1} = Hx^{k+1} + c$

(d)  $\beta^k = \frac{g^{k+1}{}^T g^{k+1}}{g^k{}^T g^k}$

(e)  $d^{k+1} = -g^{k+1} + \beta^k d^k$

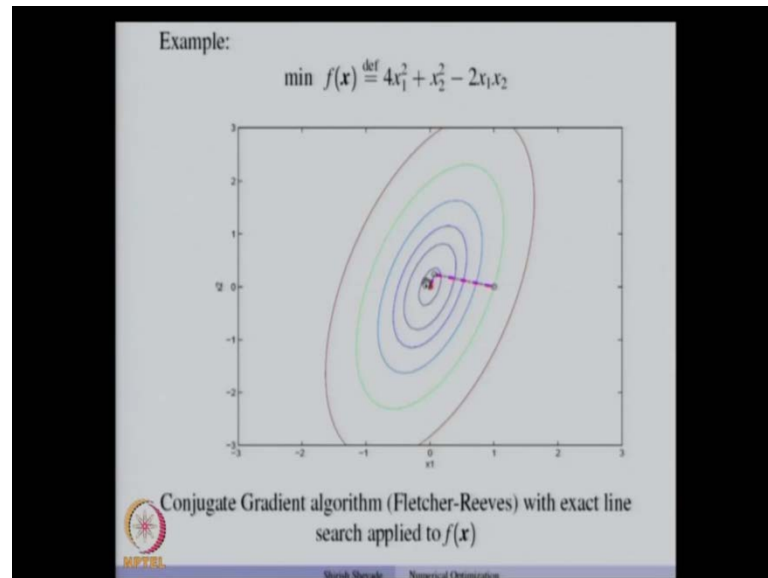
(f)  $k := k + 1$

endwhile

Output :  $x^k$ , global minimum of  $f(x)$ .

So, here is the Conjugate Gradient Algorithm, so we start with  $d^0$  to be  $-g^0$ , initialize  $x^0$  epsilon. And remember that initially we, we have given the algorithm only for the quadratic function  $\frac{1}{2}x^T Hx + c^T x$ , where  $H$  is a symmetric positive definite matrix. So, the first step is to find  $\alpha^k$  we have a closed form expression for this if we do the exact line search then  $x^{k+1}$  is equal to  $x^k + \alpha^k d^k$  and  $g^{k+1}$  is  $Hx^{k+1} + c$ , that is the gradient. We determine  $\beta^k$  and then  $d^{k+1}$  the new direction is determined  $-g^{k+1} + \beta^k d^k$ , and the  $k$  is incremented by 1. And we check whether the the gradient of the function, and the new point is less than or equal to epsilon less than or equal to epsilon we stop. And we get a global minimum of  $f$ , because let this a global minimum  $H$  is a symmetric positive definite matrix for this quadratic function.

(Refer Slide Time: 43:31)



(Refer Slide Time: 43:58)

Now, when we apply this to, the function that we had seen earlier. So, the use of conjugate gradient method to this function results in this steps so we start from here, so this is  $x^0$   $x^1$  and  $x^*$  so within 2 steps we reach the solution. Now, if you start from a different point, initial point even then you will see that we reach the solution in 2 steps. Now, how do we extend this formula to non-quadratic functions, so this initial step is still the same. But then after every any iterations there is a need to restart the method. And the reason for this is that, when we talk about non-quadratic function we have to do this line search, and sometimes that line search also may not be exact.

And more over this directions that we get in the conjugate gradient method, where we use the Fletcher-Reeves formula ,these directions may not be descent. Now, to make sure that in the n retractions, at least 1 direction is negative or at least 1 direction is a descent direction there is a need to restart the conjugate gradient method for non-quadratic functions after every n iterations, and this is the restart. That the new point is again set as  $x_0$  and  $d_0$  is set as the negative of the gradient at that point so which will make sure, that we have the descent direction at least one descent direction in n iterations remember start with  $d_0$  equal to minus  $g_0$ .

So, at least 1 direction in the n iterations will be a descent direction, and the procedure is repeated till norm of  $g_k$  becomes less than equal to epsilon. So, this is the difference between the application of Conjugate Gradient method to Quadratic functions and non-quadratic functions. For quadratic functions, we reach the solution in the at most n iterations and we did not have to worry about the restart there. And finally, what we get here is a stationary point of f of x, Now, if we look at the algorithm we use Fletcher-Reeves method to find beta k. Now, there exist different ways to determine beta k and the different Conjugate Gradient Algorithms vary with respect to this calculation of beta k.

(Refer Slide Time: 46:23)

**Determination**

- Fletcher-Reeves method
 
$$\beta_{FR}^k = \frac{g^k{}^T g^k}{g^{k-1}{}^T g^{k-1}}$$
- Polak-Ribiere method
 
$$\beta_{PR}^k = \frac{g^k{}^T (g^k - g^{k-1})}{g^{k-1}{}^T g^{k-1}}$$
- Hestenes-Steifel method
 
$$\beta_{HS}^k = \frac{g^k{}^T (g^k - g^{k-1})}{(g^k - g^{k-1})^T d^{k-1}}$$

So, so the formula that we have seen so far is in the Fletcher-Reeves formula, which uses beta k to be  $g_k^T g_k$  by  $g_{k-1}^T g_{k-1}$ . This is another

formula for Polak-Ribiere formula, and that uses  $g^k$  transpose into  $g^k$  minus  $g^k$  minus 1 by  $g^k$  minus 1 transpose  $g^k$  minus 1.

Now, the only difference between these 2 formulas is that, there is a extra quantity here. Now, if we use quadratic functions with exact line search, then we know that  $g^k$  is perpendicular to all the previous gradients. So this quantity is 0, so the 2 formulas are same, but for non-quadratic function it was empirically observe that, this formula the Polak-Ribiere formula performs better compare to the Fletcher-Reeves formula. So, for Quadratic functions with exact line search these 2 formulas will give us the same steps, there is another formula got Hestenes-Steifel formula, and that is given here now the motivation for this formula was different, but we can look it from the B F G S method view point for Quasi-Newton directions that we saw earlier.

(Refer Slide Time: 47:50)

$$B_{BFGS}^k = B + \left( 1 + \frac{\gamma^T B \gamma}{\delta^T \gamma} \right) \frac{\delta \delta^T}{\delta^T \gamma} - \left( \frac{\delta \gamma^T B + B \gamma \delta^T}{\delta^T \gamma} \right)$$

Memoryless BFGS iteration

$$B_{BFGS}^k = I + \left( 1 + \frac{\gamma^T \gamma}{\delta^T \gamma} \right) \frac{\delta \delta^T}{\delta^T \gamma} - \left( \frac{\delta \gamma^T + \gamma \delta^T}{\delta^T \gamma} \right)$$

With exact line search,  $\delta^{k-1^T} g^k = \alpha^{k-1} d^{k-1^T} g^k = 0$ . Therefore,

$$d_{BFGS}^k = -B_{BFGS}^k g^k = -g^k + \frac{\delta \gamma^T g^k}{\delta^T \gamma} = -g^k + \underbrace{\frac{g^k (g^k - g^{k-1})}{(g^k - g^{k-1})^T d^{k-1}}}_{\delta_{BFGS}^k} d^{k-1}$$

So, recall that the update formula for B F G S method is given here, and there is a variant of this B F G S method which is call the Memoryless B F G S method. And in the memoryless variant this B is replace by identity matrix, so which means that at every  $x^k$  we do not have to use  $B^k$ . But use in this case this will be  $B^k$  minus 1, so we do not have to use  $B^k$  minus 1, but instead replace  $B^k$  minus 1 by identity matrix. So, if we do that, then what we get is called a memoryless Quasi-Newton direction. And if we do the exact line search, the delta  $k$  minus 1 transpose  $g^k$  equal to 0, so remember that this delta is a function of  $x^k$  and  $x^k$  minus 1 and  $g^k$  is a function of  $g^k$ .

So,  $g^k$  is the gradient and therefore, so if we use this update, the Quasi-Newton update BFGS which uses the Memoryless iteration. So, this quantity will be nothing, but the quantity which is given here, and you will see that this formula this update of  $d^k$  will be same as minus  $g^k$  plus some  $\beta^k$  into  $d^{k-1}$  minus 1. And that  $\beta^k$  is nothing but the formula given by Hestenes-Steifel, so this BFGS method. The DFP method these can be thought of other way to generate conjugate directions.

(Refer Slide Time: 49:40)

For nonquadratic function,  $f(x)$ :

**Conjugate Gradient Algorithm (Fletcher-Reeves)**

- (1) Initialize  $x^0, \epsilon, d^0 = -g^0$ , set  $k := 0$ .
- (2) **while**  $\|g^k\| > \epsilon$ 
  - (a)  $\alpha^k = \arg \min_{\alpha > 0} f(x^k + \alpha d^k)$
  - (b)  $x^{k+1} = x^k + \alpha^k d^k$
  - (c) Compute  $g^{k+1}$
  - (d) **if**  $k < n - 1$ 
    - $\beta^k = \frac{g^{k+1T} g^k}{g^{kT} g^k}$
    - $d^{k+1} = -g^{k+1} + \beta^k d^k$
    - $k := k + 1$
  - else**
    - $x^0 = x^{k+1}$
    - $d^0 = -g^{k+1}$
    - $k := 0$
  - endif**
- endwhile**

**Output :**  $x^* = x^k$ , a stationary point of  $f(x)$

Now, now let us look at the the difference between the DFP and the BFGS method. DFP or the Quasi-Newton method and the Conjugate Gradient method. Now, this is a algorithm for minimizing a general non-quadratic function  $f$  of  $x$  using Conjugate Gradient method. Now, this does not require any second order information, so or it does not use any matrix operations. So, the main difference between the Conjugate Gradient method and Quasi-Newton method, is that there are no matrix operations involved. Secondly the Conjugate Gradient for the Quasi-Newton method like the DFP method or the Quasi-Newton, the LBFGS method both of them in addition to generating conjugate generations.

They also update the matrix  $B^k$ , and whenever we are near the solution that  $B^k$  will be a good approximation of the Hessian Inverse, and therefore, the directions that we obtained using DFP or BFGS methods near the solution will be a close to the Newton directions, and then the convergence will be faster. So, in addition to maintaining

conjugate directions strategy, they also find the directions which are more close to the netwon directions near the solution. And therefore, they are fast and that does not happen in the Conjugate Gradient method.

Conjugated Gradient method, since it does not work with any matrices there is no guaranty for general non-quadratic function with in exact line search, that the direction that we get is a descent direction. And therefore, there is a need to restart, while that is not the case, this restarting is not necessary for the Quasi-Newton methods like D F P or B F G S. So that is a important point and more over that Conjugate Direction methods or Conjugate Gradient methods they are very sensitive to the line search.

While Quasi-Newton methods are more robust, now on the other hand the Quasi-Newton methods require, the multiplication of a matrix by a vector. And that require order  $n$  square computations as well as order  $n$  square storage, because the matrix the enter matrix needs to be store. If we do not use the limited memory B F G S kind of updates. Now, that is not the case here, if you look at the storage required here will, will require order  $n$  storage we need to store some  $n$ -dimensional vectors in the Conjugate Gradient method. So, the storage wise Conjugate Gradient method has the advantage as compare to the Quasi-Newton method.

But, both these methods Quasi-Newton method as well as the Conjugate Gradient method, they have better computational complexity as compare to Newton method. Because Newton method requires inversion of a matrix to get a Newton direction and that is computationally expensive operation. Well, here neither conjugat gradient method nor quasi-newton method require any inversion of a matrix, further Conjugated Gradient methods do not even require to store any matrix.

So, storage wise these methods clearly have a advantage, so if one wants to get a robust method one has to for Quasi-Newton methods. Because they were found to be more robust compare to conjugate gradient methods and they are less sensitive to the line search, and of the 2 methods D F P and B F G S, B F G S method of finding Quasi-Newton direction, what is found to be empirically superior to D F P method. So, based on the need, one has to decide the algorithm that can be applied to minimize given function. So, for quadratic functions methods like conjugate gradient method are quite

good, for general non-quadratic case one depending upon the requirement, one can go for the Quasi-Newton kind of update especially the BFGS kind of update.

So, this completes our discussion on unconstrained optimization. So, we saw Hestenes-Steifel method, then the Newton method, then the variants of Newton method the Quasi-Newton method and conjugate gradient method, Now, all these methods are useful in different ways to solve given optimization or a given minimization problem. And in the next class, we will study the methods for constrained optimization, and find the properties related to those methods.

Thank you.