

**Foundation of Optimization**  
**Prof. Dr. Joydeep Dutta**  
**Department of Mathematics and Statistics**  
**Indian Institute of Technology, Kanpur**

**Lecture - 3**

(Refer Slide Time: 00:25)

$\nabla f$  each  $\frac{\partial f}{\partial x_i}$ ,  $i=1, 2, \dots, n$  is continuous as a function of  $x_1, x_2, \dots, x_n$  and all second order mixed partial derivatives are cont.

$$A(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix} = \nabla^2 f(x)$$

$\downarrow$   
 Hessian Matrix  
 $\downarrow$   
 is symmetric.

- Unconstrained Optimization

$$\min f(x), \quad x \in \mathbb{R}^n$$

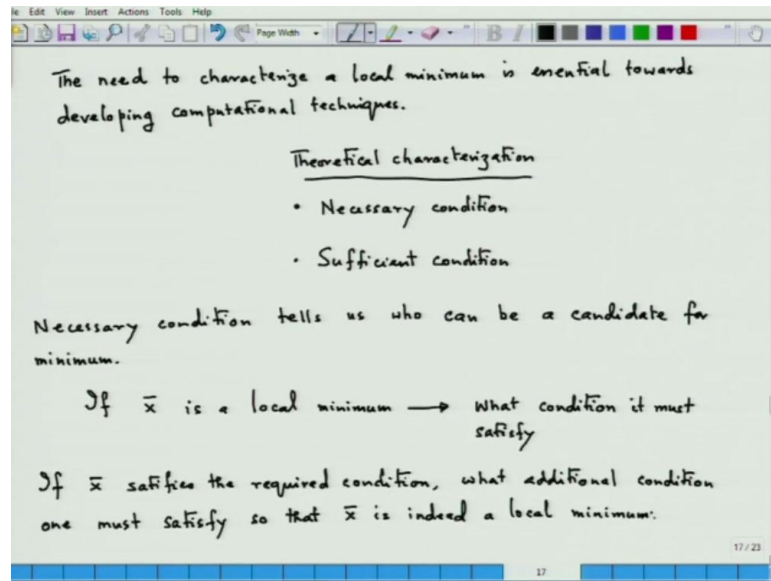
theoretical aspect

Computational aspect
- Theoretical Issues

The main is to characterize a local minimum.

So, today we are going to concentrate on unconstrained optimization. That is here we are going to look at minimizing effects where  $x$  is in  $\mathbb{R}^n$  whole  $\mathbb{R}^n$ . Now, what are the two aspects of this problem? So, any optimization problem has two aspects them two aspects; first is the theoretical aspect and the second aspect is the computational aspect. Now, in the theoretical aspect, what are we suppose to do in the computational aspect, what are we suppose to do these are the two main themes of discussion we will carry out today. So, let us see what is the, what is the theoretical issues here. The main issue is to characterize a local minimum. Why you need to characterize a local minimum?

(Refer Slide Time: 02:05)



It is important to note the need to characterize local minima, minimum is essential to develop is essential towards developing computational techniques. So, in the theoretical characterization there are 2 things that we really have to bother about these. So, there are 2 aspects of it namely the necessary condition and other is the sufficient condition. Now, what do these two conditions do? A large amount of optimization literature is bothered about really geared towards designing these sorts of conditions for various types of optimization problems. Now, why you need a necessary condition, why you need a sufficient condition? Necessary condition tells us how to compute a candidate point for a minimum.

So, necessary condition tells us who can be a candidate for minimum, who can be a candidate for minimum. Now, if he tells you who can be a candidate for minimum. So, what it should do? So, necessary condition will do this. If  $\bar{x}$  is a local minimum, what condition it must satisfy? Means a point, which does not satisfy that condition cannot be a local minimum. So, that is very important. So, in order that your point that you say that possibly be a local minimum has to satisfy this condition. First you have to find the point which satisfies this condition. Then you check whether it is a local minimum or global minimum or whatever.

If a point does not satisfy this condition, then it is ruled out it is it can never be a local minimum. So, that is why necessary conditions are very important or rather central to

the study of optimization. Now, if  $\bar{x}$  satisfies the condition the required condition, what additional assumption would guarantee that it is really a local minimum that is where sufficient condition comes. It tells you if these conditions are satisfied along with the fact that  $\bar{x}$  has come from the computation that you do with the necessary condition then that  $\bar{x}$  is a minimum. What condition it must satisfy? So that,  $\bar{x}$  is indeed a local minimum. Now once you know this. So, our first job is let us layout our job properly.

(Refer Slide Time: 07:38)

Our first job is  $\rightarrow$  To find the necessary condition

**Theorem 1:** If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable and  $\bar{x}$  is a local minimum then  $\nabla f(\bar{x}) = 0$

**Proof:** Since  $\bar{x}$  is a local minimum,  $\exists \delta > 0$ , s.t.  
 $\forall x \in B_\delta(\bar{x}), f(x) \geq f(\bar{x})$ .

For any  $w \in \mathbb{R}^n$ ,  $\exists \lambda_0 > 0$ , s.t.  $\forall \lambda \in (0, \lambda_0)$   
 $\bar{x} + \lambda w \in B_\delta(\bar{x})$   
 $\Rightarrow f(\bar{x} + \lambda w) \geq f(\bar{x})$   
 $\Rightarrow f(\bar{x}) + \lambda \langle \nabla f(\bar{x}), w \rangle + o(\|\lambda w\|) \geq f(\bar{x})$   
 $\Rightarrow \lambda \langle \nabla f(\bar{x}), w \rangle + o(\|\lambda w\|) \geq 0$

The slide includes a diagram of a 2D coordinate system with a point  $\bar{x}$  and a small circle representing the neighborhood  $B_\delta(\bar{x})$ . A vector  $w$  is shown originating from  $\bar{x}$ , and a point  $\bar{x} + \lambda w$  is marked on the line segment from  $\bar{x}$  in the direction of  $w$ .

Our first job to find the necessary condition, this we state as follows, which we write as theorem one of our study. It says that if, if from  $\mathbb{R}^n$  to  $\mathbb{R}$  is differentiable and  $\bar{x}$  is a local minimum, then my next question is, then what would happen then gradient of  $f$  at  $\bar{x}$  must vanish. Like any other part of mathematics, optimization in a mathematical subjects it demands proof of what statement you are making. You have heard this when  $n$  is equal to 1 that is when  $f$  is from  $\mathbb{R}$  to  $\mathbb{R}$  you have studied this in school if you are trying if you are given a ordinary function say  $x^2$  plus two  $x$  plus theta when I ask you to find the minimum. You will immediately go and find the derivative you got it to 0 find a point and try to do something. So, here so what we have to do to find the local minima is to first find a point which satisfies this.

So, let us see whether this condition, how this condition gets satisfied? So, what we are writing here is formally the first proof of this talk. How does we how does we do we

proceed to do the proof? No one has to understand, we are talking about local minima so around  $\bar{x}$  we are looking at a neighborhood around  $\bar{x}$ , we are not looking around, around the looking at the nature of the function very far from  $\bar{x}$ . So, then we can employ Taylor's theorem or which is another way of talking about differentiability. So, we will use this idea of differentiability in a very, very sensible way. So, what does it mean by a local minimum. Since  $\bar{x}$  is a local minimum. So, I will write it down completely. So, you can follow the proof step by step.

Now, once you know that  $\bar{x}$  is a local minimum what you have there exists  $\delta$  greater than 0, such that for all  $x$  in the  $\delta$  neighborhood of  $\bar{x}$ , the ball around  $\bar{x}$   $f$  of  $x$  must be bigger than  $f$  of  $\bar{x}$ . Now, I will try to draw the picture of this scenario in order to which can, which will make proof much clear. So, here is my  $\bar{x}$  and here is the ball here is my vector  $\bar{x}$ . So, this is my  $B_\delta$  of  $\bar{x}$ . What I can do is let me take any vector  $w$  in any direction. Then I can choose I can choose to move from  $\bar{x}$  along the direction  $w$ , say this is my  $w$  vector and so this was my  $\bar{x}$  vector.

So, I can choose to move from  $w$  from  $\bar{x}$  long a direction parallel to  $w$ . So, let this point where I come and stop the  $\bar{x} + \lambda w$ . Now, you see if I increase  $\lambda$  there will be a threshold beyond which it will come out of the ball and then I cannot do anything. So, if I am within the ball I can relate the function values of this and with that of  $\bar{x}$  that this function value at this point must be bigger than function value at  $\bar{x}$ . So, for any  $w$  element of  $\mathbb{R}^n$  there exists  $\lambda_0$  strictly greater than 0, such that for all  $\lambda$  between 0 and  $\lambda_0$  you see even if I  $\lambda_0$  is the point where  $\bar{x} + \lambda w$  this vector comes and touches the boundary, because once you are in the boundary you are out of the neighborhood.

So, that is why you have a opened here. For all of these, so that is for all  $\lambda$  of these  $\bar{x} + \lambda w$  remains in the ball this naturally implies  $f$  of  $\bar{x} + \lambda w$  is bigger than  $f$  of  $\bar{x}$ , but knowing the definition of the derivative. Because the function is differentiable at  $\bar{x}$  I can write I would just like to recall where I had written the definition of the derivative, this is the definition of the derivative. So, I can write here  $f$  of  $\bar{x} + \lambda w$  is  $f$  of  $\bar{x}$  plus  $\lambda$  times  $\text{grad } f$  at  $\bar{x}$  plus  $o(\lambda)$ . And this is bigger than  $f$  this also means the following this means that  $\lambda$  times  $\text{grad } f$  at  $\bar{x}$ , this is what you get, because I have cancelled  $f$  of  $\bar{x}$  from both sides.

(Refer Slide Time: 15:08)

$$o(\|\lambda w\|) = o(\lambda \|w\|) \approx o(\lambda)$$

$$\text{i.e. } \lim_{\lambda \rightarrow 0^+} o\left(\frac{\|\lambda w\|}{\lambda}\right) = 0$$
 Dividing both sides by  $\lambda$  we have
 
$$\langle \nabla f(\bar{x}), w \rangle + \frac{o(\|\lambda w\|)}{\lambda} \geq 0$$
 As  $\lambda \rightarrow 0^+$  (i.e.  $\lambda > 0, \lambda \rightarrow 0$ ) we have
 
$$\langle \nabla f(\bar{x}), w \rangle \geq 0.$$
 Since  $w$  was chosen arbitrarily we have
 
$$\langle \nabla f(\bar{x}), w \rangle \geq 0, \forall w \in \mathbb{R}^n.$$
 Put  $w = -\nabla f(\bar{x}) \Rightarrow \langle \nabla f(\bar{x}), -\nabla f(\bar{x}) \rangle \geq 0 \Rightarrow -\|\nabla f(\bar{x})\|^2 \geq 0$   

$$\Rightarrow \|\nabla f(\bar{x})\|^2 \leq 0$$

Now, let us look at this term  $o(\lambda w)$ . So, whatever be the powers of  $o(\lambda w)$ ,  $\lambda w$  would have the same powers. So, if you have, because you can bring out  $\lambda$  outside here. So, you can write this as  $o(\lambda w)$ . So, basically here if  $w$  is fixed if I vary the  $\lambda$  I can keep on changing the things. So, basically this is nothing but a  $o(\lambda w)$  quantity that is  $o(\lambda w)$  by  $\lambda$  as  $\lambda$  tends to 0 plus, I have taken  $\lambda$  as positive that should be equal to 0. So, here what I have observed now I have got if I divide by  $\lambda$  because  $\lambda$  here is positive, because  $\lambda$  is positive I can divide both sides by  $\lambda$ .

So, let me write down more clearly, we have now as  $\lambda$  tends to 0 plus that is  $\lambda > 0$ . And  $\lambda$  is going to 0 so as  $\lambda$  tends to 0 plus we have because this will go to 0. Now,  $w$  was chosen arbitrarily,  $w$  was not fixed thing fix vector. Since  $w$  was chosen arbitrarily, we have for all  $w \in \mathbb{R}^n$  for all  $w \in \mathbb{R}^n$ . So, once this is known, now put  $w$  is equal to the negative of  $\nabla f$  of  $\bar{x}$ . So, this would imply from this condition that  $\nabla f$  of  $\bar{x}$  minus  $\nabla f$  of  $\bar{x}$  is greater than equal to 0 which would imply minus of because I take the minus outside and normal inner product of this a vector with itself is nothing but square of its own is greater than equal to 0 which would imply at the known, but it is this is the length. So, the length of a vector from the length of a vector cannot be less than equal to 0.

(Refer Slide Time: 15:08)

$\Rightarrow \|\nabla f(\bar{x})\|^2 = 0$   
 $\Rightarrow \|\nabla f(\bar{x})\| = 0 \downarrow (\|x\|=0 \Leftrightarrow x=0)$   
 $\Rightarrow \nabla f(\bar{x}) = 0 \rightarrow \text{critical point}$

•  $f(x) = x^3$        $f'(\bar{x}) = 0$   
 $\Rightarrow 3\bar{x}^2 = 0$   
 $\Rightarrow \bar{x}^2 = 0$   
 $\Rightarrow \bar{x} = 0$   
The only critical point is  $\bar{x} = 0$   
But  $\bar{x} = 0$  is not a local minimum.

So, this would imply that and you know from the property of the norm this will imply that the gradient of  $f$  at  $\bar{x}$  is 0. So, here to here we have applied the formula norm of  $x$  is equal to 0 if and only if  $x$  is equal to 0, this is the property of the norm and here at this point from here to here we have employed the fact that  $x \cdot x$ , which have been a product of this is nothing but square of the norm. So, we have this condition. Now we need to test it up for example, you take first the function  $f(x) = x^3$ . We want to show that it is truly a necessary condition at a point  $\bar{x}$  which satisfies this need not be a point of minimum. Any point which satisfies, this is called a critical point.

Now, you have  $f(x) = x^3$ , let us look at the picture of  $f(x) = x^3$  a graph. Now,  $f'(\bar{x}) = 0$  would imply  $3\bar{x}^2 = 0$  and that would imply  $\bar{x}^2 = 0$  and that would imply  $\bar{x} = 0$ . So, the only critical point in this case is  $\bar{x} = 0$ . But  $\bar{x} = 0$  is not a point or not a local  $\bar{x} = 0$  is not a minimizer  $\bar{x} = 0$  is not a local minimum. So, this is very important that this is just a necessary condition and not a sufficient condition.

(Refer Slide Time: 21:48)

$f(x, y) = x^2 + y^2, (x, y) \in \mathbb{R} \times \mathbb{R} = \mathbb{R}^2$

$\nabla f(\bar{x}, \bar{y}) = 0$

$\begin{pmatrix} 2\bar{x} \\ 2\bar{y} \end{pmatrix} = 0 \Rightarrow \begin{matrix} \bar{x} = 0 \\ \bar{y} = 0 \end{matrix} \Rightarrow \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  is a critical point

But  $\begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  is also a point of global minimum of  $f$

$x^2 + y^2 \geq 0 \quad \forall (x, y) \in \mathbb{R}^2$

$x^2 + y^2 \geq \bar{x}^2 + \bar{y}^2, \quad \forall (x, y) \in \mathbb{R}^2$

$\Rightarrow \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix}$  is a global minimizer.

If I have a critical point, how do I know that it is a minimum?

- $f$  is twice continuously differentiable
- $\nabla^2 f(x)$  is positive-definite at  $x = \bar{x}$ , where  $\bar{x}$  is a critical pt.

Let us look at another example where this is possibly alright. So,  $x, y$  here is in  $\mathbb{R}$  cross  $\mathbb{R}$  each individually from the real line which is nothing but  $\mathbb{R}^2$ . Now, if I want to find the gradient of  $f(x, y)$ , I want to find this means I would have two of  $\bar{x}$  and  $\bar{y}$ . This vector would be equal to the 0 vector. This would imply  $\bar{x}$  is equal to 0 vector and  $\bar{y}$  is equal to 0 vector. So, it is  $\bar{x}, \bar{y}$  is equal to 0 0 is the critical point. But  $\bar{x}, \bar{y}$  equal to 0 0 is also a point of global minima of the function. How do I know this fact? Because you observe that  $x^2 + y^2$  is always greater than equal to 0. So,  $x^2 + y^2$  because this is all non negative is bigger than  $\bar{x}^2 + \bar{y}^2$  greater than 0. this is true for all  $x, y$  in  $\mathbb{R}^2$  where this is 0 this is 0.

So, this would imply that  $\bar{x}, \bar{y}$  is truly a global minimizer. So, here you see the two different notions of differentiable function in 1 case, I am having the critical point to be a global minimizer this forms an important class of functions called convex which will focus on little later for some part of for a little part of the course. But which is none the less important and there is a complete course on optimization problems with convex functions which are already been delivered and you will see it quite soon. But here we are dealing with all types of functions which are differentiable could be convex could be non convex whatever. So, here is so there are two examples; one example here for example, this function is really not convex.

Now if, now the question is what is the point here? Actually here you see the curve it changes the nature. So, if you have a function from  $\mathbb{R}$  to  $\mathbb{R}$  then a point is either a local or global minima a critical point or it is a like this where the curve is changing shape. In general, I could find the critical point which is neither a local minimum, local maxima nor nothing it is just a critical point for and it is not also a saddle point. Now, I am not going to get into the saddle point issue right now, because they might just confuse you. So, what I am going to get going at this movement is that if I have a critical point, how do I know that it is a minimum. To answer this question we have to additionally assume certain conditions on function  $f$ .

So, what are those conditions? so trying making an attempt to answer it just by having the knowledge that  $f$  is differentiable or even continuously differentiable does not help me at all it does not tell me anything, because that is where we get stuck. But now if I put some more additional conditions on the function can I devise a condition which if satisfied by a critical point with guarantee that such a point is a local minimum. First condition is that if  $f$  is twice continuously differentiable.

And second condition is the Hessian matrix is positive definite at  $x$  equal to  $\bar{x}$ , where  $\bar{x}$  is a critical point. Now, what do you mean by the term positive definite and what do you mean by the term positive semi definite these are terms which comes from matrix theory, but which are very important in optimization. And so I would just like to take a minute to explain what these terms are. So, our matrices by because we have taken the function do be twice continuously differentiable the hessian matrix is a symmetric matrix.



(Refer Slide Time: 28:23)

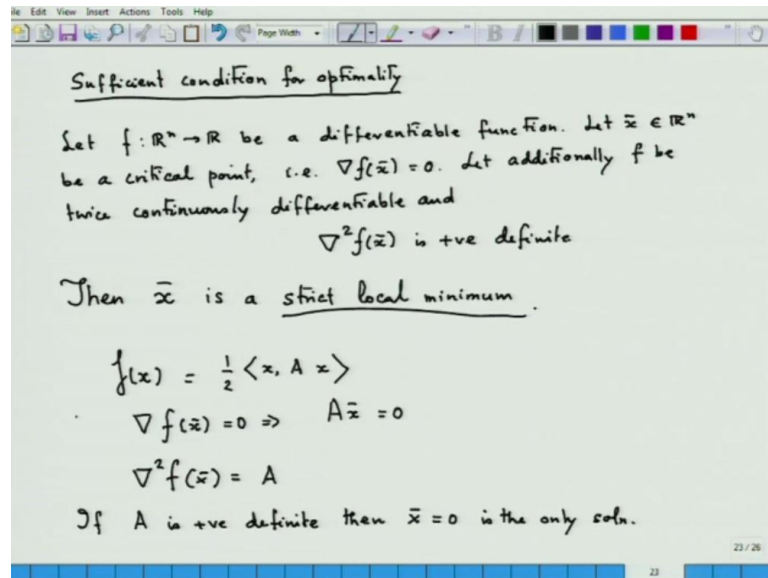
The image shows a digital whiteboard with handwritten mathematical definitions and properties. The text is as follows:

- A to be a  $n \times n$  symmetric matrix
- Now A is said to be positive semi-definite if
$$\langle x, Ax \rangle \geq 0, \forall x \in \mathbb{R}^n$$
- A is said to be positive definite if
$$\langle x, Ax \rangle > 0, \forall 0 \neq x \in \mathbb{R}^n$$
- $\exists$  if A is positive definite we can show that
$$\langle x, Ax \rangle \geq \lambda_{\min}(A) \|x\|^2 \quad \forall x \in \mathbb{R}^n, x \neq 0$$
where  $\lambda_{\min}(A)$  is the minimum of eigen-value of A.
- $\exists$  if A is +ve semi-definite all eigenvalues are non-negative
- $\exists$  if A is +ve definite then all eigenvalues are positive.

The whiteboard interface includes a menu bar (File, Edit, View, Insert, Actions, Tools, Help) and a toolbar with various drawing tools. The page number '22 / 23' is visible in the bottom right corner.

So, let us consider a to be A n cross n symmetric matrix. Now, now A is said to be positive definite positive, semi definite positive semi definite. If the inner product of x with A x is always bigger then equal to 0 for every x in R n. So, this is condition is always taken and now A is said to be positive definite if x A x is strictly bigger that 0 for all non zero x in R n. In fact if A is we can show that, this is where lambda min of A is the minimum Eigen value. It is important, however to note that if A is positive semi definite all Eigen values are non negative. If A these are also 2 bullet points, if A is positive definite, then all Eigen values are positive. Of course, these Eigen values are all real, because A is a symmetric matrix. Now, this is what you have now we go to the main question the second order result.

(Refer Slide Time: 32:25)



So, we are now going to state the sufficient condition for optimality and you would be amazed how strong a result we get. Now, here is the result. You might ask me for a proof of this, but the proof keeping in view that this optimization course is really for a very broad audience from very different disciplines in engineering in the sciences, so in mathematics of course. But would not like to bog you down with the proof of this which is quite technical and needs a slightly more deeper understanding on mathematics I would rather give you the result. And you would see for yourself how to apply it, then we would go more into the algorithmic aspect computational aspect which would be really of useful to you, and and you would be really able to use it in several things specially those who are in engineering sciences. So, let us just state this result.

So, let  $f$  from  $\mathbb{R}^n$  to  $\mathbb{R}$  be a differentiable function. Let  $\bar{x}$  elemental of  $\mathbb{R}^n$  be a critical point that is gradient of  $f$   $\bar{x}$  is equal to 0. Let additionally  $f$  be twice continuously differentiable which you had added here also I just did not do it. And the hessian matrix is positive definite, then  $\bar{x}$  is a strict local minimum proof of this would be added to the f A Q of this subject. But here we just do not do the proof. So, you can see it later on in the website of this course main conclusion here is that it is not just a local minimum it is a strict local minimum. So, it forbids global maximum to take this position so global maximum can never be a strict local minimum. So, we are really getting a local minimum and not been fooled by flat type of functions where you have you can say a local minimum can be also global maximum.

So, those sorts of anomalies won't come because it will give you this. So, this is a very important result and this result has to be really appreciated. Now, if you look at for example, I take an example of a quadratic function. Now, then grad of  $f$   $\bar{x}$  equal to 0 would imply  $A$  of  $\bar{x}$  would be equal to 0. But if  $A$  is positive definite, then you will observe that the hessian matrix is a here for any  $x$  actually the hessian matrix is  $A$ . So, if  $A$  is positive definite here the answer would be of course, 0 as we will soon see if  $A$  is positive definite. Then  $\bar{x}$  equal to 0 is the only solution it is a strict local minimum there are many other cases where a your second order condition you have learnt in school actually gives you the strict local minima. So, this is very simply case actually which is to illustrate you that how to how do you find the hessian.

(Refer Slide Time: 38:13)

$f: \mathbb{R}^n \rightarrow \mathbb{R}$   
 $n = 2 \quad f(x, y) = (x-1)^3 + y^2$   
 $\nabla f(\bar{x}, \bar{y}) = 0 \Rightarrow \begin{cases} 3(\bar{x}-1)^2 = 0 \\ 2\bar{y} = 0 \end{cases}$   
 The only critical point is  $\bar{x} = 1, \bar{y} = 0$   
 $\nabla^2 f(\bar{x}, \bar{y}) = \begin{pmatrix} 3(\bar{x}-1)^2 & 0 \\ 0 & 2\bar{y} \end{pmatrix} \nabla^2 f(\bar{x}, \bar{y})$   
 $\nabla^2 f(1, 0) = \begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix}$  is +ve semi definite  
 $\det \begin{pmatrix} 0-\lambda & 0 \\ 0 & 2-\lambda \end{pmatrix} = 0 \Rightarrow \lambda = 0, \lambda = 2$   
 $(\bar{x}, \bar{y}) = (1, 0)$  is not a local min / Not a local max is just a critical point

For example, if you have a say a more slightly strong slightly different looking scenario, put  $n$  equal to 2. Now, try to figure out say let me put  $f$  of  $x, y$ . So, let me put out try to figure out this. So, this would imply  $3x$  power minus 1 whole square is equal to 0 and  $2y$  bar is equal to 0. So, the only critical point is, point is  $\bar{x}$  is 1 and  $\bar{y}$  is 0. Now if I want to compute the hessian of this. So, I compute the hessian of this. So, you I take the gradient of this right. So, I take the gradient of this in  $\bar{x}$ , then in  $\bar{y}$  and write it as a row vector. Then this will become  $6\bar{x}$  minus 1 and the next term would be 0, because there is no  $y$  term. And here first I will differentiate with respect to  $x$  term which is 0 and then I have 2. Now, grad square  $f$  at  $1, 0$  would give me 0 0 0 2.

So, let me try to find the Eigen values of this matrix, which Eigen values of this matrix means you have to find, so that determinant of this is equal to 0 which would imply  $\lambda$  is 0 and  $\lambda$  is 2. So, this matrix is positive definite. So, this is positive semi definite, but not positive definite. In fact, in fact  $\bar{x}$ ,  $\bar{y}$  equal to 1 0 is not a local min. In fact, it is not a local max is just a critical point, see it is positive semi definite does not give you, what you want? What you require it is positive definiteness, which will give you what you require. So, this example illustrates that.

So, with this I end today's discussion. And tomorrow we would talk about decent direction as to how to really compute a point. But you know in optimization we will learn a very important lesson tomorrow that real optimization problems cannot be solved exactly. So, we will get in to that issue tomorrow and in detail try to develop computational methods to solve or minimize a differentiable function over  $\mathbb{R}^n$ .

Thank you very much.