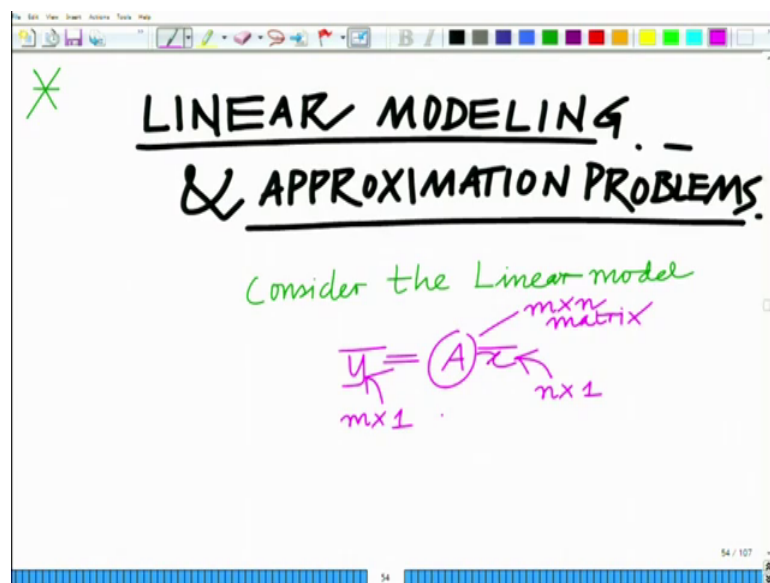


**Applied Optimization for Wireless, Machine Learning, Big Data**  
**Prof. Aditya K. Jagannatham**  
**Department of Electrical Engineering**  
**Indian Institute of Technology, Kanpur**

**Lecture – 41**  
**Linear modeling and Approximation Problems: Least Squares**

Hello, welcome to another module in this massive open online course. So, in this module, let us look at another class of problems or another class of optimization problems, specifically pertaining to Linear modeling and Approximation, which arise very frequently in various applications of course both engineering, science and so on so. This is a very important class of problems.

(Refer Slide Time: 00:38)



This is termed as linear wherever you have models, and typically more most of them or frequently they are modeled as linear model. So, we have linear modeling and approximation problems, linear modeling and approximation problems. And what you can see in this is that if you consider now consider the linear model, general linear model can be described as follows. Consider the linear model  $\bar{y}$  equals  $A \bar{x}$ , where we have  $A$ , this is in general an  $m$  cross  $n$  matrix, and which implies of course that if  $\bar{x}$  is a vector,  $\bar{x}$  is a  $n$  cross  $1$  vector. And  $\bar{y}$  is a naturally an  $m$  cross  $1$  vector.

(Refer Slide Time: 02:18)

The diagram shows the linear system  $\bar{y} = A\bar{x}$  written in purple ink. The vector  $\bar{y}$  is annotated with  $m \times 1$ . The matrix  $A$  is circled and annotated with  $m \times n$  matrix. The vector  $\bar{x}$  is annotated with  $n \times 1$ . A blue circle encloses the entire equation. Below the equation, the text  $\bar{x}$  is unknown To be Determined. is written in purple. At the bottom, the text  $m = \# \text{ Equations. } n = \# \text{ unknowns}$  is written in blue, with arrows pointing from  $m$  to the number of equations and from  $n$  to the number of unknowns. The diagram is presented on a whiteboard with a toolbar at the top and a status bar at the bottom showing '54 / 107'.

Now, what we assume or typically in this model, this vector  $\bar{x}$  is unknown, which has to be determined. So,  $\bar{x}$  is unknown and it has to be determined so,  $\bar{x}$  is to be determined. And now frequently what you have also in such a system, now of course let us start by considering a simple example. Let us say  $A$  is square matrix that is  $m$  is equal to  $n$  ok. If  $m$  is a number of rows,  $n$  is a number of columns of  $A$ . Let us now if you look at this system, you can see what is  $m$  in this system,  $m$  is basically the number of equations, it is the dimension of  $y$ . So,  $m$  equals number of equations. And  $n$  equals number of unknowns ok.

(Refer Slide Time: 03:31)

$m = \# \text{ Equations}$     $n = \# \text{ unknowns}$

if  $m = n$   
and  $A$  is invertible  
 $\Rightarrow \hat{x} = A^{-1} \bar{y}$

$\bar{y} = A \bar{x}$   
 $m > n$   
 $\Rightarrow \# \text{ Equations} > \# \text{ unknowns}$

$m \times n$   
matrix

Now, what happens in this is that let us assume a simple scenario to begin with if  $m$  equals  $n$ , and this everyone would know  $m$  equals  $n$ . And  $A$  is invertible, it is if equals  $m$  and  $A$  is invertible, remember this is as to be given. Now, this implies I can find  $\bar{y}$  is equals to  $\bar{x}$ , I can determine  $\hat{x}$  or the estimate of  $x$  equals  $A$  inverse  $\bar{y}$  ok. This I think is a typical solution for the linear system, which most people would know almost students would know.

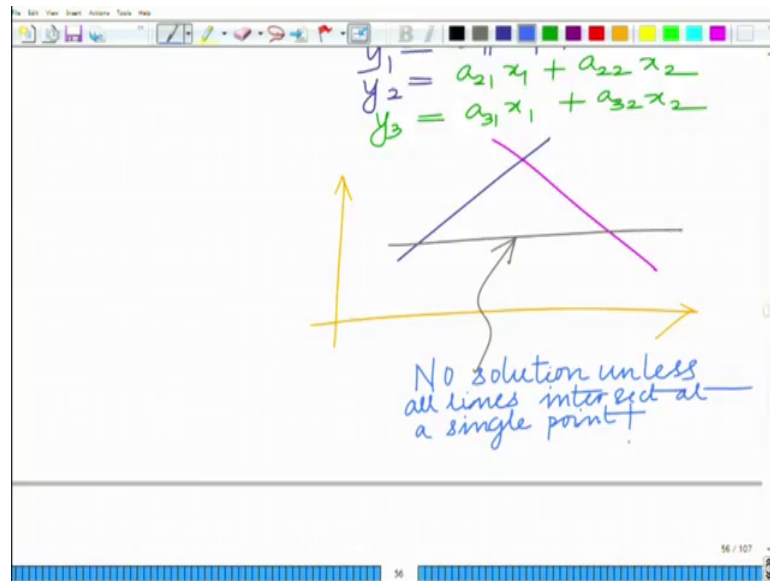
Now, however frequently what we have is that you have this  $\bar{y}$  is well, it is not exactly equal to  $A \bar{x}$  its. Actually,  $A \bar{x}$  plus  $n$  when  $n$  is the noise, I write it as  $\bar{y} = A \bar{x} + n$  this is the linear model,  $\bar{y}$  equals  $A \bar{x}$ . And  $A$  is your  $m$  cross  $n$  matrix, and  $m$  is greater than  $n$ ,  $m$  is greater than  $n$  ok. So, what that means is this implies that number of equations is greater than number of unknowns.

(Refer Slide Time: 04:57)

The image shows a whiteboard with handwritten text. At the top, a green arrow points to the words "Overdetermined System". Below this, it says "Typically solution does NOT exist". A horizontal line separates this from the equations below. The equations are:  $m = 3$ ,  $n = 2$ ,  $A = 3 \times 2$ ,  $y_1 = a_{11}x_1 + a_{12}x_2$ ,  $y_2 = a_{21}x_1 + a_{22}x_2$ , and  $y_3 = a_{31}x_1 + a_{32}x_2$ . The whiteboard has a toolbar at the top and a footer with "56 / 107".

And which also implies that this is basically an over determined system, because number of equations is greater than this is basically an over determined system. Now, of course over determined system frequently, this has no frequently no solution exists, unless the vector  $\bar{y}$  belongs the column space of  $A$ . So, frequently typically because of the noise in this system, typically does not does not exist. For instance, you can take a simple example. Let us take  $m$  equal to 3, what we have is we have  $y_1$  equals a  $1 \times 1$ . So,  $A$  equals 3 cross 2 matrix. So, we have  $y_1$  equals a  $1 \times 1$  plus a  $1 \times 2$ ;  $y_2$  equals a  $2 \times 1$  plus a  $2 \times 2$ ;  $y_3$  equals a  $3 \times 1$  plus a  $3 \times 2$ .

(Refer Slide Time: 06:31)



Now, you can see these represent three lines. So, what happens when you plot these things, you have three lines. So, you basically have three lines each representing a line. Now, the solution now unless all the lines intersect, there is no solution unless all the lines intersect at a common point all right. So, you see this point. So, you have three lines that is three equations two unknowns. So, there is no solution unless all the lines intersect at a point, there is no solution unless all the lines intersect at a single point all right.

So, you have three equations basically, two unknowns. And frequently if you take three lines at random, then they will naturally not intersect. It is highly unlikely that they all intersect at a single point, which means there will be no solution. So, in such a scenario you will have to find an approximate solution or what I mean by that is some solution that best fits the model or best explains the observed vector  $y$ . This is also known as the maximum likelihood vector  $x$ , so it is interesting.

(Refer Slide Time: 08:00)

The image shows a whiteboard with handwritten mathematical equations and annotations. At the top, it says "Typically.  $\bar{y} = A\bar{x} + \bar{n}$ ". A pink arrow points from  $\bar{n}$  to the text "Model Noise". Below this, it shows  $\bar{y} \neq A\bar{x}$ . An arrow points to the difference  $\bar{y} - A\bar{x}$ , which is enclosed in a circle and labeled  $\bar{e}$ . A green arrow points from  $\bar{e}$  to the text "Approximation Error". Another green arrow points from  $\bar{e}$  to the text "Error Vector". The whiteboard interface includes a toolbar at the top and a status bar at the bottom showing "57 / 107".

So, typically and this is because this happens because of the noise in the model.

Typically, what happens is you have  $\bar{y}$  equals  $A\bar{x}$  plus  $\bar{n}$  ok. So, this is what is this is your model noise. And what happens because of this model noise this means, your  $\bar{y}$  not equals is not equal to  $A\bar{x}$ , which means if you form  $\bar{y}$  minus  $A\bar{x}$ , there for any  $\bar{x}$  there is always, because there is no solution for any  $\bar{x}$ , there is always an error. This is termed this is basically your model error or your approximation error vector. So, you can also call this as a error vector, which is basically the approximation error. This is basically the approximation error.

(Refer Slide Time: 09:11)

The image shows a whiteboard with handwritten mathematical derivations. At the top, it says "To find best  $\bar{x}$  minimize approximation error" and "min.  $\|\bar{e}\|$ ". Below this, it shows the equation  $= \min \|\bar{y} - A\bar{x}\|$  and then  $\equiv \min \|\bar{y} - A\bar{x}\|_2^2$ . A green arrow points from the second equation to the text "Least Squares Solution" written in green. A purple arrow points from the same equation to the text "Least Squares Problem" written in purple. The whiteboard interface includes a toolbar at the top and a status bar at the bottom showing "58 / 107".

And we have to find  $\bar{x}$ , you have to find best  $\bar{x}$  that is which best explains  $\bar{y}$ , what do you do is minimize the approximation error. Now, error is of course a vector. What does it mean to minimize the error, we will simply minimize the norm of this error vector. So, we will minimize the norm of the error vector, which is basically equal to minimizing norm of  $\bar{y}$  minus  $A\bar{x}$ , which is basically equivalent to minimizing norm square of  $\bar{y}$  minus  $A\bar{x}$  square ok.

And this problem where you are and this is the two norm square, the  $l_2$  norm. And when you are minimizing the norm square, so this is the vector  $\bar{x}$ , which gives you the least squared error norm. So, this is known as the least square solution this is known as the least square solution or this is known as the least squares problem, in fact this is known as the least squares problem.

(Refer Slide Time: 10:55)

The image shows a whiteboard with handwritten mathematical notes. At the top, the equation  $\equiv \min \| \bar{y} - A \bar{x} \|_2^2$  is written. Below it, the text "Least Squares solution" is written in green. Underneath that, "Least Squares Problem" is written in purple. A horizontal line separates this from the text "Arises very Frequently in communication & signal processing applications." written in blue. Below the line, the equation  $\min \| \bar{y} - A \bar{x} \|_2^2$  is written again, followed by "= Quadratic Program (QP)." in green. The whiteboard has a toolbar at the top and a status bar at the bottom showing "58 / 107".

And this is very popular in communication. This arises, when we are going to express the show some example, so this arises very frequently in communication and signal processing applications ok. So, this arises very frequently in communication and (Refer Time: 11:39). And in fact if you look at, this is nothing but a this is a quadratic objective function or this also termed as a quadratic program. So, this minimize norm of  $\bar{y}$  minus  $A \bar{x}$  square. This is termed as a it is a quadratic objective function. It is termed as a quadratic program or basically a QP ok. And finding the solution of this QP gives the best estimate this is also known as the maximum likelihood estimate. So, the solution of this gives the least squares problem, gives the maximum likelihood estimate. In fact, strictly speaking the maximum likelihood estimate in Gaussian noise ok.



(Refer Slide Time: 12:28)

$\min \|y - Ax\|$   
= Quadratic Program (QP)  
Maximum Likelihood Estimate

## LEAST SQUARES SOLUTION

$$\min \|y - Ax\|^2$$

59 / 107

So, just gives the what is termed as a or basically if one asks the question what is the vector  $\bar{x}$ , which has the maximum likelihood right, which best explains the  $y$ , which has the maximum likelihood of having occurred that is the vector  $\bar{x}$ , which basically minimizes the least squares, this which is the solution to the least squares problem that is it minimizes the square or it gives the least squared error least squared norm of the error vector ok. And this can be solved as follows and it is not very complicated. So, what we do is that to find the solution. So, we want to find the least square solution, and that can be found as follows. We want to minimize norm  $y$  bar minus  $A x$  bar square.

(Refer Slide Time: 14:12)

## SOLUTION

$$\begin{aligned} & \min \|y - Ax\|^2 \\ &= (y - Ax)^T (y - Ax) \\ &= (y^T - x^T A^T) (y - Ax) \\ &= y^T y - x^T A^T y - y^T A x + x^T A A x \\ &= y^T y - 2x^T A^T y + x^T A^T A x \end{aligned}$$

59 / 107

So, what we do is basically remember norm square of a vector, this is nothing but vector transpose times itself. So, this is  $y$  minus  $A$   $x$  bar transpose times  $y$  bar minus  $A$   $x$  bar, which is equal to  $y$  bar minus  $x$  bar transpose  $A$  transpose into  $y$  bar minus  $A$   $x$  bar, which is equal to  $y$  bar transpose  $y$  bar minus  $x$  bar transpose  $A$  transpose  $y$  bar minus, now this is to be transpose minus  $y$  bar transpose  $A$   $x$  bar plus  $x$  bar transpose  $A$  transpose  $A$  into  $x$  bar.

Now, if you can look at this these two quantities are the same, these are the transpose of each other and scalar quantities  $x$  bar transpose  $A$  transpose  $y$  bar  $y$  bar transpose  $A$   $x$  bar. So, you can take the twice of one of these, so this will be finally simplified as  $y$  bar transpose  $y$  bar minus twice  $x$  bar transpose  $A$   $x$  bar transpose  $A$   $y$  bar plus  $x$  bar transpose  $A$  transpose  $A$   $x$  bar.

(Refer Slide Time: 15:42)

The image shows a whiteboard with handwritten mathematical equations. At the top, there is a toolbar with various drawing tools. The main content is as follows:

$$f(x) = y^T y - 2x^T A^T y + \frac{x^T A^T A x}{P}$$


---


$$\nabla_x f(x) = 0 - 2A^T y + 2A^T A x$$

$\Rightarrow 0$   
 set gradient = 0  
 to find optimal value

$$\Rightarrow A^T A x = A^T y$$

60 / 107

And now if you call this as now this is an objective function, because here you have no constraint, you have only the objective function as  $x$  bar. So, you take the gradient of  $F$  with respect to  $x$  bar. And now, we can see  $y$  bar transpose  $y$  bar gradient of that with respect to  $x$  is 0 minus twice;  $x$  bar transpose  $A$  transpose you can treat this as  $x$  bar transpose  $c$  bar, so the gradient is simply  $c$  bar. So, minus twice  $A$  transpose  $y$  bar plus  $x$  bar transpose  $A$  transpose  $x$  bar that is you can treat this as  $x$  bar; you can treat this as matrix  $P$ , it is positive semi definite,  $x$  bar transpose  $P$   $x$  bar. So, the gradient is twice  $P$   $x$  bar or twice  $A$  transpose  $A$   $y$  bar.

Now, you said this equal to 0, to find the optimal value. So, you set gradient equal to 0 to find optimal value. This implies now if you said this equal to 0, what you now this 2's cancel ok, you cancel the 2. So, what you get is A transpose A y bar equals A I am sorry A transpose A, this has to be A transpose x bar, so we have A transpose A x bar equals A transpose y bar.

(Refer Slide Time: 17:16)

Handwritten mathematical derivation on a slide:

$$\Rightarrow A^T A \hat{x} = A^T \bar{y}$$

$$\Rightarrow \hat{x} = (A^T A)^{-1} A^T \bar{y}$$

Dimensions:  $A$  is  $m \times n$ ,  $A^T A$  is  $n \times n$ .

Least Squares solution  
 Assuming  $A^T A$  to be invertible.

And what this implies is basically you have x bar or you can call it x hat, typically it used in the context of estimation. The optimal value of x is x hat equal to A transpose A inverse into A transpose y. Assuming A transpose A to be invert, because we can see A transpose A will always be a square matrix correct, because A is m cross n, and A transpose is n cross m. So, A transpose A will be n cross n matrix I am sorry, this is a transpose A will be n cross n. A transpose A inverse will also be n cross n. So, this is the least square solution.

Assuming, A transpose A is invertible, it is a very compact and elegant form. Least squares problem is also termed as the L S. And this assumes that A transpose A is invertible. This assumes A transpose A to be invertible, so that basically gives us the least square solution. And like this is one of the most fundamental problems in signal processing, and also for that matter in estimation, and communication and so far as so frequently.

A solution is very well known, and it is thought in a lot of courses. And in fact this forms this analysis of this problem from one of the staples of several course alright. And I think, this is one of the most important optimization problems with various applications that we are going to encounter in this course. So, we will stop here, and continue in other modules.

Thank you very much.