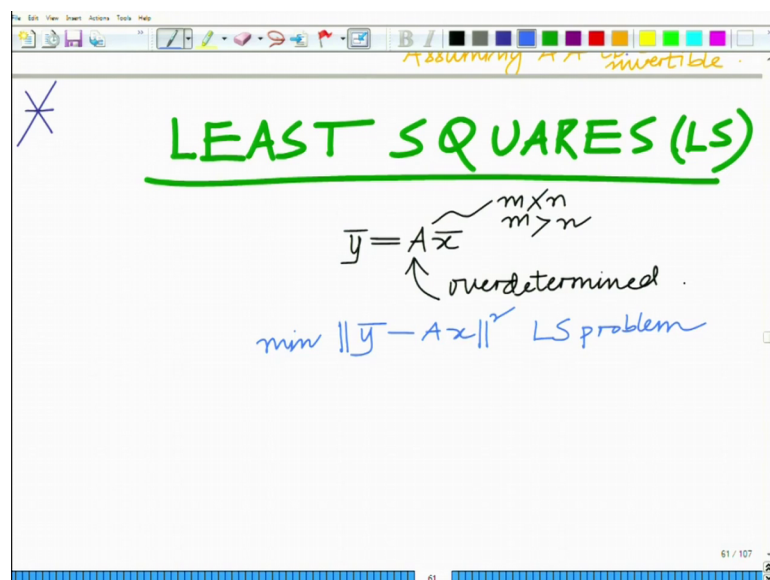


Applied Optimization for Wireless, Machine Learning, Big Data
Prof. Aditya K. Jagannatham
Department of Electrical Engineering
Indian Institute of Technology, Kanpur

Lecture – 42
Geometric Intuition for Least Squares

Hello, welcome to another module in this massive of online course. So, we are looking at the least squares optimization problem and we also derived the least squares solution that is the solution to the least squares optimization problem.

(Refer Slide Time: 00:26)



So, let us continue our discussion so, we are looking at the least squares which is a very important optimization problem. You can also think this is an approximation or modeling problem.

And what we have seen is that when we have an over determined system of equations \bar{y} equals $A \bar{x}$ with A being an m cross n matrix and m greater than n this is over determined correct. And therefore, to summarize and this cannot be solved exactly therefore, what we do is we minimize $\| \bar{y} - A \bar{x} \|^2$ this is termed as the least squares problem.

And the solution to this least square solution we have derived.

(Refer Slide Time: 01:30)

$\bar{y} = A \hat{x}$ $m > n$
↑ overdetermined
 $\min \| \bar{y} - Ax \|$ LS problem
 $\hat{x} = (A^T A)^{-1} A^T \bar{y}$
Consider $(A^T A)^T A^T$
 $\frac{n \times m \quad m \times n}{n \times n} \quad n \times m$

This previously that is \hat{x} equals $A^T A^{-1} A^T$ into \bar{y} ok. And now, if you look the now let us look at some aspect salient aspects of this solution, now consider this matrix $A^T A^{-1} A^T$ into $A^T A$ ok.

Now, consider the matrix $A^T A^{-1} A^T$ into $A^T A$. Now, we can see that you can easily say that this size of this matrix. So, A^T is m cross n , A is n cross m . So, $A^T A$ is m cross n . A^T is n cross m . So, this is basically your n cross, this is n cross m . $A^T A$ is this n cross m . So, this is n cross m , this is m cross n .

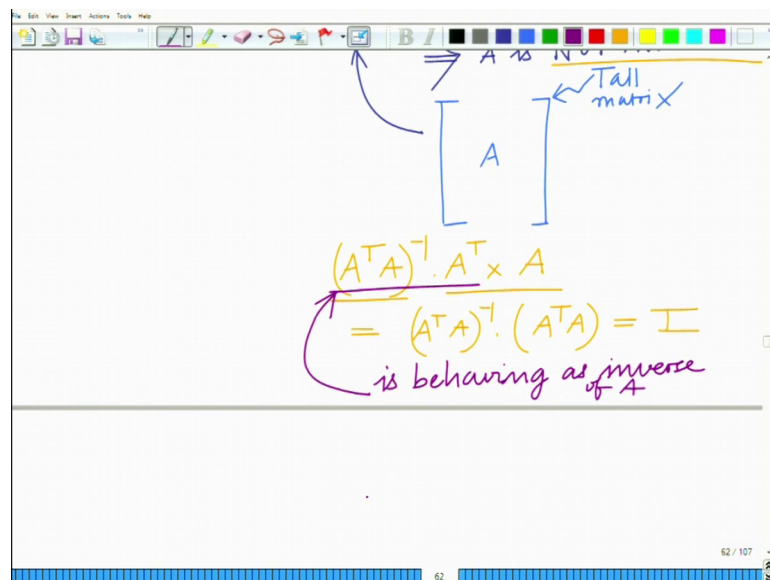
(Refer Slide Time: 03:05)

$(A^T A)^T A^T$
 $n \times m$
 $A \leftarrow m > n$
 $\Rightarrow A$ is NOT invertible
↑ Tall matrix

So, $A^T A^{-1}$ is $n \times n$ and $A^{-1} A$ is $n \times m$ implies if you look at $A^T A^{-1}$ that is therefore, $A^T A^{-1} A$ is $n \times m$ matrix ok $m > n$. So, this has more columns than rows. But, the interesting thing here is now if you look at A for $m > n$ non square matrix this implies that A is not invertible.

Now, we are considering a scenario in which $m > n$ which means the number of rows is much greater than the number of columns so, A looks like this. So, this is also known as a tall matrix that is the height number of rows of the matrix is much larger than the number of columns the matrix looks like a tall matrix. Now, for non square matrix; obviously, does not have inverse, but you can observe a very interesting aspect.

(Refer Slide Time: 04:14)

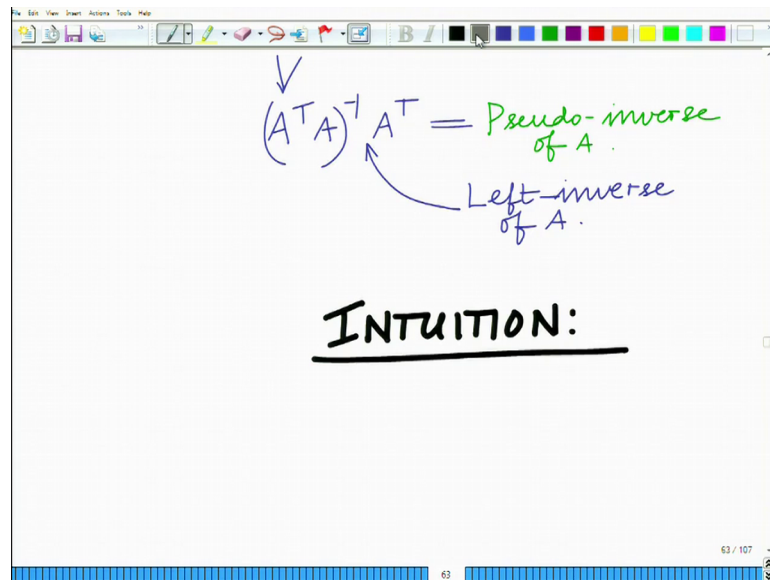


That is if you look at $A^T A^{-1}$ if you look at this matrix $A^T A^{-1} A$ is $n \times m$ and you multiply you take it is product with a now you look at this. So, A if $m > n$ it implies that it is not invertible, but if you look at this $A^T A^{-1} A$ is $n \times m$ inverse into $A^T A^{-1}$ multiply this with A .

Now, look at this we have $A^T A^{-1}$ we have a transpose. So, this is basically your $A^T A^{-1}$ times $A^T A$ which is identity. So, it is as if all the A is not invertible, it is as if this matrix $A^T A^{-1}$ is inverse into $A^T A$ this matrix acts is acting, is behaving you can say is behaving as an inverse of A right. When multiplied on the left with A it is giving identity.

Let us it behaving it not an inverse because A is not invertible when m is greater than n, but it is behaving as an inverse this is therefore, known as the pseudo inverse of A. Pseudo is an quantity pseudo when it is not actually the quantity, but it gives the appearance of that quantity ok.

(Refer Slide Time: 05:40)



So, that is basically termed as this quantity A transpose A inverse into A transpose this is termed as the pseudo. This is termed as pseudo inverse of A and this also remember a left inverse because, it is only true when you multiply it on the left.

This is also the left not the it is a left inverse of the matrix. Now, to understand the explore of this further this nature and to understand to get some intuition behind the solution. So, what you want to do is want to get some intuition behind the least square solution and the intuition is very interesting.

(Refer Slide Time: 06:48)

INTUITION:

$$\bar{y} = A \bar{x} \quad \text{No solution}$$
$$\Rightarrow \bar{y} - A \bar{x} = \bar{e}$$

Approximation Error

$$\bar{y} = \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

The intuition behind the least square solution now if you look at our problem we have \bar{y} equals $A \bar{x}$ this is our problem.

But, we know that they does not exist any \bar{x} that is this does not have any solution. So, this is does not have any solution which implies that no matter that \bar{x} you choose they it will not satisfy \bar{y} equal to \bar{y} equals $A \bar{x}$. Which means \bar{y} minus $A \bar{x}$ will always be non-zero that can be denoted by vector \bar{e} .

So, there is no vector \bar{x} so, this is an over determined system. Remember that is what we said unless and if you consider 3 equations 3 equations that 2 variables and you have 3 lines and unless they all intersected a point, it does not have any solutions. So, there will always be an approximation here let us denote that by \bar{e} . So, this is the approximation error.

So, this is the approximation error. Now, therefore, \bar{y} minus $A \bar{x}$ equals \bar{e} this is a approximation error now let me write it a little explicitly. So, this is will be \bar{y} minus A is an m cross n matrix which means it has n columns. So, these I can denote by a_1 a_2 upto a_n times x_1 x_2 up to x_n . So, this is your matrix A and this is \bar{x} .

(Refer Slide Time: 08:38)

$$\vec{y} - (x_1 \vec{a}_1 + x_2 \vec{a}_2 + \dots + x_n \vec{a}_n) = \vec{e}$$

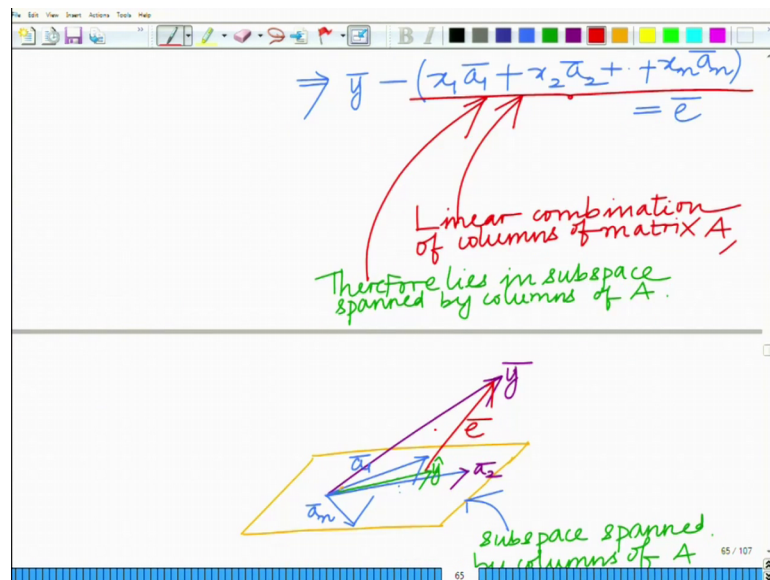
Linear combination of columns of matrix A

So, these are the n columns of the matrix A and so, this is equal to \vec{e} and. So, this implies that \vec{y} minus now we should multiply this out. So, you will get a 1 bar times x_1 times \vec{a}_1 plus x_2 times \vec{a}_2 bar.

So, this minus x_1 times \vec{a}_1 plus x_2 times \vec{a}_2 plus so on x_n times \vec{a}_n this is equal to \vec{e} where \vec{e} is a vector and now if you look at this, if you look at this x_1 times \vec{a}_1 explicit x_2 times \vec{a}_2 times so, on until x_n times \vec{a}_n . And now this is nothing, but a linear combination of the columns of matrix A . What this is? It represents a linear combination of the columns of the matrix A .

So, now, you have this linear combination of the columns which implies that this approximation that is $A \vec{x}$ always lies in the subspace spanned by the columns of A . Remember linear combination you can consider all the linear combinations of these columns you get the subspaces.

(Refer Slide Time: 10:27)

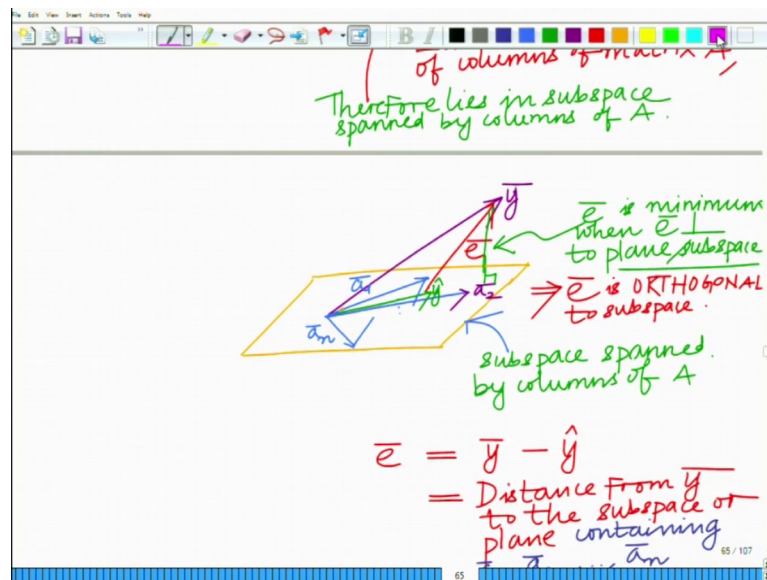


So, which means therefore, this lies in the subspace therefore, this lies in the subspace spanned by columns of A . what you have over here now to represent this pictorially if you take let us say this plane remember a subspace a plane is nothing, but a subspace. So, let us say you have the subspace which is spanned just to give you an idea I am presenting the subspace by a plane.

So, let us say this is the subspace spanned by columns of A . Let us say this subspace by columns of A and you have your vector \bar{y} which does not necessarily lie in the subspace. And you are trying to form an approximation which lies in this subspace. So, this is your approximation let us denote this approximation this approximation by \hat{y} .

Ah this approximation let us say this approximation is let us say you denote this by \hat{y} and this approximation and let us say you are let me just. So, this is your \bar{a}_2 and this approximation is let us say you denote this by \hat{y} . Now, this $\bar{y} - \hat{y}$ this is a if you look at this, this is the corresponding error.

(Refer Slide Time: 12:51)



Ok, So, this is your corresponding error vector ok. So, this \bar{y} minus \hat{y} this is your corresponding error. And therefore, now \bar{e} equals \bar{y} minus \hat{y} and what is this? This is the distance from \bar{y} to the subspace or you can say the plane, the plane that is spanned plane or plane let us make it simple plane containing a a_1 a_2 up to a_n . So, what do you think of this as basically you have a vector \bar{y} and you have this plane that contains a_1, a_2, a_n already in this plane contain different possible approximations \hat{y} .

Now, what is the error? Error is the distance between this vector \bar{y} and \hat{y} which lies in the plane. And we want to find the error the vector error I mean we want to minimize this error that is we want to make the error vector \bar{e} such that it has the minimum. Or in other words the distance of \bar{y} from this plane has to be minimum.

And now you can see this error is the distance of \bar{y} from this plane is minimum when the error is perpendicular to the plane that is the whole point. So, this error which is nothing, but the distance geometrically you can see this error is minimum when \bar{e} is perpendicular to the. This is the important ideas so this error vector is minimum when it is perpendicular to the subspace that is spanned by a_1, a_2 up to a_n . Or we can also say that this error vector is orthogonal to the subspace this is a big idea.

(Refer Slide Time: 15:35)

The image shows a whiteboard with handwritten notes. At the top, there is a definition: "= Distance from \bar{e} to the subspace or plane containing $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ ". Below this, a pink arrow points to the equation $\bar{e} \perp \bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$. A horizontal line separates this from the bottom section, where another pink arrow points to the equations $\bar{a}_1^T \bar{e} = 0$ and $\bar{a}_2^T \bar{e} = 0$. The whiteboard has a toolbar at the top and a status bar at the bottom showing "66 / 107".

So, this \bar{e} is orthogonal when is this orthogonal to the subspace this implies that \bar{e} has to be orthogonal to each of the vectors in the subspace this means \bar{e} has to be perpendicular to \bar{a}_1 , \bar{a}_2 , up to \bar{a}_n . So, \bar{e} has to be perpendicular to the subspace and.

We know the condition for orthogonality the condition 2 vectors are orthogonal when their inner product is 0. Which means we must have we must have we must have a $\bar{a}_1^T \bar{e}$ equal to 0, $\bar{a}_2^T \bar{e}$ equal to 0, $\bar{a}_3^T \bar{e}$ equal to 0, $\bar{a}_4^T \bar{e}$ equal to 0, $\bar{a}_5^T \bar{e}$ equal to 0, $\bar{a}_6^T \bar{e}$ equal to 0, $\bar{a}_7^T \bar{e}$ equal to 0, $\bar{a}_8^T \bar{e}$ equal to 0, $\bar{a}_9^T \bar{e}$ equal to 0, $\bar{a}_{10}^T \bar{e}$ equal to 0.

(Refer Slide Time: 16:23)

The whiteboard shows the following steps:

$$\begin{aligned} \Rightarrow & \begin{cases} \mathbf{a}_1^T \bar{\mathbf{e}} = 0 \\ \mathbf{a}_2^T \bar{\mathbf{e}} = 0 \\ \vdots \\ \mathbf{a}_n^T \bar{\mathbf{e}} = 0 \end{cases} \\ \Rightarrow & \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix} \bar{\mathbf{e}} = 0 \\ \Rightarrow & \mathbf{A}^T \bar{\mathbf{e}} = 0 \end{aligned}$$

So, on a n bar transpose e bar equal to 0, that is e bar has to be orthogonal to all these vector a_1 bar, a_2 bar upto a n bar. And now you can write this you can put write this as a matrix. So, this implies basically now you can concatenate these condition as a matrix. So, this implies a_1 bar transpose, a_2 bar transpose so, on upto a n bar transpose into e bar equal to 0.

And this implies well this is nothing, but now you can see this in nothing, but the matrix A transpose. So, this implies A transpose e bar equal to 0.

(Refer Slide Time: 17:13)

The whiteboard shows the following steps:

$$\begin{aligned} \Rightarrow & \mathbf{A}^T (\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}) = 0 \\ \Rightarrow & \mathbf{A}^T \mathbf{y} = (\mathbf{A}^T \mathbf{A}) \hat{\mathbf{x}} \\ \Rightarrow & \boxed{\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}} \end{aligned}$$

Least Squares (LS) solution

But, \bar{e} is \bar{y} minus $A \bar{x}$ so, this implies $A^T \bar{y}$ minus $A^T A \bar{x}$ equal to 0. This implies that now you observe something interesting $A^T \bar{y}$ equals $A^T A \bar{x}$ which is a condition that we have already seen which implies that \bar{x} or the best vector \bar{x} that minimizes the error that is \hat{x} equals $A^T A^{-1} A^T \bar{y}$ ok. So, this implies that \hat{x} equals $A^T A^{-1} A^T \bar{y}$. So, this implies this is nothing, but again you get the least square solution which is basically exactly.

So, intuitively what the least square solution is doing is basically finding the vector \hat{y} which is the best approximation to \bar{y} in the subspace that is spanned by the vectors $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ which are basically nothing, but the columns of this matrix A . And therefore, what now what is this so, therefore, that so, therefore, which implies the error vector \bar{e} which is the distance of \bar{y} to \hat{y} or basically distance of \bar{y} to the plane is minimum when the error vector is perpendicular to the plane which implies that the error vector has to be perpendicular to all the vector $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ which spanned the plane.

Now, in addition what is \hat{y} ? Now, we can see \hat{y} that is.

(Refer Slide Time: 19:15)

The slide contains a diagram on the left showing a 3D coordinate system with a blue plane. A yellow vector \bar{y} is shown above the plane, and its projection \hat{y} is shown on the plane. A red vector \bar{e} connects \bar{y} to \hat{y} and is perpendicular to the plane. To the right of the diagram, the following equations and text are written:

$$\hat{y} = A \hat{x}$$

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

\hat{y} = Projection of \bar{y} in subspace of $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$

$\Rightarrow A(A^T A)^{-1} A^T$ = Projection matrix

For subspace spanned by $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$

Approximation to \bar{y} , \hat{y} is nothing, but $A \hat{x}$ which is equal to $A^T A^{-1} A^T \bar{y}$. Now, what is \bar{y} ? Now, what is \hat{y} ? Remember, \hat{y} is the best approximation. Now, if you look at this plane again, go back

and look at this plane again this is your vector \bar{y} , this is your vector \hat{y} and the error vector is orthogonal. And the resulting error vector is orthogonal ok.

And therefore, what is \hat{y} now \hat{y} you can see is the best approximation to \bar{y} in the plane or you can also say \hat{y} is the projection of \bar{y} in the plane or subspace containing. So, \hat{y} equals projection in the best approximation, in the subspace of \bar{y} . \hat{y} is a projection of \bar{y} in the subspace or spanned by $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ this matrix which is giving you the projections. So, when you multiply this matrix by \bar{y} you get the projection. So, this implies that this matrix is the projection matrix.

This implies that $A(A^T A)^{-1} A^T$ is the projection matrix and the projection matrix for what? Projection matrix for the subspace that is spanned by $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ that is it.

(Refer Slide Time: 21:40)

The image shows a handwritten derivation of the projection matrix P_A and its idempotent property. The derivation is as follows:

$$P_A = A(A^T A)^{-1} A^T$$

$$P_A \cdot P_A = A(A^T A)^{-1} \cancel{A^T A} (A^T A)^{-1} A^T$$

$$= \underbrace{A(A^T A)^{-1}}_{P_A} \underbrace{(A^T A)^{-1} A^T}_{P_A}$$

$$= A(A^T A)^{-1} A^T = P_A$$

The final result is boxed: $P_A^2 = P_A$.

So, this matrix which is very interesting, this matrix which is the projection matrix that is spanned this is basically given as $A(A^T A)^{-1} A^T$. This is the projection matrix corresponding to the subspace that is spanned by the columns of the matrix A ok.

So, this is basically your projection matrix and this you can see this is very interesting properties. One of the most interesting properties of the projection matrix is that if you

look $P A$ times $P A$ that gives you $A A^T A^{-1} A^T$ times multiplied by this $A A^T A^{-1}$ into A^T .

So, this is A multiplying it by $P A$ and you can see now you have $A^T A^T A^{-1}$. So, these things cancel and what you are left with is again you can see $A A^T A^{-1}$ into A^T which is $P A$. So, this satisfies the property $P A^2 = P A$.

In fact, $P A^2 = P A$, $P A$ raised to the power of n for any integer n greater than or equal to 1 is $P A$. So, this is the projection matrix and this is basically the intuition behind the least square solution which sheds which basically very conducive to sort of intuitively understanding the reasoning and the methodology behind the least square solution all right. So, we will stop here and continue in the subsequent modulus.

Thank you very much.