

## Applied Multivariate Analysis

Prof. Amit Mitra

Prof. Shramishtha Mitra

Department of Mathematics and Statistics

Indian Institute of Technology, Kanpur

Lecture No. # 39

Factor Analysis

In this lecture, we continue our discussion on methods of estimation of  $L$  and  $\psi$  of a factor model. Now in the last lecture, we had discussed about the principal component method for estimation of  $L$  and  $\psi$ ; and we had come up to the point of looking at the closeness of comparison of approximating  $S$  by  $\hat{L}\hat{L}' + \psi$ .

(Refer Slide Time: 00:42)

Closeness of approximation

Result:  $\Delta = S - (\hat{L}\hat{L}' + \Psi) = (\delta_{ij})$

Then  $\left( \sum_{i,j} \delta_{ij}^2 \right) = \text{tr } \Delta^2 \leq \sum_{i=1}^p \lambda_i^2$

P.S: The diagonal entries of  $\Delta$  matrix are zero

Sum of squares of entries of  $(S - \hat{L}\hat{L}' - \Psi)$   
 $\leq$  Sum of squares of entries of  $(S - \hat{L}\hat{L}')$

i.e.  $\text{tr } \Delta^2 \left( = \sum_{i,j} \delta_{ij}^2 \right) \leq$  Sum of squares of entries of  $(S - \hat{L}\hat{L}')$ .

So, we had stated this result in the last lecture that denoting by  $\Delta$ , this matrix which is  $S$  minus  $\hat{L}\hat{L}' + \psi$   $\Delta_{ij}$  being the  $i, j$ th element of this  $\Delta$  matrix. Now, as we had discussed in the last lecture that, the way that this  $\hat{L}$  and  $\psi$  is formed, the diagonal entries of  $S$  minus  $\hat{L}\hat{L}' + \psi$  this element, the diagonal entries will be equal to 0; and only the off diagonal entries of this  $\Delta$  matrix is non-zero. So, since we had approximated  $S$  by  $\hat{L}\hat{L}' + \psi$ , we have that particular observation. Now under such a situation, we will have this sum of squares

of the entries of this delta matrix bounded by the sum of squares of the Eigen values, from  $m + 1$  to up to  $p$ , the smallest  $p - m$  Eigen values to be the quantity, which actually bounds this term here, which is the sum of squares of all the entries of this delta matrix.

So, the results that states that summation over  $i, j$  of  $\delta_{ij}^2$ , this is lesser or equal to  $\sum_{i=m+1}^p \lambda_i^2$ . So, the first thing that we will do in this lecture is to prove this particular result, which gives us a bound on the approximation, which we get from the factor, the principal component method. Now, first we realize that the diagonal entries of **the diagonal entries of** delta matrix are 0, and thus the sum of squares of entries of  $S - L \hat{L}^T$  minus  $\psi$ , this is lesser or equal to the sum of squares of entries of this  $S - L \hat{L}^T$ .

Because we are going to have the diagonal entries of these to be equal to 0, the diagonal entries of  $S - L \hat{L}^T$  are non-zero, and hence the sum of squares of the of all the entries of  $S - L \hat{L}^T$  minus  $\psi$  will be lesser or equal to the sum of square **entries sum of squares** of entries of  $S - L \hat{L}^T$ . Now, if we look at this left hand side here that is nothing but trace of  $\delta^2$  which of course, as we have denoted that is equal to double summation over  $i, j$  this  $\delta_{ij}^2$ . So, this is less than or equal to sum of squares of entries of this  $S - L \hat{L}^T$  matrix. Now, we will look at what this sum of squares entries of  $S - L \hat{L}^T$  is, and then we will have that quantity to be equal to this upper bound.

(Refer Slide Time: 04:18)

$$S = (\sqrt{\lambda_1} \hat{e}_1, \dots, \sqrt{\lambda_p} \hat{e}_p) \begin{pmatrix} \sqrt{\lambda_1} \hat{e}_1^T \\ \vdots \\ \sqrt{\lambda_p} \hat{e}_p^T \end{pmatrix}$$

$$\text{i.e. } S = (\hat{\lambda}_1 \hat{e}_1 \hat{e}_1^T + \dots + \hat{\lambda}_m \hat{e}_m \hat{e}_m^T) + (\hat{\lambda}_{m+1} \hat{e}_{m+1} \hat{e}_{m+1}^T + \dots + \hat{\lambda}_p \hat{e}_p \hat{e}_p^T)$$

$$\text{i.e. } S = \hat{L} \hat{L}^T + \sum_{j=m+1}^p \hat{\lambda}_j \hat{e}_j \hat{e}_j^T$$

$$\Rightarrow (S - \hat{L} \hat{L}^T) = \sum_{j=m+1}^p \hat{\lambda}_j \hat{e}_j \hat{e}_j^T$$

Now, let us look at how we had defined  $S$  through the Eigen value Eigen vector pairs. So,  $S$  as it was defined, its  $\lambda_1$ ,  $e_1$ , then this was of order  $p$ . So, this is  $\lambda_p$ ,  $e_p$ , this in to the transpose matrix, which is  $\lambda_1$ ,  $e_1$  transpose, and the last entry being  $\lambda_p$  times  $e_p$  transpose vector **right**; that is this  $S$  is equal to  $\lambda_1$ ,  $e_1$ ,  $e_1$  transpose plus I will just write this  $m$  eth term here, the  $m$  eth term is  $\lambda_m$ ,  $e_m$ ,  $e_m$  transpose, this plus  $\lambda_{m+1}$  plus  $e_{m+1}$ ,  $e_{m+1}$  transpose up to the last term, which is  $\lambda_p$ ,  $e_p$ ,  $e_p$  transpose **right**. Now, if you look at this part here, there is a reason why I have written all these terms here up to  $m$  and beyond  $m$ .

So, if we look at the first  $m$  terms, these first  $m$  terms are going to come, if we look at  $L$  hat  $L$  hat transpose matrix. In other words, this  $S$  can be written as this is the first  $m$  terms contribution to  $S$  is basically  $L$  hat  $L$  hat transpose, this plus I will just put it as a sum summation  $j$  equal to  $m+1$  to up to  $p$ , this is  $\lambda_j$  hat  $e_j$  hat in to  $e_j$  hat transpose. So, this will imply that this  $S$  minus  $L$  hat  $L$  hat transpose, this matrix is equal to the terms, which we have neglected; assuming that  $\lambda_{m+1}$  up to  $\lambda_p$  they are having a negligible contribution, this is  $\lambda_j$  hat  $e_j$  hat  $e_j$  hat transpose. So, we have this  $S$  minus  $L$  hat  $L$  hat transpose to be given by this right hand side here. Now, we will use this form of  $S$  minus  $L$  hat  $L$  hat transpose in the previous equation, sum of the **squares** square entries of  $S$  minus  $L$  hat  $L$  hat transpose; now this of course, is a symmetric matrix.

(Refer Slide Time: 07:28)

$$\begin{aligned}
 \text{tr } \Delta^2 &\leq \text{tr } (S - \hat{L}\hat{L}')^2 \\
 \text{tr } (S - \hat{L}\hat{L}')^2 &= \text{tr } (S - \hat{L}\hat{L}')(S - \hat{L}\hat{L}') \\
 &= \text{tr} \left[ \left( \sum_{j=m+1}^p \hat{\lambda}_j \hat{e}_j \hat{e}_j' \right) \left( \sum_{j=m+1}^p \hat{\lambda}_j \hat{e}_j \hat{e}_j' \right) \right] \\
 &= \text{tr} \sum_{j=m+1}^p (\hat{\lambda}_j \hat{e}_j \hat{e}_j') (\hat{\lambda}_j \hat{e}_j \hat{e}_j') \\
 &= \text{tr} \sum_{j=m+1}^p \hat{\lambda}_j^2 \hat{e}_j \hat{e}_j' \\
 &= \sum_{j=m+1}^p \hat{\lambda}_j^2 \text{tr} (\hat{e}_j \hat{e}_j') \\
 &= \sum_{j=m+1}^p \hat{\lambda}_j^2 \text{tr} \begin{pmatrix} \hat{e}_j' & \hat{e}_j \\ \hat{e}_j & \hat{e}_j' \end{pmatrix} \\
 &= \sum_{j=m+1}^p \hat{\lambda}_j^2
 \end{aligned}$$

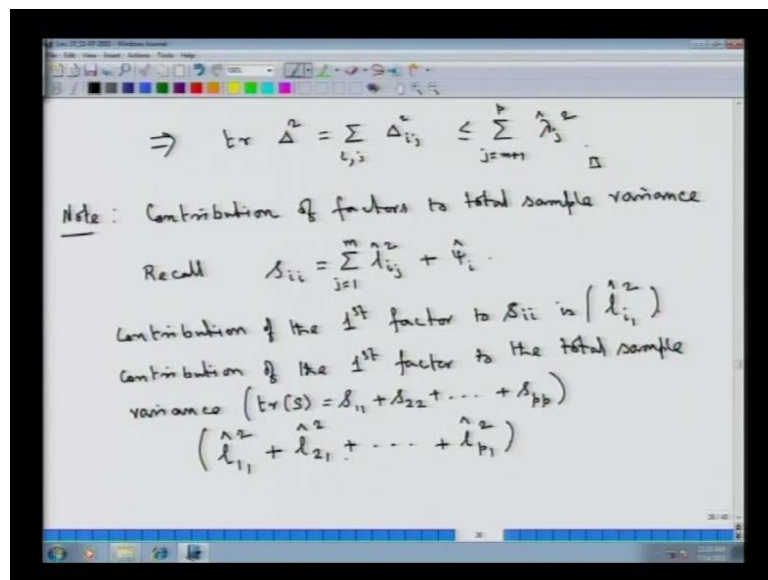
So, what we are going to have is that this trace of delta square from the previous equation is less than or equal to trace of this  $S$  minus  $L$  hat  $L$  hat transpose this square why is that so? Because we are looking at sum of squares entries of this symmetric matrix, and hence we are writing that in terms of the trace of this  $S$  minus  $L$  hat  $L$  hat transpose square. Now, what is this equal to let us look at what is the trace of  $S$  minus  $L$  hat  $L$  hat transpose square matrix, this is trace of  $I$  just write the two matrices  $L$  hat  $L$  hat transpose this into  $S$  minus  $L$  hat  $L$  hat transpose. Now, we will plug in the value of  $S$  minus  $L$  hat minus  $S$  minus  $L$  hat  $L$  hat transpose, as in this, in this expression here, and then we will be able to write it as this is summation over  $j$  equal to  $m$  plus  $1$  to up to  $p$ , this is  $\lambda_j$  hat  $e_j$  hat times  $e_j$  hat transpose. So, that is this matrix, and then product by the same thing here  $j$  equal to  $m$  plus  $1$  to up to  $p$   $\lambda_j$  hat  $e_j$  hat times  $e_j$  hat transpose **right**.

Now, if we look at this particular product here, inside the trace, what is going to happen is that we recall that this  $e_1$  hat  $e_2$  hat up to  $e_p$  hat, they are orthonormal Eigen vectors corresponding to the **corresponding corresponding** to the Eigen values. And hence, if we are looking at a particular  $j$  here, and a  $j$  prime in this summation, then the product will make that product will be equal to  $0$ , because we will be looking at multiplication of  $e_j$  hat prime times, if we look at another  $j$  prime here, then it will be  $e_j$  hat prime multiplied by  $e_j$  prime and that would be equal to  $0$ , because the vectors are orthonormal.

And hence, only the terms, for which the same index in this sum, and the same index in this sum are multiplied, and that would lead us to the following expression, which is summation  $j$  equal to  $m + 1$  to up to  $p$   $\lambda_j$  times this  $e_j$  hat  $e_j$  hat transpose times the same quantity from the other summation say in index; so, that is  $\lambda_j$  hat times  $e_j$  hat times  $e_j$  hat transpose. Now all other terms, for which the indices in these two summations differ, they will be equal to 0, because of the orthogonal properties of this  $e_j$  vectors.

So, what do we get after this multiplication? This  $\lambda_j$  hat anyways is scalar, so this  $e_j$  hat transpose  $e_j$  hat that would be equal to 1, because we have  $e_j$  hat vectors to be ortho normal. So, this is nothing but trace of the matrix, which is  $j$  equal to  $m + 1$  to up to  $p$   $\lambda_j$  hat square  $e_j$  hat  $e_j$  hat transpose. Now, we will (( )) take the trace some term by term, because trace of the sum is sum of the traces. So,  $\lambda_j$  hat square anyway is a scalar quantity. So, we take the trace straight away inside trace of  $e_j$  hat times  $e_j$  hat transpose; and furthermore the trace of a  $b$  equal to trace of  $b a$ . So, what we are going to have is summation  $j$  equal to  $m + 1$  to up to  $p$   $\lambda_j$  hat square trace of  $e_j$  hat transpose  $e_j$  hat and that is equal to 1. So, what we have is this summation to be just equal to  $j$  equal to  $m + 1$  to up to  $p$   $\lambda_j$  hat square.

(Refer Slide Time: 12:18)



So, this would imply the desired result that trace of delta square, which is summation, double summation  $i j$  delta  $i j$  square that is less than or equal to the term that we have

got in here. So, that is summation  $j$  equal to  $m + 1$  to up to  $p$   $\lambda_j$  hat square. So, this proves the result that this is quantifying the closeness of approximation that this matrix, which is the difference between  $S$  and the matrix, which approximates  $S$  in this principle component method for estimation of  $L$  and  $\psi$  that this is bounded by this quantity on the right hand side. So, after proving this particular result, we move on further, and look at another aspect of this principle component method for estimation of  $L$  and  $\psi$ , which will give us idea about the contribution of factors to total sample variance.

Now, in order to see the type of contributions that a particular factor would be having contribution of factors to total sample variance; we recall that in this principle component based method, we have this  $S_i$  to be equal to summation  $L_{ij}$  square,  $S_i$  equal to summation  $L_{ij}$  square,  $j$  equal to 1 to up to  $m$ , this plus  $L_{ij}$  hat square actually, this plus  $\psi_i$  hat. So, this is what is the expression for approximation that we are going to approximate  $S_{ii}$ , by this particular term. And from here, we can say that the contribution of the first factor to  $S_{ii}$  is going to be given by...

Now, if you look at this particular term here, this is  $L_{i1}$  hat square  $L_{i2}$  hat square and so on up to  $L_{im}$  hat square. So, the contribution of the first factor is related to the term for which  $j$  is equal to 1 in this expression. So, what we are going to have is that it is  $L_{i1}$  hat square. Now, if we look at the total sample variance, and then the contribution of this first factor to the total sample variance to the total sample variance; now what is total sample variance? Total sample variance is trace of the  $S$  matrix, which is equal to  $S_{11}$  plus  $S_{22}$  plus up to  $S_{pp}$ , this is a  $p$  dimensional random vector that we are considering.

So, the contribution of the first factor to the total sample variance would be from these expressions. So, we will have that contribution to be equal to  $L_{11}$  hat square, so this is the contribution of the first factor in  $S_{11}$ . Similarly, the contribution of the first factor to  $S_{22}$  would be given by this expression this is  $L_{21}$  hat square and so on, this up to the contribution of the first factor on the  $p$  eth component that is  $S_{pp}$  would be given by  $L_{p1}$  hat this particular term. So, this is the total contribution of the first factor to the total sample variance, which is given by trace of  $S$  which is summation of  $S_{ii}$   $i$  equal to 1 to up to  $p$  is this particular quantity.

(Refer Slide Time: 17:00)

The whiteboard shows the following mathematical expressions:

$$\hat{L}_{p \times m} = \begin{pmatrix} \sqrt{\hat{\lambda}_1} \hat{e}_1 & \dots & \sqrt{\hat{\lambda}_m} \hat{e}_m \end{pmatrix}$$

$j^{\text{th}}$  column  $\begin{pmatrix} \hat{\lambda}_{1j} \\ \vdots \\ \hat{\lambda}_{pj} \end{pmatrix} = \sqrt{\hat{\lambda}_j} \hat{e}_j$  ✓

$j=1 \rightarrow 1^{\text{st}}$  column of  $\hat{L} \rightarrow \begin{pmatrix} \hat{\lambda}_{11} \\ \vdots \\ \hat{\lambda}_{p1} \end{pmatrix}$

$$\sum_{i=1}^p \hat{\lambda}_{ij}^2 = (\sqrt{\hat{\lambda}_j} \hat{e}_j)' (\sqrt{\hat{\lambda}_j} \hat{e}_j) = \hat{\lambda}_j \hat{e}_j' \hat{e}_j = \hat{\lambda}_j \Rightarrow \sum_{i=1}^p \hat{\lambda}_{ij}^2 = \hat{\lambda}_j$$

Now, it is interesting to see what this term is actually equal to if we look at the  $\hat{L}$  hat matrix in this principal component based methods; this  $\hat{L}$  hat, which is of the dimension of  $L$  by  $p$  **I am sorry**  $p$  by  $m$ , this is given by root over lambda 1 hat  $e_1$  hat and so on; this is truncated up to the  $m$  term the number of factors, so this is lambda  $m$  hat times  $e_m$  hat **right**. So, if we look at this, then the  $j$  eth column here is this the  $j$  column is say I write that as  $L_{1j}$  hat  $L_{pj}$  hat; now that is equal to root over lambda  $j$  hat times this  $e_j$  hat vector. So, we have this as the  $j$  eth column.

Now if you look back at this contribution term what we have got this was the contribution of the first factor to the total sample variance, which was  $L_{11}$  hat square  $L_{21}$  hat square  $L_{p1}$  hat square. So, for  $j$  equal to 1 what we are going to have is the first column of this  $\hat{L}$  hat matrix, and that is nothing but  $L_{11}$  hat up to  $L_{p1}$  hat; why are we looking at this particular expression? Just to identify that this contribution of the first factor to the total sample variance is nothing but sum of squares of the entries of the first column of this  $\hat{L}$  hat matrix, this is the  $j$  column.

So in general, if we consider any  $j$  as in here, this summation which is  $L_{ij}$  summation  $i$  equal to 1 to up to  $p$  hat square that is this the norm of this particular vector is going to be given by root over lambda  $j$  hat  $e_j$  hat, this multiplied by its transpose lambda  $j$  hat  $e_j$  hat transpose. So, what is this is, this is equal to because this is  $e_j$  hat  $e_j$  hat transpose; so, we will have **I am sorry** this transpose is on the other side; so, its root over lambda  $j$

hat  $e_j$  transpose into  $\sqrt{\lambda_j}$  hat  $e_j$  hat. So, this is going to give us this  $\lambda_j$  hat, and this is  $e_j$  hat transpose times  $e_j$  hat, they are ortho normal; and then this is going to be just equal to 1, so that this particular term is equal to just  $\lambda_j$  hat.

Now, this will imply that the expression that we previously got as the contribution of the first factor to the total sample variance is this term is just going to be given by  $\lambda_1$  hat, because it concerns this first row here. So, for the particular value that this is  $L_{i1}$  hat square for  $i$  equal to 1 to up to  $p$ , this term is nothing but  $\lambda_1$  hat. So, this is the way, in which actually the contribution of the respective factors to the total sample variance can be calculated.

(Refer Slide Time: 21:00)

proportion of total sample variance explained  
 thro 1<sup>st</sup> factor  $\frac{\hat{\lambda}_1}{\text{tr } S} = \frac{\hat{\lambda}_1}{\sum_{i=1}^p \delta_{ii}}$   
 Slly Prop<sup>n</sup> of total sample variance explained by  
 1<sup>st</sup> k factors  $\frac{\sum_{i=1}^k \hat{\lambda}_i}{\sum_{i=1}^p \delta_{ii}}$

Now, using that we can further say that the proportion of total sample variance **total sample variance** explained or captured through first factor is going to be given by  $\lambda_1$  hat divided by this trace of  $S$ , which is going to be that summation, which is  $\lambda_1$  hat divided by summation of  $S_{ii}$  terms,  $i$  equal to 1 to up to  $p$ . In a similar way, the contribution of a say  $k$  factors, the first  $k$  factors can be defined, so this is, this also can be extended. So the proportion of total sample variance explained by first  $K$  factors would be given by the first  $k$  factors, would be associated with the first  $k$  Eigen values. So, this is  $i$  equal to 1 to up to  $k$  this divided by summation  $i$  equal to 1 to  $p$   $S_{ii}$  and so on **right**.



So, this gives us a way actually from the principle component based method that how we can actually quantify the contribution of different factors in explaining the total sample variance, which was the basic objective actually to capture the variance covariance structure, and as a byproduct we are looking at these outputs; that will conclude this method of principal components.

(Refer Slide Time: 23:03)

Method of estimation of  $L$  &  $\Psi$

Method II : Maximum likelihood estimation of  $L$  &  $\Psi$

Suppose from  $X - \mu = LF + \epsilon$ , where

have  $y = \begin{pmatrix} F \\ \epsilon \end{pmatrix} \sim N_{m+p} \left( 0, \begin{pmatrix} I_m & 0 \\ 0 & \Psi \end{pmatrix} \right)$

$\Rightarrow (X - \mu) = LF + \epsilon = \begin{pmatrix} L & I_p \end{pmatrix} \begin{pmatrix} F \\ \epsilon \end{pmatrix} \sim N_p(0, \cdot)$

Now, we will look at another method of estimation; method of estimation of  $L$  and  $\psi$ ; this is the first method that we had discussed was the principle component method. This is method number II; this is maximum likelihood estimation **maximum likelihood estimation** of  $L$  and  $\psi$ . Now, once we say that it is maximum likelihood method of estimation, we would require certain distributional assumption on this, the random variables concerning the system; and we will have to impose certain conditions, certain distributional assumptions. Now, suppose from this  $m$  factor model  $x$  minus  $\mu$  equal to  $L F$  plus  $\epsilon$ , we have the joint distribution of the random vectors involved on the right hand side; this we call the dimensions this is  $p$  by  $1$ , this  $L$  is a factor loading matrix  $p$  by  $n$ , this is the vector of  $m$  common factors, and this is the vector of  $p$  specific factors.

And suppose, we have the joint distribution of  $F$   $m$  dimensional and that augmented with  $\epsilon$  vector, which is  $p$  dimensional; so, this entire vector here, which is  $m$  plus  $p$  dimensional random vector; suppose the joint distribution of this is a  $N$   $p$  dimensional

multivariate normal, with a mean vector as a null vector; and the covariance matrix naturally has to satisfy this structure of the factor model. Now in the factor model, this  $F$  vector, which is a vector of the common factors; it has got a covariance matrix to be an identity matrix; and then the covariance matrix of this epsilon vector the vector of specific factors is having a diagonal structure, which is  $\psi$  matrix, and the covariance between  $F$  and epsilon in such a factor model needs to be a null matrix.

And hence, we assume that the joint distribution of  $F$  and epsilon has got this  $m$  plus  $p$  dimensional multivariate normal with a mean vector  $0$ , and this as the covariance matrix **right**. Now this would imply that our  $X$  minus  $\mu$ , which is  $L F$  plus epsilon; this can be written as  $L$  times an identity matrix of order  $p$ , this multiplied by this random vector, this  $F$  and epsilon. So, if you multiply this what we are going to get is  $L F$  plus this epsilon vector. Now, if we denote this  $m$  plus  $p$  dimensional random vector by  $Y$ , we have a result from multivariate distribution theory that if  $Y$  follows a multivariate normal; now this is a matrix of constants. So, this will, the distribution of this will also be a multivariate normal of the order that is determined through this what is the order - this is  $p$  by  $m$  and this is a  $p$  by  $p$  **right**.

So, this has got  $p$  rows and  $m$  plus  $p$  columns. So, the order of this matrix, which is  $L$  augmented with  $i$   $p$  is  $p$  rows and  $m$  plus  $p$  columns; and hence the dimension of this  $X$  minus  $\mu$  vector, multivariate normal, it will be a multivariate normal distribution  $N$   $p$  with a mean vector given by this multiplied by the expectation vector of  $F$  and epsilon what is that that is a null vector. So, this is a null vector and a covariance matrix, which we are going to calculate.

(Refer Slide Time: 27:36)

The image shows a whiteboard with the following handwritten text and equations:

$$\text{Cov} \left( \begin{pmatrix} L & I_p \end{pmatrix} \begin{pmatrix} F \\ \epsilon \end{pmatrix} \right)$$

$$= \begin{pmatrix} L & I_p \end{pmatrix} \begin{pmatrix} I_m & 0 \\ 0 & \Psi \end{pmatrix} \begin{pmatrix} L' \\ I_p \end{pmatrix}$$

$$= \begin{pmatrix} L & \Psi \end{pmatrix} \begin{pmatrix} L' \\ I_p \end{pmatrix} = LL' + \Psi$$

i.e.  $(\underline{x} - \underline{\mu}) \sim N_p(0, \Sigma)$ .

⇒ For observation vectors  $\underline{x}_1, \dots, \underline{x}_n$ , the likelihood  $f^n$  is given by

Now, the covariance matrix this covariance matrix of  $L \ i \ p$  that multiplied by  $F$  epsilon, this is going to be given by this is a matrix of constants. So, this is going to be given by  $L \ i \ p$  multiplied by the covariance matrix of this. Now, the covariance matrix of this is  $I \ m$  null null  $\psi$ ; so what we are going to get this is  $I \ m$  null null matrix  $\psi$  matrix; and then the transpose of this is going to come here; so that is  $L$  transpose  $I \ p$  **right**. Now if we take the multiplications, what we are going to get is that this multiplied by this is going to lead us to  $L$ , and  $L \ I \ p$  multiplied by this is going to give us  $\psi$ , and then this is  $L$  transpose  $I \ p$ . So, what do we get we get  $L \ l$  transpose plus  $\psi$  now what is  $L \ L$  transpose plus  $\psi$  in a factor model that is nothing but the sigma matrix.

So, we can fill in this particular dot here, and say that our  $x$  minus  $\mu$  has got multivariate normal  $p$  dimension with mean vector  $0$ , and a covariance matrix equal to sigma. Well, you have the covariance matrix of sigma straighter, covariance matrix of  $x$  straightaway equal to sigma anyway; but we had derived that covariance matrix of this  $L \ F$  plus epsilon through the multivariate normality distribution of this augmented  $F$  with epsilon vector. Under the assumption of joint multivariate normality of  $F$  and epsilon vector, what we have realized is that this  $X$  minus  $\mu$  has got a  $p$  dimensional multivariate normal with a mean vector  $0$  and a covariance matrix as sigma. Now, we are in a position to frame the likelihood estimation, because we can now write the form of the likelihood function for observations **observation** vectors  $x_1, x_2, x_n$ ; **for observation** **vectors  $x_1 \ x_2 \ x_n$** ; the likelihood function is given by...

(Refer Slide Time: 30:22)

The image shows a whiteboard with handwritten mathematical formulas for the likelihood function of a multivariate normal distribution. The formulas are as follows:

$$L(\mu, \Sigma | X) = (2\pi)^{-np/2} |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \text{tr} \Sigma^{-1} A - \frac{n}{2} (\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu)\right)$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

$$A = \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})'$$

log likelihood function

$$\ell(\mu, \Sigma | x_1, \dots, x_n) = -\frac{np}{2} \log 2\pi - \frac{n}{2} \log |\Sigma| - \frac{1}{2} \text{tr} \Sigma^{-1} A - \frac{n}{2} (\bar{x} - \mu)' \Sigma^{-1} (\bar{x} - \mu)$$

Let us denote that likelihood function as L; now this is going to be written as a function of mu and sigma, given the data  $x_1, x_2, \dots, x_n$ ; so given  $x$ , this is going to be given by  $2\pi$  to the power minus  $np/2$  why it is that, because we have random sample of size  $n$ , and we have each of the random samples a multivariate normal with dimension  $p$ . So, for each of these random samples, we will have a  $2\pi$  to the power minus  $p/2$  that is the dimension of the multivariate normality; and we have  $n$  such terms, and hence we have this factor as  $2\pi$  to the power minus  $np/2$ , and then we will have a determinant sigma to the power minus  $n/2$ , from where does this come; if we once again look at the structure of multivariate normal in the density of each of these, we will have determinant of sigma to the power minus half; and since we have  $n$  observations, we will have this actually leading us to determinant of sigma to the power minus  $n/2$  that multiplied by the exponent.

Now, we will straight away write the form of the exponent that we usually use, which is minus half trace of sigma inverse a this minus  $n/2$   $\bar{x}$  minus  $\mu$  transpose sigma inverse  $\bar{x}$  minus  $\mu$  **right**, where this  $\bar{x}$  vector is  $1/n$  summation  $x_i$  vectors. So, it is the sample mean vector obtained from the observations  $x_1$  vector  $x_2$  and up to  $x_n$ ; and this A matrix, this is actually the formulation of calculating the maximum likelihood estimators of multivariate normal distribution. So, that this A is  $\sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})'$ .

So, this is the standard form that we will have this as the exponent of the joint distribution of  $x_1, x_2, \dots, x_n$ ; each of them having a multivariate normal distribution, having a multivariate normal distribution with the order of the multivariate normality as  $p$ . Now if this is the likelihood function, we can also write the log likelihood, log likelihood function; let us denote that by small  $l$ , this is  $l(\mu, \Sigma | x)$  given this  $x_1, x_2, \dots, x_n$  or in short just to as we had written as  $x$  here. So, that is going to be equal to minus  $n p$  by  $2 \log 2 \pi$  this minus  $n$  by  $2 \log$  determinant of  $\Sigma$  this minus half trace of  $\Sigma$  inverse  $A$  minus  $n$  by  $2 x$  bar vector minus this  $\mu$  vector transpose  $\Sigma$  inverse  $x$  bar minus this mean vector  $\mu$  **right**.

Now, from this expression here, we what is the basic objective is to get the maximum likelihood estimators of  $L$  and  $\Psi$  **right** that is what we are aiming at. So, in order to get the maximum likelihood estimators of  $L$  and  $\Psi$ , we will write this log likelihood function in terms of those; note that although we have written the log likelihood function of  $\mu$  and  $\Sigma$ , it is in this particular  $\Sigma$  that we have  $L$  and  $\Psi$ .

(Refer Slide Time: 34:54)

Handwritten mathematical derivation on a whiteboard:

$$l(\mu, L, \Psi | x) = -\frac{np}{2} \log 2\pi - \frac{n}{2} \log |LL' + \Psi|$$

$$- \frac{1}{2} \text{tr} \left( (LL' + \Psi)^{-1} A \right)$$

$$- \frac{n}{2} (\bar{x} - \mu)' (LL' + \Psi)^{-1} (\bar{x} - \mu) \quad (*)$$

(\*) - log likelihood  $f^n$  of  $\mu, L, \Psi$  given  $(x_1, x_2, \dots, x_n)$ .

Since choice of  $L$  is not unique, we impose conditions like  $L' \Psi^{-1} L = A$  (a diagonal matrix) to make choice of  $L$  unique.

Maximization of (\*) subject to the imposed condition  $L' \Psi^{-1} L = A$ , gives the MLE of  $L \Delta \Psi$  with  $\hat{\mu}_{MLE} = \bar{x}$ .

And hence this can be written in terms of log likelihood function in terms of  $L$  and  $\Psi$  given  $x_1, x_2, \dots, x_n$ , this is going to be minus  $n p$  by  $2 \log$  of  $2 \pi$  minus  $n$  by  $2 \log$  of determinant. Now, here we are going to use the fact that  $\Sigma$  here we had determinant of  $\Sigma$  here; in place of  $\Sigma$ , we will write  $LL' + \Psi$ ; so, this is  $LL' + \Psi$  determinant of this matrix this minus **minus** half trace of  $\Sigma$  inverse  $A$ . So, this is

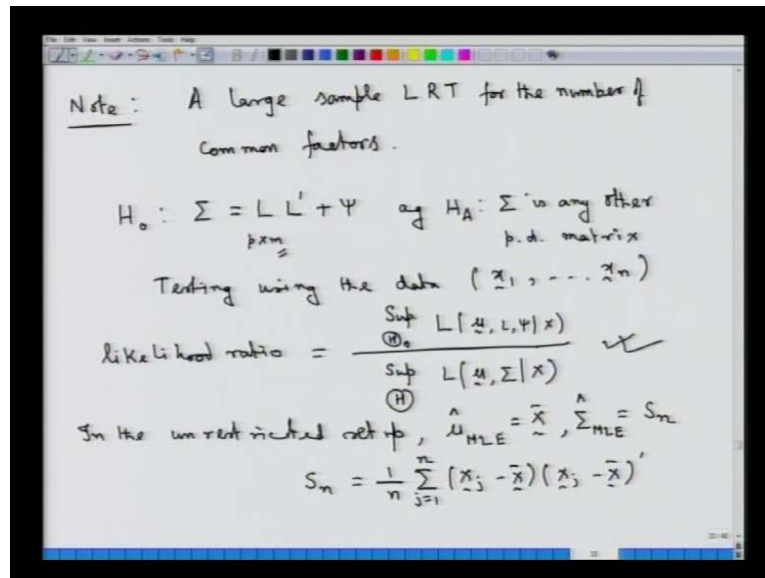
$L L^{-1} + \Psi$  whole inverse times that  $A$  matrix this minus  $n$  by  $2 \times \bar{x}$  minus  $\mu$  transpose  $\sigma$  inverse that is  $L L^{-1} + \Psi$  whole inverse times this  $\bar{x}$  minus  $\mu$ . So, what we have done is to look at this log likelihood function, this was as a function of  $\mu$  and  $\sigma$ , we transformed that actually using the factor model, just by replacing  $\sigma$  by the corresponding  $L L^{-1} + \Psi$  quantity.

So, we have this expression, which is star as the log likelihood, is the log likelihood function **log likelihood function** of this  $\mu$   $L$   $\Psi$ , given the observations set that is  $x_1, x_2, \dots, x_n$ ; and the choice of  $L$  is not unique; as we have seen in the  $n$  factor model that the choice of  $L$  is not unique, and as such we need to impose certain additional conditions in order to make  $L$  unique. And that is a type of condition, which is used in order to maximize this subject to the conditions, which would make  $L$  to be unique

So, we write this that since choice of  $L$  is not unique, we impose conditions like  $L^{-1} \Psi^{-1} L = A$  - a diagonal matrix to make  $L$  choice of  $L$  unique; and then maximization of star maximization of star, subject to the imposed condition, subject to the imposed condition  $L^{-1} \Psi^{-1} L = A$  diagonal matrix gives the maximum likelihood estimators of  $L$  and  $\Psi$ ; with of course,  $\hat{\mu}$  has the maximum likelihood estimator to be given by the usual maximum likelihood estimator, which is  $\bar{x}$ . So, using the form star as in here to be the likelihood function of  $\mu$   $L$  and  $\Psi$ , and using a condition like this to be  $A$  a given diagonal matrix; we can maximize the star with respect to this condition and arrive at the maximum likelihood estimators of  $L$  and  $\Psi$  with  $\hat{\mu}$  maximum likelihood estimator to be given by  $\bar{x}$ .

Now, there are other methods of estimation of  $L$  and  $\Psi$ , when we are talking about the factor analysis model; but we have in this course look at two such important most widely used methods of estimation of  $L$  and  $\Psi$ . The first one that we had discussed was the principal component based methods, and the second one in a more general setup, when we are looking at maximum likelihood method of estimation.

(Refer Slide Time: 40:14)



So, we will conclude this estimation part here, and to conclude actually, we will just look at 1 small note, which is important, which gives us a large sample test large sample L R T test for the number of common factors. So, what we are going to address is the following fact that the number of common factors as such is not known to us, and what is the best way to judge, what is the number of common factors that should be used in a particular x vector model. We look at a large sample likelihood ratio test; now we are doing this, because we are just now actually obtained the maximum likelihood estimators of L and psi; and this large sample likelihood ratio test is based on that maximum likelihood method of estimation. Now what we are doing here is that we are saying that a null hypothesis is of the form that sigma is equal to L L dash plus psi with the order of L to be p by n with a chosen p here.

Suppose, our interest is to test this null hypothesis against the alternate hypothesis that sigma is any other is any other positive definite matrix **any other positive definite matrix**. So, we are going to test this, if null hypothesis is accepted, then we take this m to be the number of common factors as in here; now for certain m, as we will see that this might actually get rejected, and then we will accept that this that particular choice of m does not hold good for the given observation vectors  $x_1, x_2, x_n$ . Now this testing is to be carried out using the data as in  $x_1, x_2, x_n$ ; so based on these n data vectors, we are going to test this null hypothesis against the alternate hypothesis H A using the L R T philosophy

Now, the likelihood ratio as we know, likelihood ratio is going to be given by this supremum over script theta naught of the likelihood function; now this likelihood function will be under the null hypothesis, another null hypothesis meaning thereby we will have mu L psi, given this x, the entire data vector that divided by supremum over script theta; now when we talk about script theta, it is with respect to not this factor model, it is with respect to the mean vector mu and a sigma being any **any** positive definite matrix.

So, we will have to compute this particular term, which is going to be called the likelihood ratio, when we are looking at testing of this null hypothesis; and then we will use large sample **standard large sample** theory of likelihood ratio in order to formulate this testing of H naught against H A. Now in order to do so, what we would require is these two supremum quantities. Now in the unrestricted setup **in the unrestricted setup** mu hat maximum likelihood estimator is going to be given by x bar that we have already seen time and again.

And this sigma hat M L E is going to be S n, where S n is nothing but one upon n summation x j minus x bar into x j minus x bar transpose j is equal to 1 to up to n. So, that is the maximum likelihood estimators in the unrestricted setup. So, this will actually look at the denominator part, and then in the likelihood function, we will have to plug in these values of mu as x bar, and sigma as S n, which is of this particular form.

(Refer Slide Time: 45:16)

The image shows a whiteboard with handwritten mathematical derivations. At the top, it states:  $\Rightarrow \sup_{\Theta} L(\underline{\mu}, \Sigma) \propto |S_n|^{-n/2} e^{-n/2}$ . Below this, the null hypothesis is defined as  $\Theta = \{(\underline{\mu}, \Sigma); \underline{\mu} \in \mathbb{R}^p, \Sigma > 0\}$ . The text then says "Suppose  $\hat{\underline{\mu}}_{MLE}$  &  $\hat{\Sigma}_{MLE}$  are the MLE of  $\underline{\mu}$  &  $\Sigma$  under  $\Theta_0$ ". It follows that  $\hat{\underline{\mu}}_{MLE} = \bar{x}$  and  $\hat{\Sigma}_{MLE} = \hat{L}\hat{L}' + \hat{\Psi}_{MLE}$ . The final part of the derivation shows  $\sup_{\Theta_0} L(\underline{\mu}, \Sigma) \propto |\hat{L}\hat{L}' + \Psi|^{-n/2} \exp\left[-\frac{n}{2} \text{tr}(\hat{L}\hat{L}' + \Psi)^{-1} S_n\right]$ , with a note that  $n S_n = A$ .



And then, we will be able to write this supremum over script theta, the entire parameter space; now what is entire parameter space? The entire parameter space just to recall this script theta is the set of all mu and sigma, wherein this mu belongs to  $\mathbb{R}$  to the power p, and sigma is positive definite; so that is my script theta. So, supremum over script theta of  $L(\mu, \sigma)$ , I will just drop this x, so this is going to be proportional to  $I$  am ignoring the constants. So, it is supremum over theta  $L(\mu, \sigma)$  that is going to be given by determinant of  $S_n$  to the power minus n by 2  $e$  to the power minus n p by 2. This we had of course, done when we are looking at maximum likelihood method of estimation in multivariate normal and testing; but if one wants to recall why it is so? You will have to look at this expression here. So, this is basically the likelihood function.

You plug in  $\hat{\mu}$  equal to  $\bar{x}$ , so this term would be equal to 0. So, plugging in  $\hat{\mu}$  equal to  $\bar{x}$ , this will be a null vector multiplied by whatever  $\hat{\sigma}$  you plug in this is going to be that; and this  $A$  is this matrix, which is n times. So, this

$A$  is nothing but n times  $S_n$ ; so we will have in place of sigma inverse n times  $S_n$  inverse being plugged in, and what we are going to get is the expression that written here. So, supremum over theta  $L(\mu, \sigma)$  is going to be proportional to determinant of  $S_n$  to the power minus n by 2  $e$  to the power minus n p by 2.

Now, suppose  $\hat{L}$  and  $\hat{\psi}$  are the maximum likelihood estimators of  $L$  and  $\psi$  under script theta naught. Now how we are going to obtain this  $\hat{L}$  and  $\hat{\psi}$  that is using the method of maximum likelihood estimation that we have just now touched upon. So, using those that technique in order to get to this  $\hat{L}_{MLE}$ , and  $\hat{\psi}_{MLE}$  under  $H_0$ ; under  $H_0$ , how it is going to matter? Under  $H_0$ , we will have a fixed m here; now for that fixed m, we will actually look at the m factor model; and then that is going to determine that this  $\hat{L}_{MLE}$  and  $\hat{\psi}_{MLE}$  are going to have the dimensions as what is specified through the null hypothesis that m factor.

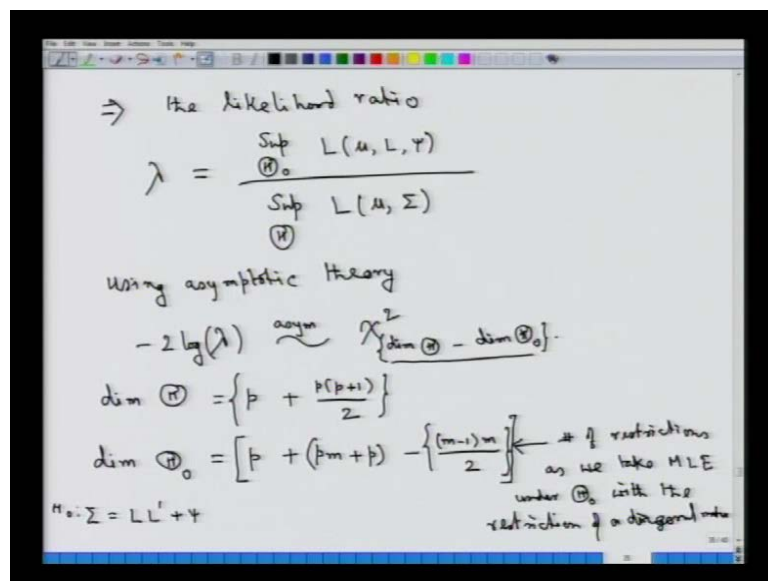
Now along with that, we will have  $\hat{\mu}_{MLE}$ , irrespective of this  $\hat{L}$   $\hat{\psi}$  that is always going to be equal to this  $\bar{x}$  term. Now if  $\hat{L}$  is a maximum likelihood estimator of  $L$ , and if  $\hat{\psi}$  is the maximum likelihood estimator of  $\psi$ , we will have the maximum likelihood estimator of  $\hat{\sigma}$  by invariance property of the maximum likelihood estimator. We will have this as  $\hat{L} \hat{L}^T + \hat{\psi}$ , where this is

the maximum likelihood estimator, as we have given in here; this also is the maximum likelihood estimator **right**.

Now, we are in a position to write this supremum over script theta naught of the likelihood functions. Now this likelihood function will be in terms of mu L psi, this given x, I will just drop it; now this one is going to be proportional to terms similar to this; now this is determinant of sigma in the likelihood function; now sigma is denoted by sigma hat, which is L hat L hat transpose plus psi hat. So, this is going to be proportional to L hat L hat transpose plus psi whole to the power minus n by 2, and here the exponents terms do not cancel out and give us a nice form like the previous one; this is going to be bit complicated, because terms do not cancel out.

The second term of course, will be equal to 0, because its x bar minus mu. So, mu hat being equal to x bar will lead the second term to be equal to 0. However, the first term is going to be trace of L hat L hat transpose plus psi hat this is sigma. So, trace of sigma inverse is this times that I will just write this as A, because it is better to write that as no if you write that as A, then this term will be given by this S hat **right**. So, we have this particular term here, so n times S n is basically A, and we have the supremum over script theta of the likelihood function given by this.

(Refer Slide Time: 50:57)



So, we can formally thus write the likelihood ratio. So, this will imply that the likelihood ratio lambda is supremum over script theta naught of L mu L **I am sorry** L mu L and psi;

this  $L$  is corresponding to the factor  $\sup_{\theta \in \Theta} L(\mu, \Sigma)$  so that this term, we can look at what terms we had got earlier, use this form here, and use this form here, in order to write this particular likelihood; and so, we can use this form, and this form here in order to write the final form of this likelihood ratio using asymptotic theory; using asymptotic theory, we know that  $-2 \log \lambda$  follows asymptotically a central chi square on the degrees of freedom, which is dimension of  $\Theta$  minus the dimension of  $\Theta_0$ .

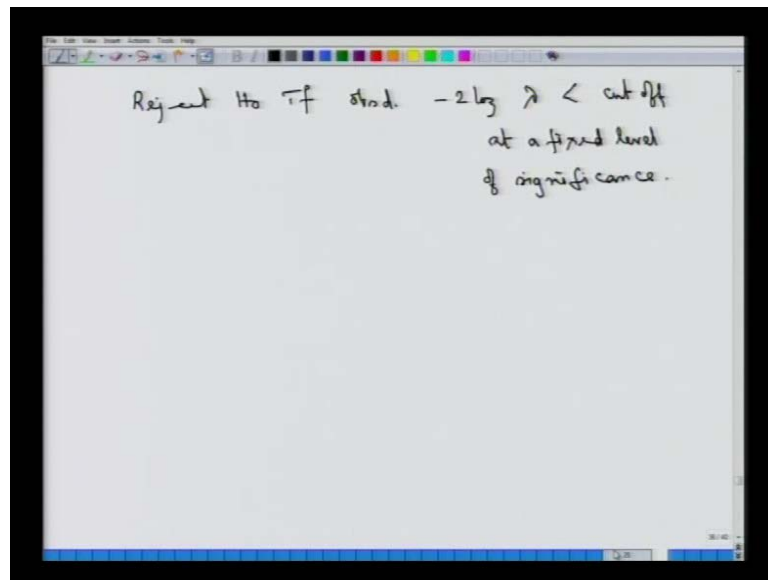
So, we will look at this likelihood ratio, and then using the asymptotic theory, we have  $-2 \log \lambda$  following asymptotically chi square central with degrees of freedom as dimension of  $\Theta$  minus dimension of  $\Theta_0$ ; from a given data, we will obtain what is observed value of  $\lambda$ , what is observed value of  $-2 \log \lambda$ , and looking at the appropriate degrees of freedom of that central chi square, we will actually look at acceptance or rejection of null hypothesis. So, what we are going to require this particular number here.

So, what we see is the dimension of  $\Theta$ ; now  $\Theta$  was this space here. So, this has got  $p$  entries and the number of distinct elements of  $\Sigma$  is going to be the dimension of this  $\Theta$ ; so that this is equal to  $p$  for  $\mu$  and  $p(p+1)/2$  for the  $\Sigma$  matrix. Now if you look at the dimension of  $\Theta_0$ ;  $\Theta_0$  says that this  $\Sigma$ , so this  $H_0$  is what we have giving us  $L_0 + \Psi$ . Now, the dimension of  $\Theta_0$ , it still has though the components of mean vector; so, that is  $p$  of them. And then we will have the elements here, and the elements here, which are going to give us, this is  $pm$ ;  $pm$  is the number of terms corresponding to the factor loading matrix  $L$ , this plus the  $p$  diagonal entries of  $\Psi$  are going to give us this.

Now, the dimension of  $\Theta_0$  is something less, because of the type of restriction that we had imposed; now this is thus going to be the number of restrictions, which is this. Now, what is or how is this one coming? This basically is coming that from the number of restrictions, as we take  $MLE$  under  $\Theta_0$  with the restriction of a diagonal matrix. So, since we have the diagonal matrix of that  $m$  by  $m$  matrix, which is you can just go back a little bit, this is basically the restriction that we plug in. So, this  $L^T \Psi^{-1} L$ , this matrix is the diagonal matrix; and hence, these are the number of restrictions that we have in computing the maximum likelihood

estimators. And hence, we will have the dimension of script theta as it is given by this and the dimensional script theta naught to be equal to this. So, there is no problem in computing the degrees of freedom of this central chi square, which is a asymptotic distribution of minus 2 log lambda.

(Refer Slide Time: 55:47)



So, we will thus reject null hypothesis; now if we look at this likelihood ratio, this is the likelihood ratio; so when we are going to reject the null hypothesis, if this particular contribution the numerator terms supremum over script theta naught is too small with respect to supremum over script theta, then we usually reject the null hypothesis. So, we will reject null hypothesis, if observed minus 2 log lambda is less than the cut off at a fixed level of significance. So, that concludes actually this testing procedure, which was based on the maximum based on the large sample theory of the likelihood ratio test for the number of **number of** common factors that we will be choosing in an m factor model. So, that concludes our discussion with about the factor analysis, there are other concepts in factor analysis like factor rotation, which also occupies an important place, but we will end the concept of factor analysis here from the next lecture, we will look at the canonical correlation analysis.