

I. Order Determination. (Ref: E. L. Lehman
Testing Statistical Hypothesis,
Wiley, N.Y.)

$$Y = X\theta^* + E$$

$$Y = (y_1, \dots, y_p)^T \quad \theta^* = (b_1, \dots, b_q)^T$$

$$E = (e_1, \dots, e_p)^T \quad \text{Cov}(E) = \sigma^2 I$$

$$X = \begin{bmatrix} u_{0-} & \dots & u_{1-q} \\ u_1 & \dots & u_{2-q} \\ \vdots & & \vdots \\ u_{p-1} & & u_{p-q} \end{bmatrix} \quad p \times q \text{ matrix (known)}$$

$p > q$; linearly independent columns.

θ^* : vector of parameters to be estimated from data.

Let d , $0 \leq d \leq q$ be given

Problem: Decide on the basis of observations

when the hypothesis

$$(H) \theta_q^* = \theta_{q-1}^* = \dots = \theta_{q-d+1}^* = 0 \text{ should be rejected.}$$

2.

Statistical Tests

Let $k_1, k_2 > 0$. A random variable X is said to have $\chi^2(k_1)$ distribution if

$$X = V_1^2 + \dots + V_{k_1}^2 \text{ where } V_i \sim N(0,1) \quad i=1, \dots, k_1$$

A random variable X is said to have a $F(k_1, k_2)$ distribution if it can be expressed as

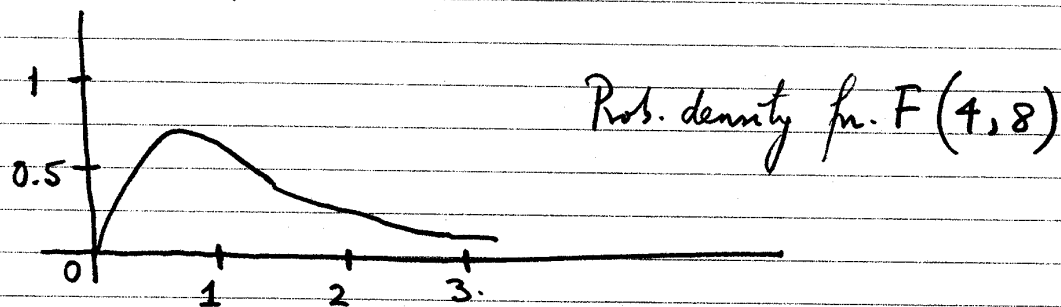
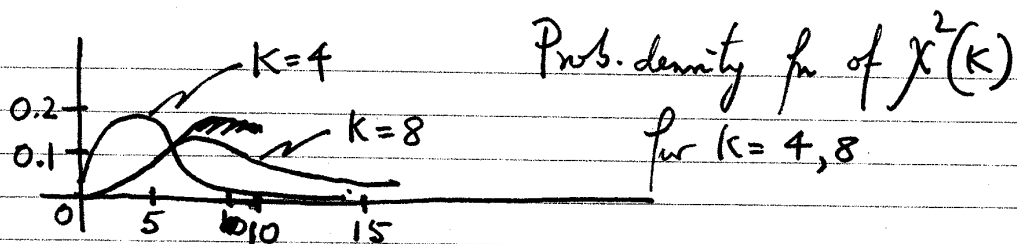
$$X = \frac{\frac{S_1}{k_1}}{\frac{S_2}{k_2}}$$

where S_1 and S_2 are independent with distributions $\chi^2(k_1)$ and $\chi^2(k_2)$ respectively.

(Define $\chi^2(0)$ to be the distribution of the degenerate random variable taking value 0, a.s.).

For large values of k_2 , if X has $F(k_1, k_2)$ distribution then to good approximation $k_1 X$ has $\chi^2(k_1)$ distribution.

3.



Let $\hat{\theta} =$ least squares estimate of θ^*
 $\hat{\theta} =$ least squares estimate of θ^* under
 hypothesis (H).

Now

$$\hat{\theta} = (\theta_0, 0, \dots, 0) \text{ where } \theta_0 = (X_0^T X_0)^{-1} X_0^T Y$$

where X_0 is $p \times q$ -d matrix obtained from X
 by removing the last d columns. (or replacing by 0).

$$\text{let } \varepsilon(\theta) = Y - X\theta, \theta \in \mathbb{R}^q \quad (1)$$

$$\text{and } S(\theta) = \varepsilon^T(\theta)\varepsilon(\theta) \quad (2)$$

Consider the statistic

$S(\hat{\theta}) - S(\hat{\theta})$ which measures the
 increase in the minimum of the least squares criterion

4.

when we decrease the number of parameters from q to $q-d$. One would reject the hypothesis if $S(\hat{\theta}) - S(\hat{\theta}_0)$ is large.

Proposition 1

Suppose that for some integer d , $0 \leq d \leq q$ hypothesis (H) is true. (For $d=0$, no restrictions on θ^*). Let $\hat{\theta}$ and $\hat{\theta}_0$ be as above and let $S(\theta)$ be defined by (2). Then $\hat{\theta}$, $S(\hat{\theta})$ and $S(\hat{\theta}_0) - S(\hat{\theta})$ are independent. Moreover

$$\frac{S(\hat{\theta})}{\sigma^2} \sim \chi^2(p-q) \text{ and } \frac{1}{\sigma^2} (S(\hat{\theta}_0) - S(\hat{\theta})) \sim \chi^2(d)$$

(Note: For $d=0$, $\hat{\theta}$ and $S(\hat{\theta})$ are independent and $\frac{1}{\sigma^2} S(\hat{\theta}) \sim \chi^2(p-q)$).

Proof: First assume: $0 < d < q$. From least squares theory

$$Y - X\hat{\theta} \perp \text{Range}(X) \quad (3)$$

$$Y - X_0\hat{\theta}_0 \perp \text{Range}(X_0) \quad (4)$$

$$\Rightarrow X(\hat{\theta} - \hat{\theta}_0) \perp \text{Range}(X_0) \quad (5)$$

5.

Now from (3) and (4)

$$\begin{aligned}\|Y - X\hat{\theta}\|^2 &= \|(Y - X\hat{\theta}) + X(\hat{\theta} - \hat{\theta})\|^2 \\ &= \|Y - X\hat{\theta}\|^2 + \|X(\hat{\theta} - \hat{\theta})\|^2\end{aligned}$$

$$\text{is. } S(\hat{\theta}) - S(\hat{\theta}) = \|X(\hat{\theta} - \hat{\theta})\|^2 \quad (6)$$

Let b_1, \dots, b_p be an orthonormal basis with the following properties

(a) b_1, \dots, b_{q-d} spans $\text{Range}(X_0)$.

(b) b_{q-d+1}, \dots, b_q are orthogonal to the $\text{Range}(X_0)$
s.t. b_1, \dots, b_q spans the range of X

(c) b_{q+1}, \dots, b_p are orthogonal to $\text{Range}(X)$.

The linear independence hypothesis on the columns of X guarantees that we can choose such a basis.

Let $B := (b_1 \dots b_p)$ and set

$$V = BE^T \quad (7)$$

(Note: to prove something has χ^2 -distribution

we have to show that the r.v. can be written

as the sum of indep. square of gaussian r.v.'s

6.

The ^{components of} V defined by (7) will be the gaussian r.v.'s

Since the columns of B are ~~independent~~ orthonormal

$$B^T = B^{-1} \quad (8)$$

and $\text{Cov}(V) = \sigma^2 I$ and hence $V \sim N(0, \sigma^2 I)$.

Now

$$E = Y - X\theta^* = (Y - X\hat{\theta}) + X(\hat{\theta} - \hat{\theta}) + X(\hat{\theta} - \theta^*)$$

Hence using (8)

$$V = B^{-1}(Y - X\hat{\theta}) + B^{-1}X(\hat{\theta} - \hat{\theta}) + B^{-1}X(\hat{\theta} - \theta^*) \quad (9)$$

Now the mapping $x \mapsto B^{-1}x$ transforms the ^{coordinates from the} standard basis of \mathbb{R}^p to the coordinates w.r. to the orthonormal basis (b_1, \dots, b_p) . Now

$$Y - X\hat{\theta} \perp \text{to } b_1, \dots, b_p \text{ and hence } B^{-1}(Y - X\hat{\theta}) \in \left\{ \xi \in \mathbb{R}^p \mid \xi_1 = \dots = \xi_p = 0 \right\} \quad (10)$$

From (5) $X(\hat{\theta} - \hat{\theta}) \in \text{Range}(X)$ but is orthogonal to range of X_0 .

7.

and hence $\bar{B}^{-1} X(\hat{\theta} - \hat{\theta}) \in \left\{ \xi \in \mathbb{R}^p \mid \xi_1 = \dots = \xi_{q-d} = \xi_{q+1} = \dots = \xi_p = 0 \right\}$

Finally, $X(\hat{\theta} - \theta^*) \in \text{Range}(X_0)$ and hence

$$\bar{B}^{-1} X(\hat{\theta} - \theta^*) \in \left\{ \xi \in \mathbb{R}^p \mid \xi_{q-d+1} = \dots = \xi_p = 0 \right\}$$

Therefore:

$$B^T (Y - X\hat{\theta}) = (0, \dots, 0, v_{q+1}, \dots, v_p)^T \quad (11)$$

$$B^T X(\hat{\theta} - \hat{\theta}) = (0, \dots, 0, v_{q-d+1}, \dots, v_q, 0, \dots, 0)^T \quad (12)$$

$$B^T X(\hat{\theta} - \theta^*) = (v_1, \dots, v_{q-d}, 0, \dots, 0)^T \quad (13)$$

Therefore from (6)

$$\begin{aligned} S(\hat{\theta}) - S(\hat{\theta}) &= \|X(\hat{\theta} - \hat{\theta})\|^2 = \|B^T X(\hat{\theta} - \hat{\theta})\|^2 \\ &= \sum_{i=q-d+1}^q v_i^2 \quad (\text{from (12)}) \quad (14) \end{aligned}$$

$$S(\hat{\theta}) = \|Y - X\hat{\theta}\|^2 = \|B^T (Y - X\hat{\theta})\|^2 = \sum_{i=q+1}^p v_i^2 \quad (15)$$

Now

$$\begin{aligned} \hat{\theta} - \theta^* &= (X^T X)^{-1} X^T X(\hat{\theta} - \theta^*) = (X^T X)^{-1} X^T B B^T X(\hat{\theta} - \theta^*) \\ &= (X^T X)^{-1} X^T B (v_1, \dots, v_{q-d}, 0, \dots, 0)^T \end{aligned}$$

that is

$$\hat{\theta} = \theta^* + (X^T X)^{-1} X^T B (v_1, \dots, v_{q-d}, 0, \dots, 0)^T. \quad (16)$$

From (14),

$$\frac{1}{\sigma^2} S(\hat{\theta}) - S(\hat{\theta}) \text{ has distribution } \chi^2(d)$$

From (15),

$$S(\hat{\theta}) \text{ has distribution } \chi^2(p-q).$$

and (14), (15) and (16) shows that

$$S(\hat{\theta}) - S(\hat{\theta}), S(\hat{\theta}) \text{ and } \hat{\theta} \text{ are independent.}$$

The case $d=0$ and q can be dealt with using similar arguments.

Case σ^2 known

The above proposition shows that the statistic

$$\frac{1}{\sigma^2} [S(\hat{\theta}) - S(\hat{\theta})] \text{ has distribution } \chi^2(d)$$

under the ~~statistic~~ hypothesis (H).

9.

Let K_α be the upper- α percentile of the $\chi^2(d)$ distribution, that is, the event $X \geq K_\alpha$ has probability α . Then the probability that the event

$S(\hat{\theta}) - S(\hat{\theta}) > \sigma^2 K_\alpha$ will occur when (H) is true is α . Therefore there is good evidence for rejecting hypothesis (H) if the above inequality holds for some pre-set value α (0.05 say).

Case σ^2 unknown

Use the property that

$$\frac{S(\hat{\theta}) - S(\hat{\theta})}{d} \bigg/ \frac{S(\hat{\theta})}{p-q} \sim F(d, p-q). \quad (17)$$

(Want ~~stat~~ something independent of σ^2).

As before, let K_α denote the upper α -percentile of the $F(d, p-q)$ distribution

The event

$$\left(\frac{S(\hat{\theta}) - S(\hat{\theta}^*)}{d} / \frac{S(\hat{\theta}^*)}{p-q} \right) > k_\alpha \quad (18)$$

has probability α if (H) is true
and there is good evidence for rejecting (H)
if $\alpha = 0.05$, say.

Procedure for Model Order Selection.

The above suggests a procedure for
model order selection.

Let $S_n = \min.$ value of least squares
criterion over vectors of parameters of dimension
 n , $n = 1, 2, \dots$

Assume $p \gg q \Rightarrow q/p$ small
 $\Rightarrow F(1, p-q) \approx \chi^2(1)$.

From the proposition, and the above, if
 n is a possible ~~best~~ model order

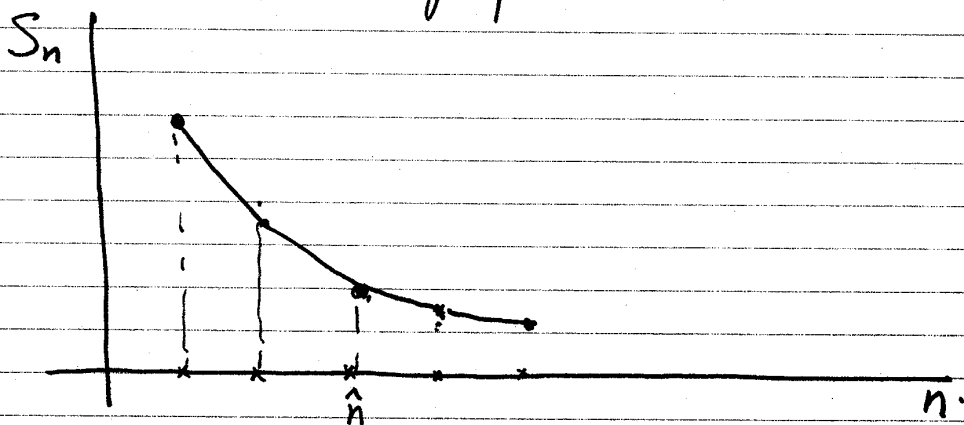
$$\frac{S_n - S_{n+1}}{\frac{S_{n+1}}{p}} \text{ has } \chi^2(1) \text{ distribution}$$

11.

The 0.01 percentile of $\chi^2(4) \approx 4$
(See Statistical tables). There are therefore
grounds then for rejecting n as a possible
model order at approx. 5% risk level
if the inequality

$$S_n - S_{n+1} > K \frac{S_{n+1}}{p} \quad \text{where } K = 4.$$

Consider now the graph



Above suggests, estimate \hat{n} of model order be
chosen to satisfy

$$S_{n-1} - S_n > K \frac{S_n}{p} \quad (19)$$

$$S_n - S_{n+1} \leq K \frac{S_{n+1}}{p}$$

for some pre-set $K = 4$ (say).

12.

Let us try to interpret these inequalities

View $S_n, n=1, \dots, p$ as a uniform

discretization of a C^1 -function $g: [0,1] \rightarrow \mathbb{R}$

i.e.

$$g\left(\frac{n}{p}\right) = S_n, \quad n=1, 2, \dots, p.$$

Now $\frac{S_n - S_{n-1}}{1/p}$ is a finite difference approx.

$$\text{to } -\frac{d}{dx} g(x) \text{ at } \hat{x} = \frac{\hat{n}}{p}.$$

Interpret (19) as

$$\left. -\frac{\frac{d}{dx} g(x)}{g(x)} - K \right|_{x = \frac{\hat{n}}{p}} = 0$$

$$\text{or } \left. \frac{d}{dx} [\log g(x) + Kx] \right|_{x = \frac{\hat{n}}{p}} = 0$$

This says that

$\log g(x) + Kx$ assumes a minimum at $x = \frac{\hat{n}}{p}$
(Necessary condition).

Interpret as :

$$\hat{n} \text{ minimizes } f(n) = \log S_n + \frac{K \cdot n}{p}$$

Now p is fixed and hence \hat{n} minimizes

$$A(n) = \log \frac{S_n}{p} + \tilde{K} \cdot n \quad \text{where } \tilde{K} = \frac{K}{p}$$

for \tilde{K} ~~set~~ at some pre-set level.

Accuracy of Estimates

Trustworthiness of estimated components $\hat{\theta}_i$ can be gauged from the α -confidence region for $\hat{\theta}_i$, which is an interval $I(\hat{\theta})$ and has the property that the event

$$\{\theta_i^* \in I(\hat{\theta})\} \text{ occurs with probability } \alpha$$

Case: σ^2 known

The estimate $\hat{\theta}$ is linear, unbiased and has covariance matrix $\sigma^2 (X^T X)^{-1}$ and $E \sim N(0, \sigma^2 I)$ it follows that

$$\hat{\theta} \sim N(\theta, \sigma^2 (X^T X)^{-1})$$

Hence

$$\hat{\theta}_i \sim N(\theta_i, \sigma^2 c_{ii}) \quad i=1, 2, \dots, q$$

where $(c_{ij}) = (X^T X)^{-1}$. Hence

$$\frac{(\hat{\theta}_i - \theta_i^*)}{\sqrt{\sigma^2 c_{ii}}} \sim N(0, 1)$$

Let k_β be the upper $\beta/2$ percentile for the distribution $N(0, 1)$. Then the confidence region can be constructed as follows:

$$(\hat{\theta}_i - k_\beta \sqrt{\sigma^2 c_{ii}}, \hat{\theta}_i + k_\beta \sqrt{\sigma^2 c_{ii}})$$

is a $(1-\beta)$ confidence region for θ_i^* , $i=1, \dots, q$.

Case 2 σ^2 Unknown

Construct regions from the statistic

$$\hat{\theta}_i / \sqrt{\hat{\sigma}^2 c_{ii}}, \text{ where}$$

$\hat{\sigma}^2$ is the unbiased estimate of σ^2 , given by

$$\hat{\sigma}^2 = (p-q)^{-1} (Y - X\hat{\theta})^T (Y - X\hat{\theta})$$

(Deduce this)

We now have to introduce another distribution. Given an integer $k > 0$, a random variable V is said to have the $t(k)$ -distribution if it can be expressed as

$$V = \frac{U}{\sqrt{W/k}}$$

where U and W are independent random variables with distributions $N(0,1)$ and $\chi^2(k)$ ~~distro~~ respectively.

Now for $i=1, 2, \dots, q$

$$\frac{\hat{\theta}_i - \theta_i^*}{\sqrt{C_{ii}}} / \sqrt{\hat{\sigma}^2} = \frac{\hat{\theta}_i - \theta_i^*}{\sqrt{\sigma^2 C_{ii}}} / \frac{\|Y - X\hat{\theta}\|^2}{\sqrt{(p-q)\sigma^2}}$$

$$\text{Now } \frac{\hat{\theta}_i - \theta_i^*}{\sqrt{\sigma^2 C_{ii}}} \sim N(0,1).$$

From the Proposition with $d=0$, we know

$\hat{\theta}_i$ is independent of $\|Y - X\hat{\theta}\|^2$ and

$$\frac{1}{\sigma^2} \|Y - X\hat{\theta}\|^2 \sim \chi^2(p-q).$$

$$\Rightarrow \frac{\hat{\theta}_i - \theta_i^*}{\sqrt{C_{ii}}} / \sqrt{\hat{\sigma}^2} \sim t(p-q).$$

Let k_β be the upper $\beta/2$ percentile for the distribution $t(p-q)$. Hence

$(\hat{\theta}_i - k_\beta \sqrt{\hat{\sigma}^2 C_{ii}}, \hat{\theta}_i + k_\beta \sqrt{\hat{\sigma}^2 C_{ii}})$ is a $(1-\beta)$ -confidence region for $\hat{\theta}_i$, $i=1, 2, \dots, q$.

II. Maximum likelihood Estimation

Predictor Model

$$Y_k = f_k(\theta; Y^{k-1}, u^{k-1}) + E_k, \quad k=1, 2, \dots, N \quad (1)$$

(E_1, \dots, E_N) has joint density function which depends on the parameter θ ; Assume

that u_0, \dots, u_{N-1} are independent of E_1, E_2, \dots, E_N

Let $p(y_N, \dots, y_1 | \theta, u^{N-1})$ be the joint density function of Y^N given θ and u^{N-1} .

ML-estimate is given by.

$$\underset{\theta}{\text{Max}} V_N(\theta, Y^N, u^{N-1}) \left[p(y_N, \dots, y_1 | \theta, u^{N-1}) \right]$$

Let $E_k \sim N(0, \Lambda_k), \Lambda_k > 0$

and E_k 's independent. Then a calculation using Bayes rule gives us

$$\log p(y_N, \dots, y_1 | u^{N-1}, \theta)$$

18.

$$\log p(y_N, \dots, y_1 | u^{N-1}, \theta)$$

$$= -\frac{Nr}{2} \log 2\pi - \frac{1}{2} \sum_{k=1}^N \log \det \Lambda_k - \frac{1}{2} \sum_{k=1}^N \|y_k - \hat{y}_k(\theta)\|_{\Lambda_k^{-1}}^2$$

Max. likelihood estimate criterion

$$\text{Max}_{\theta} L(\theta)$$

$$\text{where } L(\theta) = \sum_{k=1}^N \left\| \varepsilon_k(\theta) \right\|_{\Lambda_k^{-1}(\theta)}^2 + \sum_{k=1}^N \log \det \Lambda_k(\theta).$$