

TWO ERRORS IN STATISTICAL MODEL FITTING

NOBUO INAGAKI

(Received Apr. 26, 1976)

1. Introduction

Almost all methods of statistical estimation are carried out under the assumption that the number of parameters to be estimated is fixed and known although their values are unknown. However, in problems of statistical model fitting, it is one of the most important and difficult problems to suitably determine the number of unknown parameters. Generally, problems of statistical model fitting are treated as those of statistical test of goodness of fit and not discussed enough from the point of view of statistical estimation. That is, it is, first, tested whether or not the unknown parameters satisfy certain restrictions and next properties of estimation are discussed under either one of the null hypothesis and the alternative. (For example, Aitchison and Silvey [1]).

Recently Akaike [2] investigates problems of model fitting in the time series analysis from the point of view of statistical prediction and obtains a marvellous and reasonable criterion for deciding the numbers of unknown parameters. This criterion is based on the expectation of the Kullback-Leibler (K-L) information measure of an estimated density and the true density. Therefore we shall call it "Akaike's Information Criterion (AIC)". It is very interesting that the AIC statistic is equal to the C_p statistic due to Mallows [11] in the case of linear regression models. But it is noted that the latter is applied only in the case of linear regression models although the former is done in more general cases. Lindley [9] shows that the C_p statistic is obtained as a Bayes solution under a cost given by the number of explaining vectors to be utilized and a uniform prior distribution. But we should, indeed, discuss what cost to be adopted.

In the present paper we shall introduce two errors, one of which, K^M (say), is caused on account of "modelling" and the other, K^E (say), is done on account of "estimation". Those are based on K-L information measures and so, may be regarded to be equivalent. With this reason we shall define the risk of a pair of modelling and estimation by the sum of the above two errors: $R = K^M + K^E$ (say). (See Section

2.) In the normal linear regression case, the AIC statistic (which is equal to the C_p statistic, there,) is shown in Section 3 to be an estimator of the above risk of model of the number of parameters and estimation of their values, and furthermore a "ridge estimator" (see Hoerl and Kennard [3] and Lindley and Smith [11]) is, also, shown in Section 4 to be derived from the risk of model of the norm of parameters and estimation of their values. In Section 5 we shall propose a definition of statistical model fitting and its loss function, and obtain that the AIC statistic and a ridge estimator are Bayes solutions under the loss functions of the corresponding models and estimations stated in Sections 3 and 4, respectively. It is remarkable that this definition of statistical model fitting and its loss function are intended to enable us to test goodness of model fitting, and at the same time to estimate of parameters.

Asymptotic properties of the AIC statistic are discussed by Inagaki and Ogata [7]. The theory of the weak convergence of likelihood ratio random fields is very useful for this sake. (See LeCam [8], Ibragimov and Khas'minskii [4], [5], and Inagaki and Ogata [6], [7].)

2. Errors of modelling and estimation based on modelling

Let Θ be a parameter space and

$$(2.1) \quad \mathcal{F} = \{f(y, \theta) : \theta \in \Theta\}$$

be a family of probability density functions (p.d.f.'s) indexed by the elements of Θ . Let observation Y be distributed according to a p.d.f. $f(y, \theta_0) \in \mathcal{F}$ (that is, $\theta_0 \in \Theta$).

Let \mathcal{A} be an index space and C_α be a parameter space for each $\alpha \in \mathcal{A}$. Let

$$(2.2) \quad \mathcal{Q}_\alpha = \{g_\alpha(y, \zeta) : \zeta \in C_\alpha\}$$

be a family of p.d.f.'s for every $\alpha \in \mathcal{A}$ and let's call it a "model" for \mathcal{F} . Suppose that for each $\theta \in \Theta$ and $\alpha \in \mathcal{A}$ there exists an element $\zeta_\alpha(\theta) \in C_\alpha$ satisfying the following:

$$(2.3) \quad \int \log \{f(y, \theta)/g_\alpha(y, \zeta_\alpha(\theta))\} f(y, \theta) dy \\ = \min_{\zeta \in C_\alpha} \int \log \{f(y, \theta)/g_\alpha(y, \zeta)\} f(y, \theta) dy.$$

Then we define the error of model \mathcal{Q}_α for the true state $f(y, \theta)$, $K^\pi(\alpha|\theta)$ (say), by the Kullback-Leibler (K-L) information of $g_\alpha(y, \zeta_\alpha(\theta)) (\in \mathcal{Q}_\alpha)$ under $f(y, \theta) (\in \mathcal{F})$:

$$(2.4) \quad K^M(\alpha|\theta) = \int \log \{f(y, \theta)/g_\alpha(y, \zeta_\alpha(\theta))\} f(y, \theta) dy.$$

Hence it follows that the above definition of the error of model \mathcal{Q}_α , $K^M(\alpha|\theta)$, depends only on a subset $\tilde{\mathcal{Q}}_\alpha$ of \mathcal{Q}_α such that

$$(2.5) \quad \tilde{\mathcal{Q}}_\alpha = \{g_\alpha(y, \zeta_\alpha(\theta)) : \theta \in \Theta\} \quad (\text{say})$$

which has the common parameter space Θ for every $\alpha \in \mathcal{A}$.

Now, letting $T = T(Y)$ be any estimator for $\theta \in \Theta$, we shall call $\zeta_\alpha T = \zeta_\alpha(T(Y))$ (say) an estimator for $\zeta_\alpha(\theta)$ based on model \mathcal{Q}_α and define the error of $\zeta_\alpha T$ by the expectation of the K-L information of $g_\alpha(z, \zeta_\alpha(\theta))$, $\zeta_\alpha T(Y)$ under $g_\alpha(z, \zeta_\alpha(\theta))$:

$$(2.6) \quad K^E(T, \theta|\alpha) = E \left[\int \log \{g_\alpha(z, \zeta_\alpha(\theta))/g_\alpha(z, \zeta_\alpha T(Y))\} g_\alpha(z, \zeta_\alpha(\theta)) dz \right] \\ (\text{say}) \\ = \int \left[\int \log \{g_\alpha(z, \zeta_\alpha(\theta))/g_\alpha(z, \zeta_\alpha T(y))\} g_\alpha(z, \zeta_\alpha(\theta)) dz \right] \\ \times f(y, \theta) dy.$$

That is, $K^E(T, \theta|\alpha)$ does not mean a usual risk of estimator T of θ but a measurement of the effect of T based on modelling.

In the present paper we intend to approximate the family \mathcal{F} (which includes the true p.d.f. $f(y, \theta_0)$) by a model \mathcal{Q} (one of given models $\{\mathcal{Q}_\alpha\}$ $\alpha \in \mathcal{A}$) and at the same time, to estimate the true parameter θ_0 based on \mathcal{Q} . In general we may face the following contradiction in model fitting:

- (2.7) The larger is a fitting model \mathcal{Q} , the smaller is the error of modelling, K^M , but the larger is the error of estimation based on \mathcal{Q} , K^E , (because the vaguer is the information based on \mathcal{Q} with respect to parameter). The reverse is also true.

Note that two errors, K^M and K^E , can be regarded to be equivalent in the sense of the K-L information measure. Therefore we shall define the risk of model and estimator (\mathcal{Q}_α, T) (simply, denote it by (α, T))

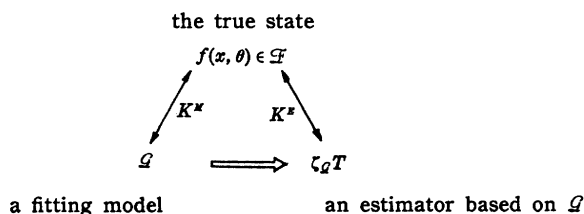


Fig. 2.1.

unless confused) by

$$(2.8) \quad R(\alpha, T|\theta) = K^u(\alpha|\theta) + K^E(T, \theta|\alpha) \quad (\text{say})$$

when θ is true. Thus, the problem in modelling and estimation based on modelling becomes to choose model and estimator (α, T) that minimize the risk $R(\alpha, T|\theta)$. Section 3 is related to AIC statistic and C_p statistic which are decision functions for dimension of parameter. Section 4 is related to ridge estimation which can be regarded as a decision function for norm of parameter.

3. Dimension and values of parameter

Let the parameter space Θ be a subset of the k -dimensional Euclidean space R^k . Suppose that we are given such a prior information that \mathcal{F} in (2.1) is known and the true parameter vector $\theta_0 = (\theta_0^{(1)}, \dots, \theta_0^{(k)})^t$ (say) satisfies

$$\theta_0 = (\theta_0^{(1)}, \dots, \theta_0^{(r_0)}, \theta_{00}^{(r_0+1)}, \dots, \theta_{00}^{(k)})^t$$

where $\theta_{00} = (\theta_{00}^{(1)}, \dots, \theta_{00}^{(k)})^t$ (say) is a given interior point of Θ but r_0 is unknown except that $1 \leq r_0 \leq k$. Without any loss of generality we may assume that

$$\theta_{00} = 0 = (0, \dots, 0)^t \in \Theta.$$

For $\theta = (\theta^{(1)}, \dots, \theta^{(r)}, \theta^{(r+1)}, \dots, \theta^{(k)})^t \in \Theta$, set

$$(3.1) \quad {}_r\theta = (\theta^{(1)}, \dots, \theta^{(r)}, \overbrace{0, \dots, 0}^{k-r})^t,$$

and assume

$$(3.2) \quad {}_r\theta \in \Theta, \quad 1 \leq r \leq k.$$

Then the prior information is restated as follows:

$$(3.3) \quad \begin{aligned} \theta_0 = {}_{r_0}\theta_0 &= (\theta_0^{(1)}, \dots, \theta_0^{(r_0)}, \overbrace{0, \dots, 0}^{k-r_0})^t, \\ \theta_0^{(r_0)} &\neq 0, \\ 1 &\leq r_0 \leq k, \quad \text{but unknown.} \end{aligned}$$

In this case, what matters is how to estimate r_0 and $\theta_0^{(1)}, \dots, \theta_0^{(r_0)}$, so called "dimension and values" of parameter θ_0 , respectively. Applying the concepts in Section 2, we shall, reasonably, make modelling in the following forms:

$$\begin{aligned}
 \mathcal{A} &= \{r \text{ integer: } 1 \leq r \leq k\}, \\
 (3.4) \quad C_r, \theta &= \{\zeta: \zeta \in \Theta\} \quad (\text{say}), \\
 \mathcal{Q}_r, \mathcal{F} &= \{f(y, \zeta): \zeta \in \Theta\} \quad (\text{say}).
 \end{aligned}$$

The error of model \mathcal{F} , $K^{\mathcal{M}}(r|\theta)$ becomes to be

$$\begin{aligned}
 (3.5) \quad K^{\mathcal{M}}(r|\theta) &= \int \log \{f(y, \theta)/f(y, \zeta_r(\theta))\} f(y, \theta) dy \\
 &= \min_{\zeta \in \Theta} \int \log \{f(y, \theta)/f(y, \zeta)\} f(y, \theta) dy.
 \end{aligned}$$

Note

$$(3.6) \quad {}_1\mathcal{F} \subset {}_2\mathcal{F} \subset \dots \subset {}_k\mathcal{F} = \mathcal{F}.$$

The error of estimator T based on \mathcal{F} is given by

$$\begin{aligned}
 (3.7) \quad K^E(T, \theta|r) &= \int \left[\int \log \{f(z, \zeta_r(\theta))/f(z, \zeta_r(T(y)))\} f(z, \zeta_r(\theta)) dz \right] \\
 &\quad \times f(y, \theta) dy.
 \end{aligned}$$

3.1. Normal linear regression case

Let

$$x_r = (x_{1r}, \dots, x_{nr})^t, \quad 1 \leq r \leq k$$

be explaining vectors and set

$$\begin{aligned}
 (3.8) \quad X_r &= (x_1, \dots, x_r), \quad n \times r\text{-matrix}, 1 \leq r \leq k, \\
 X &= X_k = (x_1, \dots, x_k).
 \end{aligned}$$

Suppose that observations $Y = (Y_1, \dots, Y_n)^t$ are distributed according to n -dimensional normal distribution $N_n(X\theta, \sigma^2 I)$ where $\sigma^2 > 0$ is known, $I = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$, the n -dimensional unit matrix, and $\theta = (\theta^{(1)}, \dots, \theta^{(k)})^t \in R^k$ is the unknown parameter. The problem of choosing the first r explaining vectors, x_1, \dots, x_r (but r is unknown) in the normal linear regression model may be regarded as a special case of modelling (3.4) with

$$(3.9) \quad f(y, \theta) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} (y - X\theta)^t (y - X\theta) \right\}$$

for $y = (y_1, \dots, y_n)^t$. Since

$$(3.10) \quad \int \log \{f(y, \theta)/f(y, \zeta)\} f(y, \theta) dy = \frac{1}{2\sigma^2} (X\theta - X \cdot \zeta)^t (X\theta - X \cdot \zeta),$$

we obtain

$$(3.11) \quad \zeta_r(\theta) = \tilde{P}_r \cdot X\theta \quad X \cdot \zeta_r(\theta) = P_r \cdot X\theta$$

where

$$(3.12) \quad \begin{aligned} P_r &= X_r(X_r'X_r)^-X_r', & n \times n\text{-matrix}, \\ \tilde{P}_r &= \begin{bmatrix} (X_r'X_r)^-X_r' \\ 0 \end{bmatrix}, & k \times n\text{-matrix}, \end{aligned}$$

and $(X_r'X_r)^-$ is a generalized inverse matrix of $X_r'X_r$. Then, P_r is the projection matrix from the n -dimensional Euclidean space R^n to the vector space generated by $\{x_1, \dots, x_r\}$, $\mathcal{V}(x_1, \dots, x_r)$ (say). Thus we have the following lemma.

LEMMA 3.1.

- (i) $\zeta_r(\theta) = \tilde{P}_r X\theta$ and $X\zeta_r(\theta) = P_r X\theta$.
- (ii) $K^u(r|\theta) = (1/2\sigma^2)(X\theta)'(P_k - P_r)(X\theta)$, and

$$(3.13) \quad \begin{aligned} K^E(T, \theta|r) &= E \frac{1}{2\sigma^2} (T - \theta)' X' P_r X (T - \theta) \\ &= \int \frac{1}{2\sigma^2} (\theta - T(y))' X' P_r X (\theta - T(y)) f(y, \theta) dy. \end{aligned}$$

Now, we obtain from (3.9) that the maximum likelihood estimator of θ , $\hat{\theta}$ (say), is given by

$$(3.14) \quad \hat{\theta} = \tilde{P}_k y,$$

and hence from (3.11) and (3.14) that

$$(3.15) \quad \zeta_r(\hat{\theta}) = \tilde{P}_r X\hat{\theta} = \tilde{P}_r P_k y = \tilde{P}_r y, \quad X\zeta_r(\hat{\theta}) = P_r y.$$

Note that $\zeta_r(\hat{\theta}) = \tilde{P}_r y$ is the maximum likelihood estimator of $\zeta_r(\theta)$. Further we have from (3.13) that

$$(3.16) \quad \begin{aligned} K^E(\hat{\theta}, \theta|r) &= E \left\{ \frac{1}{2\sigma^2} (X\theta - P_k y)' P_r (X\theta - P_k y) \right\} \\ &= E \left\{ \frac{1}{2\sigma^2} (X\theta - y)' P_r (X\theta - y) \right\} = \frac{1}{2} \text{trace } P_r \\ &= \frac{r}{2}, \quad \text{if rank } X_r = r. \end{aligned}$$

LEMMA 3.2. *The risk of model \mathcal{F} under the maximum likelihood estimation is as follows:*

$$(3.17) \quad R(r, \hat{\theta}|\theta) = \frac{1}{2\sigma^2} (X\theta)'(P_k - P_r)(X\theta) + \frac{1}{2} \text{trace } P_r \\ = \frac{1}{2} \left\{ \frac{1}{2\sigma^2} (X\theta)'(P_k - P_r)(X\theta) + r \right\}, \quad \text{if rank } X_r = r.$$

THEOREM 3.1. The C_p statistic due to Mallows [11],

$$(3.18) \quad C(r) = \frac{1}{\sigma^2} y'(P_k - P_r)y + 2 \text{trace } P_r - \text{trace } P_k \quad (\text{say})$$

is an unbiased estimator of $2R(r, \hat{\theta}|\theta)$.

Thus,

$$(3.19) \quad C(r^*) = \min \{C(r) : 1 \leq r \leq k\} \quad (\text{say})$$

can be regarded as an estimator of

$$(3.20) \quad 2R(r_1, \hat{\theta}|\theta) = \min \{2R(r, \hat{\theta}|\theta) : 1 \leq r \leq k\} \quad (\text{say}).$$

That is, we shall be able to consider that r^* is an estimator of r_1 (i.e. model $r_1\mathcal{F}$) which minimizes $R(r, \hat{\theta}|\theta)$, the risk of modelling under the maximum likelihood estimation. Set

$$(3.21) \quad v_r(X\theta) = \frac{1}{\sigma^2} (X\theta)'(P_r - P_{r-1})(X\theta), \quad 2 \leq r \leq k \\ v_1(X\theta) = \frac{1}{\sigma^2} (X\theta)'P_1(X\theta).$$

Then,

$$(3.22) \quad 2R(r, \hat{\theta}|\theta) = (1 - v_1(X\theta)) + \dots + (1 - v_r(X\theta)) + \frac{1}{\sigma^2} (X\theta)'(X\theta).$$

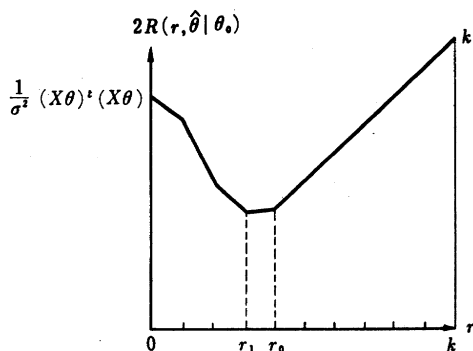


Fig. 2.2.

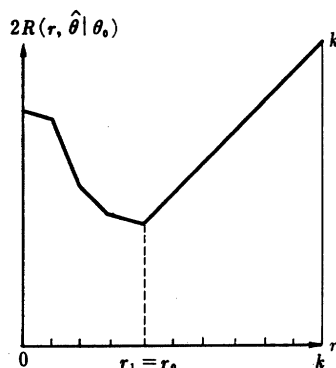


Fig. 2.3.

From the definition, (3.12), of projection P_r ,

$$(3.23) \quad v_r(X\theta_0) = 0, \quad \text{for } r \geq r_0 \text{ with } {}_{r_0}\theta_0 = \theta_0.$$

See Fig. 2.2. Further, if $v_r(X\theta_0) \leq 1$, for all r such that $2 \leq r \leq r_0$, we have $r_1 = r_0$ in (3.20). See Fig. 2.3.

3.2. Large sample (i.i.d.) case

Let \mathcal{F} be a family of n -product densities such that

$$f_n(y, \theta) = \prod_{i=1}^n f(y_i, \theta) \text{ (say)}, \quad y = (y_1, \dots, y_n)^t.$$

Then

$${}_r\mathcal{F} = \{f_n(y, {}_r\zeta) : \zeta \in \Theta\}.$$

Inagaki and Ogata [6] obtain the following results about the likelihood ratio random field and its related statistics.

LEMMA 3.3 (Theorem 3.4 in [6]). *The likelihood ratio random field*

$$(3.24) \quad h \rightsquigarrow Z_n(h) = f_n(y, \theta_0 + h/\sqrt{n}) / f_n(y, \theta_0) \quad \text{(say)}$$

converges weakly to a degenerated normal random field

$$(3.25) \quad h \rightsquigarrow Z(h) = \exp \left\{ h^t \Gamma^{1/2} \xi - \frac{1}{2} h^t \Gamma h \right\} \quad \text{(say)}$$

as $n \rightarrow \infty$, where $\Gamma = \Gamma(\theta_0)$ is the Fisher information matrix at θ_0 and ξ is a k -dimensional standard normal random vector $N_k(0, I)$.

This lemma suggests that asymptotic properties of the likelihood ratio, Z_n , and its related statistics are ascribable to properties of Z and functionals on it which are functions of a k -dimensional normal random vector ξ . Now, let

$$(3.26) \quad \begin{aligned} \Gamma^{1/2} &= (\gamma_1^{1/2}, \dots, \gamma_k^{1/2}), \\ {}_r\Gamma^{1/2} &= (\gamma_1^{1/2}, \dots, \gamma_r^{1/2}, \overbrace{0, \dots, 0}^{k-r}), \\ \mathcal{CV}({}_r\Gamma^{1/2}) &= \text{the vector space generated with vector } \gamma_1^{1/2}, \dots, \gamma_r^{1/2} \\ &\quad \text{of } {}_r\Gamma^{1/2}, \end{aligned}$$

$$P_r : R^k \rightarrow \mathcal{CV}({}_r\Gamma^{1/2}), \quad \text{projection.}$$

Since ${}_r(\theta_0 + h/\sqrt{n}) = \theta_0 + {}_r h/\sqrt{n}$ for $r \geq r_0$ with ${}_{r_0}\theta_0 = \theta_0$, we may define

$$(3.27) \quad {}_r\zeta(\theta_0 + h/\sqrt{n}) = \theta_0 + {}_r\zeta_n(h)/\sqrt{n} \quad \text{(say)}$$

for $r \geq r_0$ in (3.5):

$$(3.28) \quad K_n^M(r|h) = \int \log \{f_n(y, \theta_0 + h\sqrt{n})/f_n(y, \theta_0 + \tau\zeta_n(h)/\sqrt{n})\} \\ \times f_n(y, \theta_0 + h\sqrt{n}) dy \quad (\text{say}).$$

Denote the distribution of a random variable Z under P_r by $\mathcal{L}[Z|\theta]$.

LEMMA 3.4 (Lemma 5.5 in [6]). For $r \geq r_0$,

$$(3.29) \quad K_n^M(r|h) \rightarrow \frac{1}{2} h^t \Gamma^{1/2}(I - P_r) \Gamma^{1/2} h, \quad \tau\zeta_n(h) \rightarrow \tau\zeta(h),$$

as $n \rightarrow \infty$, where convergences are uniform on $|h| \leq M$ and $\tau\zeta(h)$ satisfies

$$(3.30) \quad \Gamma^{1/2} \tau\zeta(h) = P_r \Gamma^{1/2} h.$$

LEMMA 3.5 (see (i) of Theorem 2.1 and Application 4.3 in [6]). For maximum likelihood estimators $\hat{\theta}_{kn} = \hat{\theta}_{kn}(y)$ for θ and $\hat{\theta}_{rn} = \hat{\theta}_{rn}(y)$ for $\tau\zeta$ and under $r \geq r_0$

$$(3.31) \quad \mathcal{L}[\log \{f_n(y, \hat{\theta}_{kn})/f_n(y, \hat{\theta}_{rn})\} | \theta_0] \rightarrow \frac{1}{2} \xi^t (I - P_r) \xi,$$

and hence by the contiguity of $\{P_{n\theta_0}\}$ and $\{P_{n\theta_0+h/\sqrt{n}}\}$

$$(3.32) \quad \mathcal{L} \left[\log \{f_n(y, \hat{\theta}_{kn})/f_n(y, \hat{\theta}_{rn})\} \right. \\ \left. - \frac{1}{2} h^t \Gamma^{1/2} (I - P_r) \Gamma^{1/2} h | \theta_0 + h/\sqrt{n} \right] \rightarrow \frac{1}{2} \xi^t (I - P_r) \xi.$$

In the same way we can see by Lemma 5.5 in [6] that, for estimators $T_n = T_n(y)$ such that $\{\mathcal{L}(\sqrt{n}(T_n - \theta_0) | \theta_0)\}$ are relatively compact,

$$\int \log \{f_n(z, \theta_0 + \tau\zeta_n(h)/\sqrt{n})/f_n(z, \theta_0 + \tau\zeta_n(\sqrt{n}(T_n - \theta_0))/\sqrt{n})\} \\ \times f_n(z, \theta_0 + \tau\zeta_n(h)/\sqrt{n}) dz \\ - \frac{1}{2} (\sqrt{n}(T_n - \theta_0) - h)^t \Gamma^{1/2} P_r \Gamma^{1/2} (\sqrt{n}(T_n - \theta_0) - h) \rightarrow 0$$

in P_{θ_0} . Hence we have the following lemma for the maximum likelihood estimator $\hat{\theta}_n = \hat{\theta}_{kn}(y)$ (say) of θ_0 .

LEMMA 3.6. For $r \geq r_0$,

$$(3.33) \quad K_n^E(\hat{\theta}_n, \theta_0 + h/\sqrt{n} | r) \\ = \int \left[\int \log \{f_n(z, \theta_0 + \tau\zeta_n(h)/\sqrt{n}) \right. \\ \left. / f_n(z, \theta_0 + \tau\zeta_n(\sqrt{n}(\hat{\theta}_n(y) - \theta_0))/\sqrt{n}) \} \right]$$

$$\times f_n(z, \theta_0 + \zeta_n(h)/\sqrt{n}) dz \Big] f_n(y, \theta_0 + h/\sqrt{n}) dy \quad (\text{say})$$

$$\rightarrow r/2.$$

After all, we conclude that the AIC statistic due to Akaike [2],

$$(3.34) \quad \text{AIC}_n(r) = 2 \log \{f_n(y, \hat{\theta}_{kn}(y)) / f_n(y, \hat{\theta}_{rn}(y))\} + 2r - k \quad (\text{say})$$

is an estimator of

$$(3.35) \quad 2R_n(r, \hat{\theta}_n|h) = 2\{K_n^M(r|h) + K_n^E(\hat{\theta}_n, \theta_0 + h/\sqrt{n} | r)\} \quad (\text{say})$$

where $R_n(r, \hat{\theta}_n|h)$ is the risk of model \mathcal{F} under the maximum likelihood estimation and further, have the following.

THEOREM 3.2 (See 4.3 of Section 4 in [6]).

$$(3.36) \quad \text{AIC}_n(r) \rightarrow C(r) = \xi'(I - P_r)\xi + 2r - k, \quad \text{as } n \rightarrow \infty \text{ for } r \geq r_0$$

where $C(r)$ is the Mallows's C_p statistic in (3.18), and for $1 \leq r < r_0$,

$$(3.37) \quad \text{AIC}_n(r) \rightarrow \infty, \quad \text{as } n \rightarrow \infty.$$

Note that $R_n(r, \hat{\theta}_n|h) \rightarrow \infty$ for $1 \leq r < r_0$. Therefore, for all sufficiently large n , r_n^* such that

$$(3.38) \quad \text{AIC}_n(r_n^*) = \min \{\text{AIC}_n(r) : 1 \leq r \leq k\}$$

can be regarded as an estimator r_0 in (3.3) which minimizes $R_n(r, \hat{\theta}_n|h)$ for all sufficiently large n (in the same way as in (3.23)).

It is very interesting that the AIC is exactly equal to the C_p statistic in the case of linear regression models. But it is remarkable that the latter applied only in the case of linear regression models although the former is done in more general cases.

4. Norm and values of parameter

We return to the normal linear regression case discussed in 3.1. For symmetric matrix $X'X$, there is an orthogonal matrix S satisfying

$$(4.1) \quad S'(X'X)S = A = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k \end{pmatrix} \quad (\text{say})$$

where $\lambda_1, \dots, \lambda_k \geq 0$. Assume

$$(4.2) \quad \lambda_1, \dots, \lambda_k > 0, \quad \text{i.e. rank } X = k.$$

We shall consider another modelling in the normal linear regression case

such that

$$\begin{aligned}
 \mathcal{A} &= R^+ = \{\alpha: \text{real and positive}\}, \\
 (4.3) \quad \mathcal{C}_\alpha &= \Theta_\alpha = \{\zeta \in \Theta: |\zeta| \leq \alpha\}, \quad (\text{say}) \\
 \mathcal{G}_\alpha &= \mathcal{F}_\alpha = \{f(y, \zeta): |\zeta| \leq \alpha\},
 \end{aligned}$$

where $|\zeta|$ is a norm of ζ so that

$$\begin{aligned}
 |\zeta|_0 &= \max \{|\zeta^{(1)}|, \dots, |\zeta^{(k)}|\}, \\
 |\zeta|_1^2 &= \zeta' \zeta, \\
 (4.4) \quad |\zeta|_X^2 &= \zeta' (X' X) \zeta, \\
 |\zeta|_A^2 &= \zeta' A \zeta, \quad (\text{say})
 \end{aligned}$$

for a positive definite matrix A . However, we shall discuss the case of Euclidean norm $|\zeta|_1$ in the present paper. Remark

$$(4.5) \quad \mathcal{F}_{\alpha_1} \subset \mathcal{F}_{\alpha_2}, \quad \text{if } \alpha_1 < \alpha_2.$$

By (3.10), $\zeta_\alpha(\theta)$ (defined in (2.3)) is given as the solution of ζ that minimizes

$$(4.6) \quad (\zeta - \theta)' X' X (\zeta - \theta), \quad \text{subject to } \zeta' \zeta \leq \alpha^2.$$

Set

$$(4.7) \quad F(\zeta, \nu) = (\zeta - \theta)' X' X (\zeta - \theta) + \nu(\zeta' \zeta - \alpha^2).$$

Kuhn-Tucker condition is as follows: there exist $\nu^0 \geq 0$ and ζ^0 such that

$$(4.8) \quad \left(\frac{\partial F}{\partial \zeta} \right)_{\zeta^0} = 0, \quad \left(\frac{\partial F}{\partial \nu} \right)_{\nu^0} \leq 0, \quad \nu^0 \left(\frac{\partial F}{\partial \nu} \right)_{\nu^0} = 0.$$

That is,

$$\begin{aligned}
 (\zeta^0 - \theta)' X' X + \nu^0 \zeta^{0t} &= 0, \\
 (4.9) \quad \zeta^{0t} \zeta^0 - \alpha^2 &= 0, \quad \text{if } |\theta| > \alpha, \\
 \nu^0 &= 0, \quad \text{if } |\theta| \leq \alpha.
 \end{aligned}$$

Now, we define mappings as follows:

$$\begin{aligned}
 (4.10) \quad \tilde{Q}_\alpha &= (X' X + \nu_\alpha I)^{-1} X', \quad k \times n\text{-matrix}, \\
 Q_\alpha &= X \tilde{Q}_\alpha = X (X' X + \nu_\alpha I)^{-1} X', \quad n \times n\text{-matrix}
 \end{aligned}$$

where $\nu_\alpha = \nu(\alpha, \theta)$ is a positive number satisfying

$$(4.11) \quad \begin{aligned} (X\theta)' \tilde{Q}_\alpha \tilde{Q}_\alpha (X\theta) &= \alpha^2 & \text{if } |\theta| > \alpha, \\ \nu_\alpha &= 0 & \text{if } |\theta| \leq \alpha, \end{aligned}$$

for $\alpha > 0$ and $\theta \in \Theta$. Hence, we see from (4.6) and (4.9)–(4.11) that

$$(4.12) \quad \zeta_\alpha(\theta) = \tilde{Q}_\alpha X \theta.$$

Given and fixed θ in (4.11). According to

$$\alpha^2 = \theta' (\tilde{Q}_\alpha X)' (\tilde{Q}_\alpha X) \theta = \sum_{i=1}^k \frac{\beta_i^2 \lambda_i^2}{(\lambda_i + \nu)^2}$$

where

$$(4.13) \quad \beta = (\beta_1, \dots, \beta_k)' = S' \theta, \quad \text{say, (see (4.1)).}$$

$\nu = \nu_\alpha = \nu(\alpha, \theta)$ has such a derivative that

$$(4.14) \quad \frac{d\nu}{d\alpha} = -\alpha / \left\{ \sum_{i=1}^k \frac{\beta_i^2 \lambda_i^2}{(\lambda_i + \nu)^3} \right\} < 0, \quad \text{for } \nu \geq 0.$$

That is, $\nu = \nu_\alpha = \nu(\alpha, \theta)$ is a decreasing function of α with

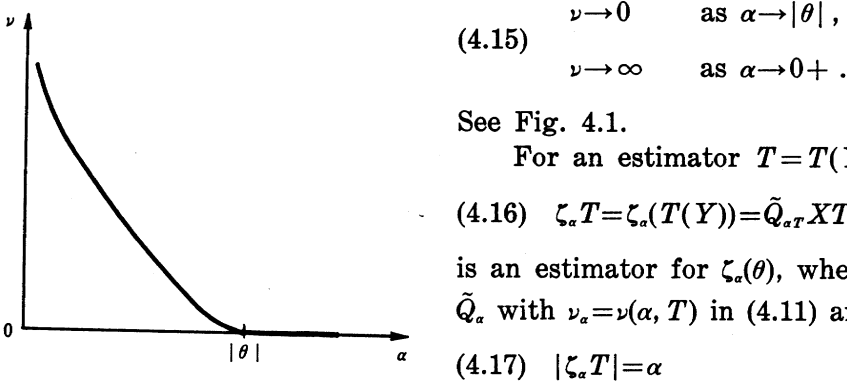


Fig. 4.1.

$$(4.15) \quad \begin{aligned} \nu &\rightarrow 0 & \text{as } \alpha \rightarrow |\theta|, \\ \nu &\rightarrow \infty & \text{as } \alpha \rightarrow 0+. \end{aligned}$$

See Fig. 4.1.

For an estimator $T = T(Y)$ for θ ,

$$(4.16) \quad \zeta_\alpha T = \zeta_\alpha(T(Y)) = \tilde{Q}_{\alpha T} X T(Y)$$

is an estimator for $\zeta_\alpha(\theta)$, where $\tilde{Q}_{\alpha T}$ is \tilde{Q}_α with $\nu_\alpha = \nu(\alpha, T)$ in (4.11) and hence

$$(4.17) \quad |\zeta_\alpha T| = \alpha \quad \text{if } |T| > \alpha \text{ (from (4.11)).}$$

See Fig. 4.2. We shall introduce another estimator for $\zeta_\alpha(\theta)$,

$$(4.18) \quad \tilde{\zeta}_\alpha T = \tilde{Q}_\alpha X T(Y), \quad \text{with } \nu_\alpha = \nu(\alpha, \theta),$$

which is so called a “ridge estimator” for θ (see Hoerl and Kennard [3]). See Fig. 4.3.

The following lemma is easily seen from (3.10) and (4.12).

LEMMA 4.1.

(i) The error of model \mathcal{F}_α is

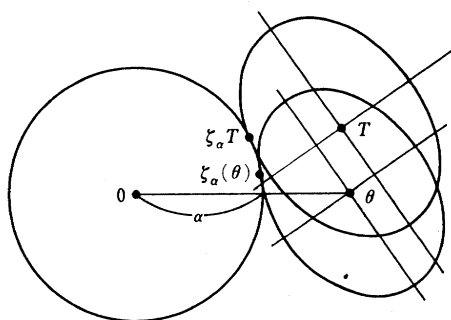
ellipse $\theta \propto$ ellipse T

Fig. 4.2.

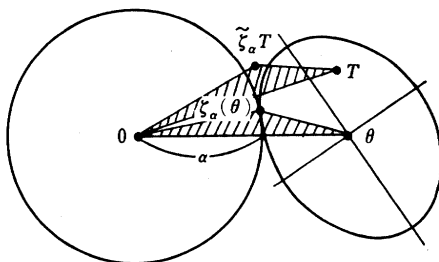


Fig. 4.3.

$$(4.19) \quad K^M(\alpha|\theta) = \frac{1}{2\sigma^2} (X\theta)'(P_k - Q_\alpha)^2(X\theta) .$$

(ii) The errors of estimators $\zeta_a T$ and $\tilde{\zeta}_a T$ are

$$(4.20) \quad K^E(T, \theta|\alpha) = E \left\{ \frac{1}{2\sigma^2} (Q_\alpha X\theta - Q_{\alpha T} XT)'(Q_\alpha X\theta - Q_{\alpha T} XT) \right\} ,$$

$$\tilde{K}^E(T, \theta|\alpha) = E \left\{ \frac{1}{2\sigma^2} (X\theta - XT)'Q_\alpha^2(X\theta - XT) \right\} ,$$

respectively.

In the sequel, we shall treat only the ridge case.

LEMMA 4.2. For the maximum likelihood estimator $\hat{\theta} = \tilde{P}_k Y$ in (3.11), the error of $\tilde{\zeta}_a \hat{\theta}$ is

$$(4.21) \quad \tilde{K}^E(\hat{\theta}, \theta|\alpha) = \frac{1}{2} \text{trace } Q_\alpha^2 = \frac{1}{2} \sum_{i=1}^k \left(\frac{\lambda_i}{\lambda_i + \nu_\alpha} \right)^2 .$$

Further, the risk of model α under the maximum likelihood estimation, $\tilde{R}(\alpha, \hat{\theta}|\theta)$ (say), is

$$(4.22) \quad \begin{aligned} \tilde{R}(\alpha, \hat{\theta}|\theta) &= \frac{1}{2\sigma^2} (X\theta)'(P_k - Q_\alpha)^2(X\theta) + \frac{1}{2} \text{trace } Q_\alpha^2 \\ &= \frac{1}{2} \left\{ \frac{\nu_\alpha^2}{\sigma^2} \sum_{i=1}^k \frac{\beta_i^2 \lambda_i}{(\lambda_i + \nu_\alpha)^2} + \sum_{i=1}^k \frac{\lambda_i^2}{(\lambda_i + \nu_\alpha)^2} \right\} \\ &= E \left\{ \frac{1}{2\sigma^2} (Q_\alpha Y - X\theta)'(Q_\alpha Y - X\theta) \right\} \end{aligned}$$

for $\beta = (\beta_1, \dots, \beta_k)' = S'\theta$ in (4.13).

PROOF. By the virtue of

$$(4.23) \quad \tilde{Q}_\alpha P_k = (X'X + \nu_\alpha I)^{-1} X' \cdot X (X'X)^{-1} X' = \tilde{Q}_\alpha,$$

we have

$$\begin{aligned} \tilde{K}^E(\hat{\theta}, \theta | \alpha) &= E \left\{ \frac{1}{2\sigma^2} (X\theta - P_k y)' Q_\alpha^2 (X\theta - P_k y) \right\} \\ &= \frac{1}{2} \text{trace } Q_\alpha^2 P_k = \frac{1}{2} \text{trace } Q_\alpha^2 \\ &= \frac{1}{2} \text{trace } \{(X'X + \nu_\alpha I)^{-1} X'X\}^2 \end{aligned}$$

and so from (4.1)

$$= \frac{1}{2} \sum_{i=1}^k \left(\frac{\lambda_i}{\lambda_i + \nu_\alpha^2} \right)^2.$$

It follows from (4.19) and (4.21) that

$$\tilde{R}(\alpha, \hat{\theta} | \theta) = \frac{1}{2} \left\{ \frac{1}{\sigma^2} \theta' (X - Q_\alpha X)' (X - Q_\alpha X) \theta + \text{trace } Q_\alpha^2 \right\}.$$

Now, from (4.10)

$$\begin{aligned} X - Q_\alpha X &= X - X(X'X + \nu_\alpha I)^{-1} (X'X + \nu_\alpha I - \nu_\alpha I) \\ &= \nu_\alpha X (X'X + \nu_\alpha I)^{-1} = \nu_\alpha \tilde{Q}^t \end{aligned}$$

and hence

$$(X - Q_\alpha X)' (X - Q_\alpha X) = \nu_\alpha^2 (X'X + \nu_\alpha I)^{-1} X'X (X'X + \nu_\alpha I)^{-1}.$$

This leads to (4.22) by using (4.13).

The following lemma is an immediate result of the last lemma.

LEMMA 4.3. $\tilde{R}(\nu) = \tilde{R}(\alpha, \hat{\theta} | \theta)$ with $\nu = \nu_\alpha$ is a function of ν as follows: for a given and fixed θ ,

$$\begin{aligned} \tilde{R}(\nu) &= \frac{1}{2\sigma^2} \left\{ \nu^2 \sum_{i=1}^k \frac{\beta_i^2 \lambda_i}{(\lambda_i + \nu)^2} + \sigma^2 \sum_{i=1}^k \frac{\lambda_i^2}{(\lambda_i + \nu)^2} \right\}, \\ \tilde{R}(0+) &= \frac{k}{2}, \quad \tilde{R}(+\infty) = \sum_{i=1}^k \beta_i^2 \lambda_i = |X\theta|^2 \\ (4.24) \quad \frac{\partial}{\partial \nu} \tilde{R}(\nu) &= \frac{\nu}{\sigma^2} \sum_{i=1}^k \frac{\beta_i^2 \lambda_i^2}{(\lambda_i + \nu)^3} - \sum_{i=1}^k \frac{\lambda_i^2}{(\lambda_i + \nu)^3}, \\ \frac{\partial}{\partial \nu} \tilde{R}(\nu) &< 0 \quad \text{if } \nu < \sigma^2 / \beta_{\max}^2, \quad \frac{\partial}{\partial \nu} \tilde{R}(\nu) > 0 \quad \text{if } \nu > \sigma^2 / \beta_{\min}^2, \end{aligned}$$

where

$$(4.25) \quad \beta_{\max}^2 = \max \{\beta_1^2, \dots, \beta_k^2\}, \quad \beta_{\min}^2 = \min \{\beta_1^2, \dots, \beta_k^2\}, \quad (\text{say}).$$

On the other hand, ridge estimator $\tilde{\zeta}_a \hat{\theta} = \tilde{Q}_a Y$ (recall (4.23)) has the following properties (see Hoerl and Kennard [3]): for a given and fixed θ ,

$$\begin{aligned} L(\nu) &= L(\nu_a) = E(\tilde{Q}_a Y - \theta)'(\tilde{Q}_a Y - \theta) \quad (\text{say}) \\ &= \nu^2 \sum_{i=1}^k \frac{\beta_i^2}{(\lambda_i + \nu)^2} + \sigma^2 \sum_{i=1}^k \frac{\lambda_i}{(\lambda_i + \nu)^2}, \\ L(0+) &= \sigma^2 \sum_{i=1}^k \frac{1}{\lambda_i}, \quad L(+\infty) = \sum_{i=1}^k \beta_i^2 = |\theta|^2 \\ (4.26) \quad \frac{\partial}{\partial \nu} L(\nu) &= 2 \left\{ \sum_{i=1}^k \frac{\beta_i^2 \lambda_i}{(\lambda_i + \nu)^3} - \sigma^2 \sum_{i=1}^k \frac{\lambda_i}{(\lambda_i + \nu)^3} \right\}, \\ \frac{\partial}{\partial \nu} L(\nu) &< 0 \quad \text{if } \nu < \sigma^2 / \beta_{\max}^2, \quad \frac{\partial}{\partial \nu} L(\nu) > 0 \quad \text{if } \nu > \sigma^2 / \beta_{\min}^2. \end{aligned}$$

Fig. 4.5 is due to Hoerl and Kennard [3] (Fig. 1).

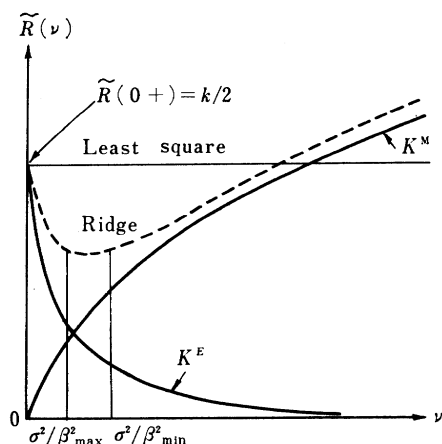


Fig. 4.4.

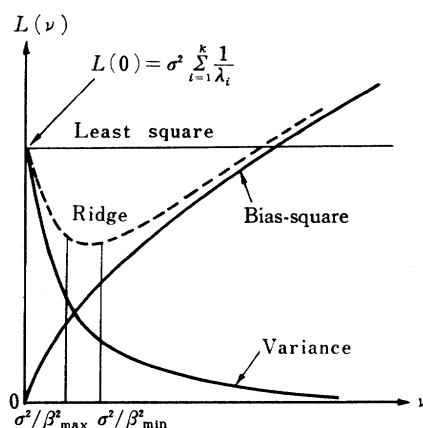


Fig. 4.5.

THEOREM 4.1. *There always exists a positive number $\nu > 0$ such that*

$$\tilde{R}(\nu) < \tilde{R}(0+) = k/2.$$

Further, if $\nu < \sigma^2 / \beta_{\max}^2$,

$$(4.27) \quad \tilde{R}(\nu) < \tilde{R}(0+) \quad \text{and} \quad L(\nu) < L(0+).$$

Now, we shall estimate $K^{\#}(\nu) = K^{\#}(\nu_a) = K^{\#}(\alpha, \hat{\theta} | \theta)$ by an unbiased estimator,

$$\frac{1}{2\sigma^2} Y'(P_k - Q)^2 Y - \frac{1}{2} \text{trace} (P_k - Q)^2.$$

Then, we see

$$(4.28) \quad B_1(\nu) = \frac{1}{\sigma^2} Y'(P_k - Q_\alpha)^2 Y + 2 \text{trace } Q_\alpha - k \quad (\text{say})$$

is an unbiased estimator of $2\tilde{R}(\nu)$. The following lemma is proved similarly as Lemma 4.3.

LEMMA 4.4.

$$(4.29) \quad \begin{aligned} B_1(\nu) &= \frac{\nu^2}{\sigma^2} \sum_{i=1}^k \frac{z_i^2 \lambda_i}{(\lambda_i + \nu)^2} + 2 \sum_{i=1}^k \frac{\lambda_i}{\lambda_i + \nu} - k \\ B_1(0+) &= k, \quad B_1(+\infty) = \frac{1}{\sigma^2} \sum_{i=1}^k z_i^2 \lambda_i - k = \frac{1}{\sigma^2} |P_k y|^2 - k \\ \frac{\partial}{\partial \nu} B_1(\nu) &= 2 \left\{ \frac{\nu}{\sigma^2} \sum_{i=1}^k \frac{z_i^2 \lambda_i^2}{(\lambda_i + \nu)^3} - \sum_{i=1}^k \frac{\lambda_i}{(\lambda_i + \nu)^2} \right\} \\ \frac{\partial}{\partial \nu} B_1(\nu) &< 0 \quad \text{if } \nu < \sigma^2 / (z^2 \lambda)_{\max}, \quad \frac{\partial}{\partial \nu} B_1(\nu) > 0 \quad \text{if } \nu > \sigma^2 / (z^2 \lambda)_{\min} \end{aligned}$$

where

$$(4.30) \quad \begin{aligned} z &= (z_1, \dots, z_k)' = S' \tilde{P}_k Y, \\ (z^2 \lambda)_{\max} &= \max \{z_1^2 \lambda_1, \dots, z_k^2 \lambda_k\}, \quad (z^2 \lambda)_{\min} = \min \{z_1^2 \lambda_1, \dots, z_k^2 \lambda_k\}. \end{aligned}$$

Thus, we can conclude that ν_1^* and ν_0 , $0 < \nu^*$, $\nu_0 < \infty$, exist which satisfy

$$(4.31) \quad \begin{aligned} B_1(\nu_1^*) &= \min \{B(\nu): 0 < \nu < \infty\}, \\ \tilde{R}(\nu_0) &= \min \{R(\nu): 0 < \nu < \infty\}, \end{aligned}$$

and that model α_1^* obtained by ν_1^* and the equation

$$(4.32) \quad (P_k Y)' \tilde{Q}_\alpha \tilde{Q}_\alpha (P_k Y) = \alpha^2$$

is an estimator α_0 obtained by ν_0 and (4.11). After all, we have a "ridge estimator"

$$\tilde{\zeta}_{\alpha_1^*} \hat{\theta} = \tilde{Q}_{\alpha_1^*} Y = (X'X + \nu_1^* I)^{-1} X'Y$$

for θ .

5. Statistical model fitting

In this section we shall define the concept of "Statistical model fitting" and its risk. Further, we shall show that the AIC estimator and the ridge estimator are obtained as Bayes solutions of the above risks under uniform prior distribution, respectively.

Denote an estimator for $\alpha \in \mathcal{A}$ by

$$(5.1) \quad \tau = \tau(Y) = \sum_{\alpha \in \mathcal{A}} \alpha \chi_{H_\alpha}(Y)$$

where for each α

$$(5.2) \quad \begin{aligned} H_\alpha &= \{y; \tau(y) = \alpha\}, \text{ the contour of } \tau, \\ \chi_{H_\alpha} &= \text{the indicator function of set } H_\alpha. \end{aligned}$$

That is, τ is a multiple decision function for models \mathcal{G}_α , $\alpha \in \mathcal{A}$ in (2.2). Define a loss function of estimator τ by

$$(5.3) \quad J(\tau(y), \theta) = \log \{f(y, \theta) / g_{\tau(y)}(y, \zeta_{\tau(y)}(\theta))\},$$

using ζ_α in (2.3). It is remarkable that $J(\tau(y), \theta)$ is not always non-negative, but it is decomposed into the sum of the following two parts, J_0 and J_1 (say), J_0 being common to all $\alpha \in \mathcal{A}$ (and so independent of τ), and J_1 nonnegative:

$$(5.4) \quad \begin{aligned} J_0(\theta) &= \log \{f(y, \theta) / g^0(y, \theta)\} \\ J_1(\tau(y), \theta) &= \log \{g^0(y, \theta) / g_{\tau(y)}(y, \zeta_{\tau(y)}(\theta))\} \end{aligned}$$

where

$$(5.5) \quad g^0(y, \theta) = \sup_{\alpha \in \mathcal{A}} g_\alpha(y, \zeta_\alpha(\theta)).$$

We call a pair of estimators for model and parameter, (τ, T) , a "statistical model fitting" and define a loss and a risk of model fitting (τ, T) by

$$(5.6) \quad \begin{aligned} W(\tau(y), T(y) | \theta) &= J(\tau(y), \theta) + \int \log \{f(z, \zeta_{\tau(y)}(\theta)) / f(z, \zeta_{\tau(y)} T(y))\} \\ &\quad \times f(z, \zeta_{\tau(y)}(\theta)) dz. \end{aligned}$$

and

$$(5.7) \quad \mathcal{P}(\tau, T | \theta) = E W(\tau(Y), T(Y) | \theta) = \int W(\tau(y), T(y) | \theta) f(y, \theta) dy,$$

respectively.

5.1. AIC statistic

Under the same situations as in 3.1 (normal linear regression), we shall consider a model fitting for dimension and values of parameter. Then, an estimator for r , $1 \leq r \leq k$, is

$$(5.8) \quad \tau = \tau(Y) = \sum_{r=1}^k r \chi_{H_r}(Y) .$$

It is easy to see by (3.9), (3.11) and (5.3) that

$$(5.9) \quad J(\tau(y), \theta) = \frac{1}{2\sigma^2} (y - P_{\tau(y)} X \theta)' (y - P_{\tau(y)} X \theta) - \frac{1}{2\sigma^2} (y - X \theta)' (y - X \theta) .$$

Therefore, from (3.13) and (5.6) through (5.9) we have the following lemma.

LEMMA 5.1. *The loss and the risk of a model fitting (τ, T) are given by*

$$(5.10) \quad W(\tau, T | \theta) = \frac{1}{2\sigma^2} (y - P_{\tau} X \theta)' (y - P_{\tau} X \theta) - \frac{1}{2\sigma^2} (y - X \theta)' (y - X \theta) \\ + \frac{1}{2\sigma^2} (XT - X \theta)' P_{\tau} (XT - X \theta) ,$$

and

$$(5.11) \quad \mathcal{P}(\tau, Y | \theta) = \int \frac{1}{2\sigma^2} \{ (y - P_{\tau(y)} X \theta)' (y - P_{\tau(y)} X \theta) \\ + (XT(y) - X \theta)' P_{\tau(y)} (XT(y) - X \theta) \} f(y, \theta) dy - n/2 ,$$

respectively.

Now, let's take the Lebesgue measure on R^k as a prior distribution of θ . Then the posterior distribution, $\Psi(\theta | y)$ (say), becomes to be the k -dimensional normal distribution $N_k(\tilde{P}_k y, \sigma^2(X'X)^{-})$, that is,

$$(5.12) \quad \Psi(\theta | y) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^k |X'X|^{1/2} \exp \left\{ -\frac{1}{2\sigma^2} (\theta - \tilde{P}_k y)' (X'X) (\theta - \tilde{P}_k y) \right\} ,$$

if $\text{rank } X = k$. Hence, the posterior risk is calculated as follows:

$$(5.13) \quad \rho(\tau, T | y) = \int W(\tau, T | \theta) \Psi(\theta | y) d\theta \quad (\text{say}) \\ = \frac{1}{2\sigma^2} y' (P_{\tau} - P_{\tau}) y + \frac{1}{2\sigma^2} \int \{ (\theta - \tilde{P}_k y)' X' P_{\tau} X (\theta - \tilde{P}_k y) \\ - (\theta - \tilde{P}_k y)' X' X (\theta - \tilde{P}_k y) + (\theta - T)' X' P_{\tau} X (\theta - T) \} \\ \times \Psi(\theta | y) d\theta$$

$$\begin{aligned}
&= \frac{1}{2\sigma^2} y'(P_k - P_\tau)y + \frac{1}{2} \text{trace } P_\tau - \frac{1}{2} \text{trace } P_k \\
&\quad + \frac{1}{2\sigma^2} (T - \tilde{P}_k y)' X' P_\tau X (T - \tilde{P}_k y) + \frac{1}{2} \text{trace } P_\tau \\
&= \frac{1}{2} \left\{ C(\tau) + \frac{1}{\sigma^2} (T - \tilde{P}_k y)' X' P_\tau X (T - \tilde{P}_k y) \right\} \quad (\text{recall (3.18)}).
\end{aligned}$$

Note, again, that the AIC statistic is the same as the C_p statistic in the normal linear regression case. We shall call a model fitting $(r^*, \hat{\theta})$ with r^* in (3.19) and $\hat{\theta}$ in (3.14), the "AIC model fitting". (5.13) leads to the following theorem.

THEOREM 5.1. *The AIC model fitting $(r^*, \hat{\theta})$ is obtained as the (generalized) Bayes solution when the loss function is $W(\tau, T|\theta)$ in (5.10) and the prior distribution is the Lebesgue measure on R^k .*

In the i.i.d. large sample case, Inagaki and Ogata [6] obtain the following theorem.

THEOREM 5.2 (Theorem 5.2 in [6]). *The AIC model fitting (r_n^*, θ_n^*) with r_n^* defined in (3.38) and $\theta^* = \theta_{r_n^*, n}$ (see (3.31)), is asymptotically equivalent to a Bayes solution under a uniform prior distribution.*

5.2. Ridge estimators

We shall discuss a model fitting for norm and values of parameter under the same situations as in Section 4. Let

$$(5.14) \quad \tau = \tau(Y) = \sum_{\alpha \in \mathcal{A}} \alpha \chi_{H_\alpha}(Y)$$

be an estimator for α , $0 < \alpha < \infty$. Then it follows from (3.9), (4.12) and (5.3) that

$$(5.15) \quad J(\tau(y), \theta) = \frac{1}{2\sigma^2} (y - Q_{\tau(y)} X \theta)' (y - Q_{\tau(y)} X \theta) - \frac{1}{2\sigma^2} (y - X \theta)' (y - X \theta).$$

When not $\zeta_\alpha T$ but $\tilde{\zeta}_\alpha T$ is utilized in Lemma 4.1, we have the following lemma from (4.20), (5.6), (5.7) and (5.15).

LEMMA 5.2. *The loss and the risk of a model fitting (τ, T) are given by*

$$\begin{aligned}
(5.16) \quad \tilde{W}(\tau, T|\theta) &= \frac{1}{2\sigma^2} (y - Q_\tau X \theta)' (y - Q_\tau X \theta) - \frac{1}{2\sigma^2} (y - X \theta)' (y - X \theta) \\
&\quad + \frac{1}{2\sigma^2} (XT - X \theta)' Q_\tau^2 (XT - X \theta),
\end{aligned}$$

and

$$(5.17) \quad \tilde{\mathcal{P}}(\tau, T|\theta) = \int \frac{1}{2\sigma^2} \{ (y - Q_{\tau(y)} X\theta)^t (y - Q_{\tau(y)} X\theta) \\ + (XT(y) - X\theta)^t Q_{\tau(y)}^2 (XT(y) - X\theta) \} f(y, \theta) dy - n/2 ,$$

respectively.

Let's take the Lebesgue measure on R^* as a prior distribution, again. As in (5.13), we see (recalling (4.23))

$$(5.18) \quad \tilde{\rho}(\tau, T|y) = \int \tilde{W}(\tau, T|\theta) \Psi(\theta|y) d\theta \quad (\text{say}) \\ = \frac{1}{2\sigma^2} \int \{ (\theta - \tilde{P}_k)^t X^t Q_{\tau}^2 X (\theta - \tilde{P}_k y) - 2y^t Q_{\tau} X\theta \\ + 2y^t Q_{\tau}^2 X\theta - y^t Q_{\tau}^2 y - (\theta - \tilde{P}_k y)^t X^t X (\theta - \tilde{P}_k y) \\ + y^t P_k y + (\theta - T) X^t Q_{\tau}^2 X (\theta - T) \} \Psi(\theta|y) d\theta \\ = \frac{1}{2} \left\{ \frac{1}{\sigma^2} y^t (P_k - Q_{\tau})^2 y + 2 \text{trace } Q_{\tau}^2 - k \right. \\ \left. + (T - \tilde{P}_k y)^t X^t Q_{\tau}^2 X (T - \tilde{P}_k y) \right\} .$$

Set

$$(5.19) \quad \nu = \nu_{\tau}(y) = \nu(\tau(y), \tilde{P}_k y) \quad (\text{say})$$

(see (4.10) and (4.11)), and

$$(5.20) \quad B_2(\nu) = \frac{1}{\sigma^2} y^t (P_k - Q_{\tau})^2 y + 2 \text{trace } Q_{\tau}^2 - k .$$

In the same way as Lemmas 4.3 and 4.4, it holds that

$$(5.21) \quad B_2(\nu) = \frac{\nu^2}{\sigma^2} \sum_{i=1}^k \frac{z_i^2 \lambda_i}{(\lambda_i + \nu)^2} + 2 \sum_{i=1}^k \frac{\lambda_i^2}{(\lambda_i + \nu)^2} - k \\ B_2(0+) = k , \quad B_2(+\infty) = \frac{1}{\sigma^2} |P_k y|^2 - k \\ \frac{\partial}{\partial \nu} B_2(\nu) = 2 \left\{ \frac{\nu}{\sigma^2} \sum_{i=1}^k \frac{z_i^2 \lambda_i^2}{(\lambda_i + \nu)^3} - 2 \sum_{i=1}^k \frac{\lambda_i^2}{(\lambda_i + \nu)^3} \right\} \\ \frac{\partial}{\partial \nu} B_2(\nu) < 0 , \quad \text{if } \nu < 2\sigma^2/z_{\max}^2 , \quad \frac{\partial}{\partial \nu} B_2(\nu) > 0 \quad \text{if } \nu > 2\sigma^2/z_{\min}^2 ,$$

where $z_{\max}^2 = \max_{i=1, \dots, k} z_i^2$ and $z_{\min}^2 = \min_{i=1, \dots, k} z_i^2$ (see (4.30)). Therefore there exists $\nu_2^* > 0$ such that

$$(5.22) \quad B_2(\nu_2^*) = \min \{B_2(\nu) : 0 < \nu < \infty\}$$

and hence, we can obtain α_2^* by the virtue of ν_2^* and the equation (4.32). Consequently, we have the following theorem from (5.18) and (5.22).

THEOREM 5.3. *The model fitting $(\alpha_2^*, \hat{\theta})$, (recall $\hat{\theta} = \tilde{P}_k y$), based on ridge estimation, is the Bayes solution under loss function $W(\tau, T|\theta)$ of (5.16) and Lebesgue prior distribution.*

6. Remarks

(a) It is apparent that ridge estimators obtained by using norms $|\zeta| = |\zeta|_0$ and $|\zeta|_x$ in (4.4) correspond to a general ridge estimator and a so called shrinkage one, respectively.

(b) It may be necessary that properties of estimators with decision rule (for example, AIC, C_p , $B(\nu)$ and so on) are investigated from the view of estimation theory.

(c) We might have to argue, first of all, what is statistical model fitting.

Acknowledgements

The author wishes to thank Professor H. Akaike for his suggestions and kind guidance of the problems with respect to the Akaike's Information Criterion. The author also wishes to thank Professor M. Okamoto and Mr. Y. Takada for several discussions and valuable comments.

OSAKA UNIVERSITY

REFERENCES

- [1] Aitchison, J. and Silvey, S. D. (1958). Maximum-likelihood estimation of parameters subject to restraints, *Ann. Math. Statist.* **29**, 813-828.
- [2] Akaike, H. (1972). Information theory and an extension of the maximum likelihood principle, *Proc. 2nd Intl. Symp. on Information Theory*, Supplement to Problems of Control and Information Theory, 267-281.
- [3] Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: Biased estimation for non-orthogonal problems, *Technometrics*, **12**, 55-67.
- [4] Ibragimov, I. A. and Khas'minskii, R. Z. (1972). Asymptotic behavior of statistical estimators in the smooth case, *Theory Prob. Appl.*, **17**, 443-460.
- [5] Ibragimov, I. A. and Khas'minskii, R. Z. (1973). Asymptotic behavior of some statistical estimators II. Limit theorems for the a posteriori density and Bayes' estimators, *Theory Prob. Appl.*, **18**, 76-91.
- [6] Inagaki, N. and Ogata, Y. (1975). The weak convergence of likelihood ratio random fields and its applications, *Ann. Inst. Statist. Math.*, **27**, 391-419.
- [7] Inagaki, N. and Ogata, Y. (1975). The weak convergence of likelihood ratio random

fields for Markov observation, *Research Memorandum*, No. 79, The Institute of Statistical Mathematics.

- [8] LeCam, L. (1970). On the assumptions used to prove asymptotic normality of maximum likelihood estimates, *Ann. Math. Statist.*, **41**, 802-828.
- [9] Lindley, D. V. (1968). The choice of variables in multiple regression, *J. R. Statist. Soc. B*, **30**, 31-53.
- [10] Lindley, D. V. and Smith, A. F. M. (1972). Bayes estimates for linear model, *J. R. Statist. Soc. B*, **34**, 1-18.
- [11] Mallows, C. L. (1973). Some comments on C_p , *Technometrics*, **15**, 661-675.