

ON MODEL SELECTION AND THE ARC SINE LAWS¹

BY MICHAEL WOODROOFE

University of Michigan

Generalizations of the arc sine laws are shown to provide insight into the operating characteristics of certain techniques for selecting models to fit a given data set, when the available models are nested. As a corollary, one sees that a popular technique may be expected to include about one superfluous parameter, even if the sample size is large.

1. Introduction. There are several techniques which may be used to select an appropriate model from a class of available models to fit a given data set. In particular, Mallows's (1964, 1973) C_p criterion and Akaike's (1974) entropy maximization criterion have been recommended for use in model selection. Assuming that one of the available models is correct, one may inquire about such operating characteristics as the probability of finding the smallest correct model and the distribution of the number of superfluous parameters in the model selected. Here these operating characteristics are studied in the special case that the models are nested—as, for example, in polynomial regression and moving average models for time series. In this case there is a natural relation between the operating characteristics and certain generalizations of the arc sine law, as described by Feller (1966, Chapter 12), for example. Briefly, the selection techniques choose the model for which a criterion is maximum; and the generalized arc sine laws determine the distribution of the index for which sums of i.i.d. random variables attain a maximum. So, if the model selection criteria form sums of i.i.d. random variables, the arc sine laws may be used to determine the distribution of the index of the model selected.

In Section 2 the relation between the operating characteristics of Akaike's technique and the generalized arc sine laws is indicated in the simplest case—when the data are independent, normally distributed random variables with unknown means and unit variances. Then the generalized arc sine laws are reviewed in Section 3, and applied to the normal case in Section 4. A natural extension then yields the operating characteristics for Mallows' C_p in Section 5. In Sections 6 and 7, the simple normal example is shown to provide an asymptotic distribution for the number of superfluous parameters for models with a well-behaved likelihood function and a large sample size. As a corollary, it is noted that Akaike's technique is inconsistent in large samples. The same is true of Mallows' C_p , but not of Schwarz's (1978) Bayesian criterion. These remarks are detailed in Section 8.

There are several related articles. The formulation of the problem as one of selecting from multiple hypotheses is similar to that in Anderson's (1962) determination of the degree to use in polynomial regression. Anderson developed optimal procedures. The emphasis here is on the properties of suboptimal, though closely related, procedures. That Akaike's technique is inconsistent in large samples was shown by Shibata (1976) for autoregressive processes and more generally by Hinkley (1976) in an unpublished manuscript. Much of the present Section 6 was anticipated in the latter. Recently, Shibata (1980, 1981) has studied model selection techniques under a different limiting operation. Under this operation, Akaike's technique and Mallows's C_p are asymptotically efficient.

2. Akaike's Criterion. This technique starts with a large model which is assumed to be correct, but possibly redundant, and eliminates parameters which appear to be

Received November 1981; revised May 1982.

¹ Research supported by the National Science Foundation under MCS-8101897.

AMS 1980 subject classifications. 62F99; 62J05.

Key words and phrases. Akaike's criterion, Asymptotic distributions, Mallows C_p , Random walks.

superfluous. Thus, let $X = (X_1, \dots, X_n)'$ denote a random column vector having a density $f(\cdot; \theta)$, where $\theta = (\theta_1, \dots, \theta_k)'$ is a vector of unknown parameters taking values in an open subset $\Omega \subset R^k$; for $j = 0, \dots, k$, let

$$\Omega_j = \{\theta \in \Omega: \theta_i = 0, i = j + 1, \dots, k\};$$

and let $H_j: \theta \in \Omega_j$ be the assertion that the smaller model, with Ω replaced by Ω_j , contains the true distribution of X . Thus, the models are assumed to be nested, as in polynomial regression or moving average models for time series. Next, let

$$\Lambda_j = \sup_{\theta \in \Omega_j} \log f(X; \theta)$$

be the maximum value of the log likelihood, assuming $\Omega_j, j = 0, \dots, k$. Then Akaike's technique selects the model for which

$$AIC(j) = \Lambda_j - j = \max.$$

Now let θ^0 denote the true value of the parameter. That is, suppose that X has density $f(\cdot; \theta^0)$. Then Akaike's technique selects the model for which

$$AIC^*(j) = \Lambda_j^* - j = \max,$$

where

$$\Lambda_j^* = \sup_{\theta \in \Omega_j} \log f(X; \theta) - \log f(X; \theta^0)$$

is the log likelihood ratio statistic for testing $\theta = \theta^0$ vs. $\theta \in \Omega_j, j = 0, \dots, k$. Thus, the index of the model selected is

$$(1) \quad J_k = \min\{j: AIC(j) = \max_{0 \leq i \leq k} AIC(i)\};$$

and AIC may be replaced by AIC^* .

To see the relation between Akaike's Criterion and the generalized arc sine laws, consider the simple special case in which $n = k, X_1, \dots, X_k$ are independently, normally distributed random variables with unknown means $\theta_1, \dots, \theta_k$ and unit variances, and $\theta^0 = (0, \dots, 0)'$. Then,

$$AIC^*(j) = \frac{1}{2} \sum_{i=1}^j X_i^2 - j = \sum_{i=1}^j \left(\frac{1}{2} X_i^2 - 1 \right) = S_j, \text{ say,}$$

for $j = 0, \dots, k$. Observe that S_1, \dots, S_k form an initial segment of a random walk with negative drift, since $E(\frac{1}{2} X_i^2 - 1) = -\frac{1}{2}, i = 1, \dots, k$. Thus,

$$(2) \quad J_k = \min\{j: S_j = \max_{1 \leq i \leq k} S_i\};$$

and J_k is the number of superfluous parameters included, since $\theta^0 = 0$. The distribution of random variables of the form J_k have been studied extensively in the context of general random walks S_1, S_2, \dots . Some of the relevant results are reviewed in the next section.

3. The generalized arc sine laws. Let Y_1, Y_2, \dots be any sequence of i.i.d. random variables, and let $S_j, j \geq 0$, denote the associated random walk, $S_0 = 0$ and $S_j = Y_1 + \dots + Y_j, j \geq 1$. Further, let $p_0 = q_0 = 1$,

$$p_j = P\{S_1 > 0, \dots, S_j > 0\}$$

and

$$q_j = P\{S_1 \leq 0, \dots, S_j \leq 0\}, \quad j \geq 1.$$

Then

$$(3) \quad P\{J_k = j\} = p_j q_{k-j}, \quad \text{for } 0 \leq j \leq k, k \geq 1,$$

where $J_k, k \geq 1$, are defined by (2). The quantities $p_j, j \geq 0$, and $q_j, j \geq 0$, may be determined from the generating functions

$$(4) \quad P(s) = \sum_{j=0}^{\infty} p_j s^j \quad \text{and} \quad Q(s) = \sum_{j=0}^{\infty} q_j s^j, \quad 0 < s < 1,$$

in view of the following, remarkable identities. Let

$$a_n = P\{S_n > 0\}, \quad n \geq 1;$$

then

$$(5a) \quad P(s) = \exp\left\{\sum_{n=1}^{\infty} \frac{1}{n} a_n s^n\right\}, \quad 0 < s < 1,$$

and

$$(5b) \quad Q(s) = \exp\left\{\sum_{n=1}^{\infty} \frac{1}{n} (1 - a_n) s^n\right\}, \quad 0 < s < 1.$$

For example, if Y_1 has a continuous, symmetric distribution, then $a_n = 1/2$ for all $n \geq 1$, so that $P(s) = Q(s) = 1/\sqrt{1-s}$, $0 < s < 1$, and $P\{J_k = j\} = (-1)^k \binom{-1/2}{j} \binom{-1/2}{k-j}$ for $0 \leq j \leq k$ and $k \geq 1$. That is, J_k has the classical, discrete arc sine distribution as described by Feller (1968, Chapter 3), for example. The results of this paragraph are taken directly from Feller (1966, Chapter 12).

The mean and variance of J_k may be computed easily from (3). In fact, $E(J_k) = \sum_{j=0}^k j p_j q_{k-j}$ is the coefficient of s^{k-1} in the expansion of $P'(s)Q(s)$; and the latter is easily seen to be $a_1 + \dots + a_k$. So,

$$(6) \quad E(J_k) = a_1 + \dots + a_k, \quad k \geq 1;$$

and

$$D(J_k) = \sum_{j=1}^k j a_j - \sum_{1 \leq i < j \leq k} a_i a_j, \quad k \geq 1,$$

by a similar argument.

There is natural interest in the distribution of J_k when k is large; and, if $E(Y_1) < 0$, then the latter is easily determined. In fact, if $E(Y_1) < 0$, then the series $\sum_{n=1}^{\infty} n^{-1} a_n$ is convergent and

$$(7) \quad \lim_{k \rightarrow \infty} q_k = \lim_{s \uparrow 1} \{(1-s)Q(s)\} = \exp\left(-\sum_{n=1}^{\infty} \frac{1}{n} a_n\right) = q_{\infty}, \quad \text{say};$$

so,

$$P\{J_k = j\} \rightarrow q_{\infty} p_j,$$

as $k \rightarrow \infty$ for all $j \geq 0$.

4. Normal case. Now reconsider the simple normal case in which X_1, \dots, X_k are independent normally distributed random variables with means $\theta_1, \dots, \theta_k$ and unit variances. If $\theta_i = 0$ for all $i \leq k$, then $AIC^*(j) = Y_1 + \dots + Y_j$ for $1 \leq j \leq k$ with $Y_i = 1/2 X_i^2 - 1$ for $1 \leq i \leq k$; and the generalized arc sine laws may be applied directly to find the distribution of J_k , with

$$a_j = P(\chi_j^2 > 2j), \quad 1 \leq j \leq k.$$

For example, the limiting distribution of J_k as $k \rightarrow \infty$ is $\lim P_0\{J_k = j\} = q_{\infty} p_j$ for $j \geq 0$, where $p_j, j \geq 0$, are determined from (5) and q_{∞} is as in (7).

Table 1 lists the exact distribution of J_k for $k = 5, 10$ and the limiting distribution as $k \rightarrow \infty$. Observe that the probabilities approach their limits quite quickly, but that the convergence of $E(J_k)$ is slower. Observe also that the limiting distribution assigns slightly more than 1% of its mass to integers $j > 10$. Two of the numbers in Tables 1 and 2 are of special interest. If $\theta = 0$ and k is large, then the probability of correctly determining that $\theta = 0$ is approximately 0.712, but the expected number of superfluous parameters included is approximately 0.946.

The exact computations for the special case $\theta = 0$ provide bounds for the general case. To see how, let

TABLE 1
The Distribution of J_k for $k = 5, 10, \infty$.

j	$k = 5$	$k = 10$	$k = \infty$
0	.736	.718	.712
1	.117	.113	.112
2	.061	.058	.057
3	.038	.035	.035
4	.027	.023	.023
5	.022	.016	.016
6		.012	.011
7		.0088	.0083
8		.0067	.0061
9		.0054	.0046
10		.0047	.0035
> 10			.0115
$E(J_k)$.571	.791	.946

The computations were done on an Apple II microcomputer, using formula (24.4.6) of Abramowitz and Stegun (1970) to compute the Chi squared distribution function. $E(J_k)$ was computed from (6).

$$(8) \quad \gamma_k(j; \theta) = P_\theta(J_k > j),$$

for $j = 0, \dots, k, k \geq 1$, and $\theta \in R^k$. If $\theta \in \Omega_r - \Omega_{r-1}$, where $1 \leq r < k$, then $\gamma_k(r + j; \theta)$ is the probability that the criterion includes more than j superfluous parameters in the model for $j = 0, \dots, k - r - 1$.

THEOREM 1. For $1 \leq r < k$ and $0 \leq j < k - r$,

$$\sup_{\theta \in \Omega_r - \Omega_{r-1}} \gamma_k(r + j; \theta) = \gamma_{k-r}(j; 0).$$

PROOF. Let $S_0 = 0$ and $S_j = Y_1 + \dots + Y_j$ for $1 \leq j \leq k$. Then $J_k > r + j$ iff $S_{r+i} > \max(S_0, \dots, S_{r+j})$ for some $i \leq k - r$. If $\theta \in \Omega_r - \Omega_{r-1}$, then

$$(9) \quad \begin{aligned} P_\theta(J_k > r + j) &= P_\theta\{S_{r+i} > \max(S_0, \dots, S_{r+j}), \exists i \leq k - r\} \\ &\leq P_\theta\{S_{r+i} - S_r > \max(S_r - S_r, \dots, S_{r+j} - S_r), \exists i \leq k - r\} \\ &= P_0\{S_i > \max(S_0, \dots, S_j), \exists i \leq k - r\} = \gamma_{k-r}(j; 0). \end{aligned}$$

Moreover, the difference between the first and second lines in (9) is at most

$$P_\theta\{\max(0, S_1, \dots, S_{r+j}) > \max(S_r, \dots, S_{r+j}), \exists j \leq k - r\} \leq P_\theta\{S_r < \max(0, \dots, S_r)\}$$

which tends to zero as $\theta_r \rightarrow \infty$.

5. Mallows's criterion. Now consider a linear model

$$X = M\beta + \varepsilon,$$

where M is an $n \times k$ matrix of full rank $k < n$ and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)'$ has the n -variate normal distribution with mean vector 0 and covariance matrix $\sigma^2 I_n$. Here the unknown parameters are $\beta \in R^k$ and $\sigma^2 > 0$; and the nested models are

$$\Omega_j = \{\beta \in R^k: \beta_i = 0, i = j + 1, \dots, k\}, \quad 0 \leq j \leq k.$$

Let SSE_j denote the error sum of squares when the model $H_j: \beta \in \Omega_j, 0 < \sigma^2 < \infty$ is fit; and let $\hat{\sigma}^2 = SSE_k / (n - k)$ denote the unbiased estimator of σ^2 when all k regression parameters

are fit. Then

$$C_p = \hat{\sigma}^{-2} \cdot \text{SSE}_p + (2p - n)$$

has been suggested as a criterion for judging the adequacy of the models H_p , $0 \leq p \leq k$. See Mallows (1964, 1973) and, for example, Daniel and Wood (1980, pages 86-90). The distribution of the index

$$\mathcal{J}_k = \min\{p: C_p = \min_{1 \leq j \leq k} C_j\}$$

may be found by a natural extension of the techniques of Sections 3 and 4.

Let $L_j \subset R^k$ be the linear subspace spanned by the first j columns of M , $1 \leq j \leq k$; and let e_1, \dots, e_n be an orthonormal basis for R^n for which e_1, \dots, e_j is an orthonormal basis for L_j for $1 \leq j \leq k$. Then

$$\text{SSE}_j = \sum_{i=j+1}^n Z_i^2, \quad 1 \leq j \leq k,$$

where

$$Z_i = e_i' Y, \quad 1 \leq i \leq n,$$

are independent normally distributed random variables with means $\theta_i = e_i' M \beta$ for $1 \leq i \leq k$, and $e_i' M \beta = 0$ for $k < i \leq n$, and common variance σ^2 . See, for example, Lehmann (1959, Section 7.2). It follows that

$$C_0 - C_p = \hat{\sigma}^{-2} \sum_{i=1}^p Z_i^2 - 2p = \hat{\sigma}^{-2} S_p^*,$$

where

$$S_p^* = \sum_{i=1}^p (Z_i^2 - 2\hat{\sigma}^2), \quad 1 \leq p \leq k;$$

and

$$\mathcal{J}_k = \min\{j: S_j^* = \max_{1 \leq t \leq k} S_t^*\}.$$

As in Section 4, the distribution of \mathcal{J}_k in the special case that $\beta = 0$ provides a bound for the general case. Let

$$\Delta_{n,k}(j; \beta) = P_{\sigma,\beta}(\mathcal{J}_k > j)$$

for $0 \leq j \leq k$, $1 \leq k < n$, $\beta \in R^k$.

THEOREM 1'. For $1 \leq r < k$ and $0 \leq j < k - r$,

$$\sup_{\beta \in \Omega_r - \Omega_{r-1}} \Delta_{n,k}(r + j; \beta) = \Delta_{n-r,k-r}(j; 0).$$

The proof of Theorem 1' is similar to that of Theorem 1 and has been omitted.

Suppose now that $\beta_j = 0$ for $j = 1, \dots, k$, so that Z_1, \dots, Z_n are i.i.d. normally distributed random variables with common mean 0 and common variance σ^2 . Then, since Z_1, \dots, Z_k are independent of $\hat{\sigma}^2$, the conditional distribution of \mathcal{J}_k , given $\hat{\sigma}^2$, may be found from the techniques of Sections 3 and 4. Let

$$a_j(t) = P(\chi_j^2 > 2jt), \quad t > 0, j \geq 1;$$

let $p_j(t)$ and $q_j(t)$, $j \geq 0$, be defined by (4) and (5) with a_j , $j \geq 1$, replaced by $a_j(t)$, $j \geq 1$, for each $t > 0$; and define $q_\infty(t)$ by (7) with a_j replaced by $a_j(t)$ for $t > 1/2$. Then

$$(10) \quad P_{\sigma,0}(\mathcal{J}_k = j | \hat{\sigma}^2) = p_j(\hat{\sigma}^2/\sigma^2) q_{k-j}(\hat{\sigma}^2/\sigma^2)$$

for $0 \leq j \leq k$, $1 \leq k < n$, and $\sigma^2 > 0$; and the unconditional distribution of \mathcal{J}_k may be found by integrating over the possible values of $\hat{\sigma}^2$. In particular, the (unconditional) expectation of \mathcal{J}_k is

$$E_{\sigma,0}(\mathcal{J}_k) = \sum_{j=1}^k E_{\sigma,0}\{a_j(\hat{\sigma}^2/\sigma^2)\} = \sum_{j=1}^k P\{\mathcal{F}(j, n - k) > 2\},$$

where $\mathcal{F}(j, n - k)$ denotes a random variable having the F -distribution on j and $n - k$ degrees of freedom, $1 \leq j \leq k$.

TABLE 2
Expected Values of \mathcal{J}_k

k	n						
	12	24	36	48	96	192	∞
6	1.274	.873	.780	.739	.683	.657	.633
12		1.762	1.315	1.156	.972	.899	.836
18		3.712	1.952	1.526	1.132	1.003	.904
24			3.102	1.995	1.245	1.056	.926

The computations were done on an Apple II microcomputer, using formula (26.6.5) of Abramowitz and Stegun (1970) to compute the F distribution function and formula (26.4.6) to compute the Chi squared distribution function. The last column gives the values of $E(\mathcal{J}_k)$ from Section 4.

Selected values of $E_{\sigma,0}(\mathcal{J}_k)$ are listed in Table 2. Observe that the convergence of $E_{\sigma,0}(\mathcal{J}_k)$ to its limit is quite slow as $n \rightarrow \infty$.

If $n - k$ is large, then the distribution of $\hat{\sigma}^2/\sigma^2$ is nearly degenerate at 1, suggesting that the conditional probabilities in (10) might be expanded in a Taylor series about 1. This expansion is detailed in Theorem 2. First, the derivatives of $p_j(t)$ and $q_j(t)$ w.r.t. t are studied.

LEMMA 1. For $j \geq 1$, the derivatives of $p_j(t)$ and $q_j(t)$ w.r.t. $t > 0$ are

$$(11) \quad \dot{p}_j(t) = \sum_{i=1}^j \frac{1}{i} \dot{a}_i(t) p_{j-i}(t)$$

and

$$(12) \quad \dot{q}_j(t) = - \sum_{i=1}^j \frac{1}{i} \dot{a}_i(t) q_{j-1}(t),$$

where

$$\dot{a}_i(t) = -i^{(1/2)} t^{(1/2)i-1} e^{-it} / \Gamma(1/2 i)$$

denotes the derivative of $a_i(t)$ w.r.t. $t > 0$ for $i \geq 1$. Moreover, (12) holds when $j = \infty$ and $t > 1/2$ too.

PROOF. Let $P(t, s)$ denote the generating function of $p_j(t)$, $j \geq 0$, for $0 < s < 1$ and $t > 0$. Then $p_j(t)$ is the coefficient of s^j in $\dot{P}(t, s) = \partial P(t, s) / \partial t$; and

$$\dot{P}(t, s) = \left\{ \sum_{i=1}^{\infty} \frac{1}{i} \dot{a}_i(t) s^i \right\} P(t, s), \quad 0 < s < 1, t > 0.$$

Relation (11) follows immediately; and (12) may be established similarly for $1 \leq j < \infty$.

Observe that $|\dot{a}_1(t)| + |\dot{a}_2(t)| + \dots$ is convergent uniformly in $t \geq 1/2 + \epsilon$ for any $\epsilon > 0$, by Stirling's Formula. Thus, for $t > 1/2$,

$$\log q_{\infty}(t) = - \sum_{j=1}^{\infty} \frac{1}{j} a_j(t)$$

may be differentiated term by term. That (12) holds when $j = \infty$ follows.

When iterated, (11) and (12) yield expressions for the second derivatives of $p_j(t)$ and $q_j(t)$, $j \geq 1$. For example,

$$\ddot{q}_j(t) = - \sum_{i=1}^j \frac{1}{i} \{ \ddot{a}_i(t) q_{j-i}(t) + \dot{a}_i(t) \dot{q}_{j-i}(t) \}$$

for $t > 0$ and $j \geq 1$; and, using the uniform convergence of $|\dot{a}_1(t)| + |\dot{a}_2(t)| + \dots$ and $|\ddot{a}_1(t)| + |\ddot{a}_2(t)| + \dots$, one finds easily that $\ddot{q}_j(t) \rightarrow \ddot{q}_{\infty}(t)$ uniformly in t on compact

subintervals of $(\frac{1}{2}, \infty)$ as $j \rightarrow \infty$.

THEOREM 2. *Let $k = k_n, n \geq 1$, be integers for which $0 < n - k \rightarrow \infty$ as $n \rightarrow \infty$. Then*

$$P_{\sigma,0}(\mathcal{J}_k = j) = p_j q_{k-j} + \{\ddot{p}_j q_{k-j} + 2\dot{p}_j \dot{q}_{k-j} + p_j \ddot{q}_{k-j}\} \left(\frac{1}{n-k}\right) + o\left(\frac{1}{n-k}\right)$$

as $n \rightarrow \infty$ for each fixed $j \geq 0$, where $p_j = p_j(1), \dot{p}_j = \dot{p}_j(1)$, etc.

PROOF. Since the probabilities are independent of σ^2 , there is no loss of generality in supposing that $\sigma^2 = 1$. Let F_n denote the distribution function of $\hat{\sigma}^2$; and, for fixed $j \geq 0$, let $g_n(t) = \ddot{p}_j(t)q_{k-j}(t) + 2\dot{p}_j(t)\dot{q}_{k-j}(t) + p_j(t)\ddot{q}_{k-j}(t)$ for $t > 0$ and $n \geq 1$. Then, for any $\delta > 0$,

$$\begin{aligned} P_{\sigma,0}(\mathcal{J}_k = j) &= \int_{1-\delta}^{1+\delta} p_j(t)q_{k-j}(t)dF_n(t) + o\left(\frac{1}{n-k}\right) \\ (13) \qquad &= p_j q_{k-j} + g_n(1)\left(\frac{1}{n-k}\right) \\ &\quad + \int_{1-\delta}^{1+\delta} \frac{1}{2} \{g_n(t_n^*) - g_n(1)\}(t-1)^2 dF_n(t) + o\left(\frac{1}{n-k}\right), \end{aligned}$$

where t_n^* denotes an intermediate point between t and 1. Indeed, (13) follows from elementary properties of the Chi squared distribution. Finally, it follows from Lemma 1 that $g_n, n \geq 1$, are equicontinuous in t on compact subintervals of $(\frac{1}{2}, \infty)$; so, given $\epsilon > 0$, the last integral in (13) may be made less than $\epsilon \int (t-1)^2 dF_n(t)$ by taking $\delta > 0$ sufficiently small; and since $\int (t-1)^2 dF(t) = 2/(n-k)$, the theorem follows.

6. Asymptotics with large n. Now consider a sequence of problems, indexed by the sample size $n \geq 1$. Suppose first that the parameter space Ω is the same for all sample sizes. Thus, Ω is an open subset of R^k for some $k \geq 1; 0 \in \Omega$; and the nested models of interest are $\Omega_j = \{\theta \in \Omega: \theta_i = 0 \text{ for } i = j + 1, \dots, k\}$. The data X_1, \dots, X_n are assumed to have a joint density $f_n(\cdot; \theta)$ w.r.t. a dominating (sigma-finite) measure for all $\theta \in \Omega$ for each $n \geq 1$; so, the log-likelihood function may be written

$$L_n(\theta) = \log f_n(X_1, \dots, X_n; \theta), \quad \theta \in \Omega,$$

for each $n \geq 1$. As in Section 2, the true value of the parameter is denoted by $\theta^0 = (\theta_1^0, \dots, \theta_k^0)'$; and the index of the model selected may be written

$$(14) \qquad J_{nk} = \min\{j: \Lambda_{nj}^* - j = \max_{0 \leq i \leq k} \Lambda_{ni}^* - i\},$$

where

$$(15) \qquad \Lambda_{ni}^* = \sup_{\theta \in \Omega_i} L_n(\theta) - L_n(\theta^0)$$

for $i = 1, \dots, k$ and $n \geq 1$. It is shown that the simple normal example of Section 4 provides an asymptotic distribution for J_{nk} as $n \rightarrow \infty$, under some regularity conditions.

It is most efficient to state the regularity conditions directly in terms of the likelihood function. In their statements, P denotes a probability measure under which X_1, \dots, X_n have joint density $f_n(\cdot; \theta^0)$ w.r.t. the dominating measure. Thus, the dependence of P on n and θ^0 is suppressed in the notation.

CONDITION C1. For every $\epsilon > 0$,

$$\sup_{\|\theta - \theta^0\| \geq \epsilon} L_n(\theta) - L_n(\theta^0) \rightarrow_P -\infty$$

as $n \rightarrow \infty$, where $\|\cdot\|$ denotes the Euclidean norm.

If a maximum likelihood estimator exists, then C1 guarantees that it converges to θ^0 in

probability as $n \rightarrow \infty$.

CONDITION C2. For some $\varepsilon_0 > 0$, $L_n(\theta)$ is twice continuously differentiable in $\|\theta - \theta^0\| < \varepsilon_0$ w.p.1 (P) for all sufficiently large n .

If C2 is satisfied, then the gradient and Hessian,

$$Z_n(\theta) = \left[\frac{\partial}{\partial \theta_1} L_n(\theta), \dots, \frac{\partial}{\partial \theta_k} L_n(\theta) \right],$$

and

$$M_n(\theta) = \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} L_n(\theta); i, j = 1, \dots, k \right],$$

are well defined for $\|\theta - \theta^0\| < \varepsilon_0$ and sufficiently large n . It is convenient to write $Z_n = Z_n(\theta^0)$.

CONDITION C3. Condition C2 is satisfied; and there are $\varepsilon_1, 0 < \varepsilon_1 < \varepsilon_0$, positive constants $\alpha_n, n \geq 1$, and non-random matrices $M_\theta, \|\theta - \theta^0\| \leq \varepsilon_1$ for which (i) $\alpha_n \rightarrow \infty$, as $n \rightarrow \infty$, (ii) M_θ are continuous in θ and positive definite, (iii) $Z_n/\sqrt{\alpha_n}$ is asymptotically normal with mean 0 and covariance matrix $M = M_{\theta^0}$, and (iv) $\sup_{\|\theta - \theta^0\| \leq \varepsilon_1} \|\alpha_n^{-1} M_n(\theta) - M_\theta\|_{tr} \rightarrow 0$ in probability as $n \rightarrow \infty$, where $\|\cdot\|_{tr}$ denotes the trace norm.

Conditions C1, C2, and C3 imply that a maximum likelihood estimator $\hat{\theta}^n$ exists with probability approaching one, and that $\sqrt{\alpha_n}(\hat{\theta}^n - \theta^0)$ is asymptotically normal with mean vector 0 and covariance matrix M^{-1} , where $M = M_{\theta^0}$.

THEOREM 3. Suppose that conditions C1, C2, and C3 are satisfied. Suppose also that $\theta^0 \in \Omega_h - \Omega_{h-1}$ for some $h, 1 \leq h < k$. Then

$$\lim_{n \rightarrow \infty} P(J_{nk} - h > j) = \gamma_{k-h}(j; 0)$$

for $j = 0, \dots, k - h - 1$, where γ_k is as in (8).

Observe that $J_{nk} - h$ is the number of superfluous parameters included in the model selected.

PROOF. Since the distance from θ^0 to Ω_i is positive for all $i \leq h - 1$, it follows directly from C1 that $\max_{0 \leq i \leq h} \Lambda_{ni}^* \rightarrow_p -\infty$ as $n \rightarrow \infty$; so,

$$J_{nk} - h = \min\{j : \Lambda_{n,h+j}^* - (h + j) = \max_{0 \leq i \leq k-h} \Lambda_{n,h+i}^* - (h + i)\}$$

with probability approaching one as $n \rightarrow \infty$. Let ε_1 be as in the statement of Condition C3; define stochastic processes $W_n(t), t \in R^k, n \geq 1$, by

$$W_n(t) = \begin{cases} L_n\left(\theta^0 + \frac{t}{\sqrt{\alpha_n}}\right) - L_n(\theta^0), & \|t\| \leq \varepsilon_1 \sqrt{\alpha_n}, \\ L_n\left(\theta^0 + \frac{\varepsilon_1 t}{\|t\|}\right) - L_n(\theta^0), & \|t\| > \varepsilon_1 \sqrt{\alpha_n}, \end{cases}$$

and let $K_j = \{t \in R^k : t_i = 0 \text{ for } i = j + 1, \dots, k\}$ for $j = 1, \dots, k$. Then for $h \leq j < k$,

$$\Lambda_{n,j}^* = \sup_{t \in K_j} W_n(t)$$

with probability approaching one as $n \rightarrow \infty$. Now, for each fixed $t \in R^k, W_n(t)$ converges in distribution to $W(t) = t'Z - \frac{1}{2}t'Mt$ as $n \rightarrow \infty$, where Z has the normal distribution with mean vector 0 and covariance matrix M ; in fact, it follows from C3 that the joint distribution of $\Lambda_{nh}^* - h, \dots, \Lambda_{nk}^* - k$ converges to that of S_h, \dots, S_k as $n \rightarrow \infty$ where

$$S_j = \sup_{t \in K_j} (t'Z - \frac{1}{2}t'Mt) - j$$

for $1 \leq j \leq k$. Finally, by straightforward linear algebra, there are i.i.d. standard, univariate normal random variables Y_1, \dots, Y_k for which

$$S_j = \sum_{i=1}^j (\frac{1}{2} Y_i^2 - 1), \quad j = 1, \dots, k.$$

It follows easily that $J_{nk} - h$ converges in distribution to

$$J = \min\{j : S_{h+j} = \max_{0 \leq i \leq k-h} S_{h+i}\} = \min\{j : S_{h+j} - S_h = \max_{0 \leq i \leq k-h} S_{h+1} - S_h\},$$

which has the distribution (8), with k replaced by $k - h$.

7. Asymptotics with large n and large k . Let R^∞ denote the set of all infinite sequences of real numbers $x = (x_1, x_2, \dots)$, endowed with the product topology; let Θ be an open subset of R^∞ for which $(0, 0, \dots) \in \Theta$; let

$$\Omega_k = \{\theta \in \Theta : \theta_i = 0 \text{ for all } i > k\}, \quad k \geq 1;$$

and let

$$\Omega = \bigcup_{k=1}^\infty \Omega_k.$$

For each $n \geq 1$, let X_1, \dots, X_n be random vectors with joint density $f_n(\cdot; \theta)$ w.r.t. a dominating measure for some unknown $\theta \in \Omega$; let $k_n, n \geq 1$, be a non-decreasing sequence of positive integers for which $k_n \rightarrow \infty$ as $n \rightarrow \infty$; and suppose that Akaike's technique is applied with Ω replaced by Ω_{k_n} at the n th stage for each $n \geq 1$. Then the index of the model selected is

$$J_n = \min\{j : \Lambda_{n,j}^* - j = \max_{0 \leq i \leq k_n} \Lambda_{n,i}^* - i\},$$

where $\Lambda_{n,1}^*, \Lambda_{n,2}^*, \dots$ are defined by (15). The results of Sections 4 and 6 suggest that J_n may have the limiting distribution $q_\infty p_j, j \geq 0$, under general conditions.

As in Section 6, the true value of the parameter is denoted by θ^0 , and P denotes a probability measure under which X_1, \dots, X_n have density $f_n(\cdot; \theta^0)$ w.r.t. the dominating measure. It is assumed below that $\theta^0 \in \Omega$, and that Conditions C1, C2, and C3 are satisfied when Ω is replaced by Ω_k for all large k . In addition, the following condition is needed.

CONDITION C4. For every $\epsilon > 0$ there is an integer $\ell_0 \geq 1$ for which

$$(16) \quad P\{\max_{\ell_0 \leq \ell \leq k_n} \Lambda_{n,\ell}^* - \frac{3}{4} \ell \geq 0\} < \epsilon$$

for all sufficiently large n .

THEOREM 4. Suppose that $\theta^0 \in \Omega_h - \Omega_{h-1}$, where $h \geq 1$. Suppose also that Conditions C1, C2, C3, and C4 are satisfied. Then

$$\lim_{n \rightarrow \infty} P(J_n - h = j) = q_\infty p_j$$

for all $j \geq 0$, where q_∞ and $p_j, j \geq 0$, are as in (4) and (5).

PROOF. Given $j \geq 0$ and $\epsilon > 0$, there is an integer $\ell_0 \geq h$ for which $|q_{\ell_0-j} - q_\infty| < \epsilon$ and $P^n(J_n = J_{n,\ell_0}) \geq 1 - \epsilon$ for all sufficiently large n , where J_{n,ℓ_0} is defined by (14). Since $\lim P^n\{J_{n,\ell_0} - h = j\} = q_{\ell_0-j} p_j$ as $n \rightarrow \infty$, by Theorem 3,

$$\begin{aligned} q_\infty p_j - 2\epsilon &\leq \liminf_{n \rightarrow \infty} P(J_n - h = j) \\ &\leq \limsup_{n \rightarrow \infty} P(J_n - h = j) \leq q_\infty p_j + 2\epsilon; \end{aligned}$$

and, since $\epsilon > 0$ was arbitrary, the theorem follows.

Condition C4 is related to the order of consistency of the maximum likelihood estimator when $k = k_n \rightarrow \infty$ with n , a difficult question. See, for example, Huber (1973) and Yohai and Maronna (1979) for discussions of this question for M estimators. The following example indicates that C4 may be replacable by a growth condition on k_n .

EXAMPLE. Suppose that n is of the form $n = km$, and that (with the obvious conventions)

$$X_{i,j} = \theta_i + U_{i,j}$$

where U_{11}, \dots, U_{km} are i.i.d. with a common distribution G . Suppose further that G has a positive, bounded, twice continuously differentiable density g (w.r.t. Lebesgue measure) which has finite Fisher information and satisfies some other mild conditions, described in the Appendix. Then Condition C4 is satisfied if

$$(17) \quad \log k = o(m) \quad \text{or} \quad k \log k = o(n) \quad \text{as } n \rightarrow \infty.$$

The proof of this assertion is given in the Appendix. The point of the example is this: in this simple regression problem, Condition C4 may be replaced by the mild growth condition (17).

8. Remarks. In the simple normal example of Section 3, it was shown that the probability of including no superfluous parameters is at least 0.712 and that the expected number of superfluous parameters is at most one for all values of n . This seemed reassuring—better than the author expected at the beginning of this study. These numbers seem much less reassuring in the context of Theorems 3 and 4, however, since it is possible to find the correct model with probability approaching one as $n \rightarrow \infty$. In fact, as explained below, Schwarz’s (1978) Bayesian criterion will do so, under appropriate regularity conditions.

To understand the behavior of Schwarz’s technique, suppose that X_1, \dots, X_n are i.i.d. random k vectors with common density

$$(18) \quad g_\theta(x) = \exp\{\theta'x - \psi(\theta)\}, \quad x \in R^k, \theta \in \Omega$$

with respect to some dominating, sigma-finite measure. Let Ω denote the natural parameter space of the family (18); suppose that Ω satisfies the conditions imposed in Section 6; and let $\Lambda_{n,j}^*$, $0 \leq j \leq k$, be as in Section 6. Then Schwarz’s technique selects the model for which $\Lambda_{n,j}^* - \frac{1}{2}j \log n$ is maximum, so the index of the model selected is

$$K_n = \min\{j : \Lambda_{n,j}^* - \frac{1}{2}j \log n = \max_{0 \leq i \leq k} \Lambda_{n,i}^* - \frac{1}{2}i \log n\}.$$

If $\theta^0 \in \Omega_r - \Omega_{r-1}$, where $r \geq 0$ and $\Omega_{-1} = \emptyset$, then $\Lambda_{n,j}^*/n$ converges to a negative value for $0 \leq j \leq r - 1$ and $\Lambda_{n,j}^*/n \rightarrow 0$ for $r \leq j \leq k$ w.p.1 as $n \rightarrow \infty$. So, $\liminf_{n \rightarrow \infty} K_n \geq r$ w.p.1. To bound the distribution of the number of superfluous parameters, observe that

$$P(K_n \geq r + i) \leq \sum_{j=r+i}^k P\{\Lambda_{n,j}^* > \frac{1}{2}(j - r) \log n + \Lambda_{n,r}^*\} \leq \sum_{j=r+i}^k P\{\Lambda_{n,j}^* > \frac{1}{2}(j - r) \log n\}$$

for $i = 1, \dots, k - r$. For smooth exponential families, the Chi squared approximation may be applied in the tail to yield

$$P(\Lambda_{n,j}^* > b \log n) \sim P(\frac{1}{2}\chi_j^2 > b \log n) \sim \{\Gamma(\frac{1}{2}j)\}^{-1} (b \log n)^{1/2j-1} n^{-b} \quad \text{as } n \rightarrow \infty$$

for $i = 1, \dots, k - r$. See Woodroffe (1978) for details. Thus,

$$P(K_n \geq r + i) = O\{n^{-1/2i}(\log n)^{1/2(r+i-1)}\} \quad \text{as } n \rightarrow \infty,$$

for $i = 1, \dots, k - r$. The asymptotic behavior of Schwarz’s technique is quite different from that of Akaike.

9. Acknowledgements. Thanks to Jan Kmenta for helpful discussions; and thanks to the editor and referees for helpful criticisms and for several of the references.

APPENDIX

The assertion made in the Example in Section 7 is proved here. The notations of Section 7 are used throughout.

THEOREM 5. *Suppose that n is of the form $n = km$, and that*

$$X_{ij} = \theta_i + U_{ij}, \quad 1 \leq j \leq m, 1 \leq i \leq k,$$

where U_{11}, \dots, U_{km} are i.i.d. with a common distribution G . Suppose further that G has a positive, bounded, twice continuously differentiable density g (w.r.t. Lebesgue measure) for which the following conditions are satisfied:

(i)
$$\mathcal{J} = \int_{-\infty}^{\infty} \frac{\{g'(x)\}^2}{g(x)} dx < \infty,$$

(ii) for some $\alpha, 0 < \alpha < 1$,
$$\int_{-\infty}^{\infty} \{g(x)\}^\alpha dx < \infty,$$
 and

(iii) letting $H(x) = \log g(x), -\infty < x < \infty, H''$ is bounded above,

$\mathcal{J} = \int -H'' dG$, and $\int \sup_{|t| < \epsilon} |H''(x - t)|^\beta dG(x) < \infty$ for some $\epsilon > 0$ and $\beta > 1$. Then Condition C4 holds, if $\log k = o(m)$.

PROOF. The conditions imply that there are maximum likelihood estimators $\hat{\theta}_i = \hat{\theta}_i(X_{i1}, \dots, X_{im})$ of θ_i for which $\hat{\theta}_1 - \theta_1, \dots, \hat{\theta}_k - \theta_k$ are i.i.d. (with a common distribution which is independent of $\theta_1, \dots, \theta_k$ for each fixed k and m ; moreover, $\sqrt{m}(\hat{\theta}_1 - \theta_1)$ is asymptotically normal with mean 0 and variance $1/\mathcal{J}$ as $m \rightarrow \infty$; and, for every $\epsilon > 0$, there is a $\rho = \rho(\epsilon)$ for which $0 < \rho < 1$ and

$$P\{|\hat{\theta}_1 - \theta_1| \geq \epsilon\} \leq C\rho^m$$

for all $m \geq 1$ for some constant C (cf Wald, 1949).

To simplify the exposition, suppose that $\theta^0 = (0, \dots, 0)'$ and that $\mathcal{J} = 1$. Then, for $1 \leq \ell \leq k$,

$$\Lambda_{n^\ell}^* = \sum_{i=1}^\ell m \{L_i(\hat{\theta}_i) - L_i(0)\},$$

where

$$L_i(t) = \frac{1}{m} \sum_{j=1}^m \{H(X_{ij} - t) - H(X_{ij})\}, \quad -\infty < t < \infty.$$

Expanding L_1, \dots, L_k in Taylor series about $\hat{\theta}_1, \dots, \hat{\theta}_k$, we find that

$$\Lambda_{n^\ell}^* = -\frac{1}{2} \sum_{i=1}^\ell L_i''(t_i) Z_{mi}^2,$$

where

$$Z_{mi} = \sqrt{m} \hat{\theta}_i', \quad 1 \leq i \leq k,$$

and t_1, \dots, t_k are intermediate points between $\hat{\theta}_1, \dots, \hat{\theta}_k$ and $0, \dots, 0$. Let $\epsilon > 0$ be so small that

$$c = \int_{-\infty}^{\infty} \sup_{|t| < \epsilon} \{-H''(x - t)\} g(x) dx \leq \frac{9}{8} = \frac{9}{8} \mathcal{J};$$

and let

$$A = \{|\hat{\theta}_i| \leq \epsilon, \text{ for all } i = 1, \dots, k\}.$$

Then

$$P(A') \leq Ck\rho^m,$$

which tends to zero as $k \rightarrow \infty$ since $\log k = o(m)$. Next, the terms $L_1''(t_1), \dots, L_k''(t_k)$ are bounded. For $i = 1, \dots, k$, let

$$V_{mi} = \frac{1}{m} \sum_{j=1}^m \sup_{|t| \leq \varepsilon} -H''(X_{ij} - t), \quad W_{mi} = V_{mi} Z_{mi}^2 \cdot I_{\{|\hat{\theta}_i| \leq \varepsilon\}}.$$

Then

$$\Lambda_{n, \ell}^* I_A \leq \frac{1}{2} \sum_{i=1}^{\ell} W_{mi}, \quad 1 \leq \ell \leq k.$$

Now W_{m1}, \dots, W_{mk} are i.i.d. for each $m \geq 1$; W_{m1} converges in distribution to $c\chi_1^2$ as $m \rightarrow \infty$, where $c \leq \frac{1}{8}$, by the Central Limit Theorem and the Law of Large Numbers; and $W_{m1}, m \geq 1$, are uniformly integrable, by Lemma 2 below. In particular, the mean $\mu_m = E(W_{m1})$ converges to c as $m \rightarrow \infty$; so, there is an m_0 for which $\mu_m < \frac{1}{4}$ for all $m \geq m_0$. It follows that

$$(19) \quad P(\max_{\ell_0 \leq \ell \leq k} \Lambda_{n, \ell}^* - \frac{3}{4}\ell \geq 0, A) \leq P\{\max_{\ell_0 \leq \ell \leq k} \ell^{-1} \mid \sum_{i=1}^{\ell} \frac{1}{2}(W_{mi} - \mu_m) \mid \geq \frac{1}{8}\},$$

for all $m \geq m_0$. Finally, a simple adaptation of the proof of the strong Law of Large Numbers shows that the right side of (19) may be made arbitrarily small for all $m \geq m_0$ by taking ℓ_0 sufficiently large. See Lemma 3 below. Since $P(A) \rightarrow 0$ as $k, m \rightarrow \infty$, Condition C4 is satisfied.

LEMMA 2. *Suppose that (i)–(iii) are satisfied; then for any $\gamma < \beta$, $\sup_{m \geq 1} E(W_{m1}) < \infty$.*

PROOF. For sufficiently small $\varepsilon > 0$, there are positive constants C and η and a $\rho, 0 < \rho < 1$, for which

$$P(|\hat{\theta}_1| \leq \varepsilon_1 \mid Z_{m1} \mid > t) \leq C(e^{-\eta t^2} + \rho^m)$$

for $0 < t < \varepsilon\sqrt{m}$; see Woodroffe (1979, page 806). Thus all powers of Z_{mi} are uniformly integrable. The lemma now follows directly from (iii) and Hölder's inequality.

LEMMA 3. *For each $m \geq 1$, let Y_{m1}, \dots, Y_{mk} be i.i.d. random variables, w.r.t. probability measure $P = P_m$, for which*

$$(20) \quad E(Y_{m1}) = 0 \quad \text{and} \quad \sup_{m \geq 1} E \mid Y_{m1} \mid^\alpha < \infty$$

for some $\alpha > 1$. If $k = k_m \rightarrow \infty$ as $m \rightarrow \infty$, then for every $\delta > 0$ there is an integer $\ell_0 \geq 1$ for which

$$(21) \quad P(\max_{\ell_0 \leq \ell \leq k} \ell^{-1} \mid \sum_{i=1}^{\ell} Y_{mi} \mid > \delta) < \delta$$

for all $m \geq 1$.

PROOF. Consider ℓ_0 of the form $\ell_0 = 2^q$, where $q \geq 1$ is an integer; let $r = r_m$ be an integer for which $2^{r-1} < k \leq 2^r$; and let $Y_{mi} = 0$ for $k < i \leq 2^r$. Then, the left side of (21) does not exceed

$$\sum_{j=q}^r P(\max_{\ell \leq 2^j} \mid \sum_{i=1}^{\ell} Y_{mi} \mid > \frac{1}{2}\delta 2^j).$$

By the martingale inequality, Condition (20), and Theorem 2 of Von Bahr and Esseen (1965),

$$P(\max_{\ell \leq 2^j} \mid \sum_{i=1}^{\ell} Y_{mi} \mid > \frac{1}{2}\delta 2^j) \leq (\frac{1}{2}\delta 2^j)^{-\alpha} E \mid \sum_{i=1}^{2^j} Y_{mi} \mid^\alpha$$

and

$$E \mid \sum_{i=1}^{2^j} Y_{mi} \mid^\alpha \leq C \cdot 2^j$$

for all $j = 1, \dots, r$ for some constant C . Thus, the left side of (20) does not exceed $2^a \delta^{-a} C \sum_{j=q}^{\infty} (\frac{1}{2})^{j(a-1)}$, which is independent of m and may be made arbitrarily small by taking q sufficiently large.

REFERENCES

- ABRAMOWITZ, M. and STEGUN, I. A. (1970). *Handbook of Mathematical Functions*. National Bureau of Standards.
- AKAIKE, H. (1974). A new look at statistical model identification. *I.E.E. Trans. Auto Control* **19** 716-723.
- ANDERSON, T. W. (1962). The choice of the degree of a polynomial regression as a multiple decision problem. *Ann. Math. Statist.* **33** 255-265.
- DANIEL, C. and WOOD, F. (1980). *Fitting Equations to Data*. Wiley, New York.
- FELLER, W. (1968). *An Introduction to Probability Theory and its Applications*, vol. 1 (3rd ed.). Wiley, New York.
- FELLER, W. (1966). *An Introduction to Probability Theory and its Applications*, vol. 2. Wiley, New York.
- HINKLEY, D. (1976). Note on model selection and the Akaike criterion. Unpublished manuscript, University of Minnesota.
- HUBER, P. (1973). Robust regression. *Ann. Statist.* **1** 799-821.
- LEHMANN, E. (1959). *Testing Statistical Hypotheses*, Wiley, New York.
- MALLOWS, C. (1964). Choosing variables in a linear regression: A graphical aid. Presented at the Central Regional Meeting of the IMS, Manhattan, Kansas.
- MALLOWS, C. (1973). Some comments on C_p . *Technometrics* **15** 661-675.
- SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461-464.
- SHIBATA, R. (1976). Selection of the order of an autoregressive model by Akaike's information criterion. *Biometrika* **63** 117-126.
- SHIBATA, R. (1980). Asymptotically efficient selection of the order of the model for estimating parameters of a linear process, *Ann. Statist.* **8** 147-164.
- SHIBATA, R. (1981). An optimal selection of regression variables. *Biometrika* **68** 45-54.
- VON BAHR, B. and ESSEEN, G. (1965). Inequalities for the r th absolute moment of a sum of random variables, $1 \leq r \leq 2$. *Ann. Math. Statist.* **36** 299-303.
- WALD, A. (1949). A note on the consistency of the maximum likelihood estimator. *Ann. Math. Statist.* **20** 596-601.
- WOODROOFE, M. (1978). Large deviations of the likelihood ratio statistics with applications to sequential testing. *Ann. Statist.* **6** 72-84.
- WOODROOFE, M. (1979). A one-armed bandit problem with a concomitant variable. *J. Amer. Statist. Assoc.* **74** 799-806.
- YOHAI, V. and MARONNA, R. (1979). Asymptotic behavior of M estimators for the linear model. *Ann. Statist.* **7** 258-268.

DEPARTMENT OF STATISTICS
UNIVERSITY OF MICHIGAN
419 SOUTH STATE STREET
ANN ARBOR, MICHIGAN 48109