

A generalised learning algorithm based on guard zones

A. PATHAK and S.K. PAL

Electronics and Communication Sciences Unit, Indian Statistical Institute, Calcutta 700 035, India

Received 16 August 1985

Revised 8 January 1986

Abstract: An algorithm for learning class parameters using a restricted updating programme is described along with investigation of its convergence for optimum learning. The algorithm is a generalisation of some existing ones which were found to be useful for practical data.

Key words: Guard zone, stochastic convergence, learning, robust estimates.

1. Introduction

An adaptive pattern recognition system can be viewed as a learning machine in which the decision of the system gradually approaches the optimal decision by acquiring necessary information from observed patterns. System performance is improved as a result. In a supervised system, the machine requires an extra source of knowledge, usually of a higher order, for correcting the decision taken by a classifier. When an extra source of knowledge on which a supervisory programme could be based is not readily available, the performance of the system becomes highly unpredictable.

The most widely used tools for recursive learning of class parameters are Bayesian estimation methods and stochastic approximation (Fu, 1968; Tsypkin, 1973). In this context, we would like to single out the self-supervised learning system based on the concept of a 'guard zone', mooted by Pal et al. (1980). It was used to restrict the updating of estimates of parameters (feature means and variances) by means of 'doubtful' samples. For this purpose a guard zone was defined for each class in such a way that a training sample was used for updating only if it fell within the guard zone.

A similar algorithm was presented by Chien (1970) as a solution to the problem of identifying 'spurious', that is, possibly non-representative training samples for the case when feature means are to be learned. A threshold is defined such that if the 'distance' of the current training sample from the preceding estimate of the mean (the same 'distance' is used for defining a guard zone) exceeds it, the training sample is rejected. As such both algorithms are basically the same, and are stochastic approximation procedures of sorts.

Although these algorithms were tested with success on some practical data, the two works did not provide any proof of convergence or any theoretical investigation of the choice of the controlling parameter namely, guard zone dimension or threshold for optimum learning.

The present work describes a generalized version of the two, called the Generalized Guard Zones Algorithm (GGA). Basically, it aims to detect outliers and reject them from the parameter-updating procedure. As such, it can be looked upon as a robust estimation procedure (Andrews, 1975; Huber, 1981). (Essentially, the term 'robustness' signifies insensitivity to small deviations from the underlying assumptions.)

1981)). Under rather general assumptions, we have investigated the stochastic convergence of this algorithm for some special problems of estimation. For this purpose, results on stochastic approximation are used.

2. The Generalized Guard Zones Algorithm (GGA)

Let $X = [X_1, X_2, \dots, X_N]'$, $X \in \mathbb{R}^N$ be an N -dimensional feature vector defined over a pattern class C .

Let us make the following assumptions:

- (A1) The distribution of X over C is continuous.
 (A2) This distribution depends on a q -dimensional parameter vector θ , some or all of which need to be learned.
 (A3) The distribution of X over C is such that $E(X)$ exists and is equal to μ .
 (A4) The dispersion matrix of X , namely,

$$\text{Disp}(X) = \Sigma = (\sigma_{ij}) \text{ exists.}$$

Before stating the algorithm itself, let us define a guard zone formally as follows:

Definition. Let S be a metric space and δ a metric defined on it. Then for any point $a \in S$, a guard zone $G(a, \lambda)$ having an 'extent' λ is the subset of S defined by

$$G(a, \lambda) = \{x: \delta(a, x) \leq \lambda\},$$

where $\lambda \geq 0$.

Clearly, $G(a, \lambda)$ is nothing but a closed ball of radius λ centred at a in S with respect to the metric δ .

In the subsequent discussions we shall be taking $S = \mathbb{R}^N$ and a metric d defined as

$$d^2(x, y) = (x - y)' A (x - y), \quad x, y \in \mathbb{R}^N,$$

A being a symmetric, positive definite matrix.

Let us now proceed to the algorithm properly.

Let X_1, X_2, X_3, \dots be the sequence of learning (or training) samples, randomly selected from C , that is, assumed to be independently and identically distributed, i.e., we assume that correctly labelled training samples are available.

We restrict ourselves to the case where θ includes μ and/or elements of Σ only.

The generalized guard zones algorithm (GGA) for estimating θ recursively is as follows:

$$\theta_k = \begin{cases} f(X_k) & \text{for } k = 1, \\ \theta_{k-1} - a_k Y_k & \text{for } k > 1; \end{cases} \quad (1)$$

$$Y_k = \begin{cases} \theta_{k-1} - f(X_k) & \text{if } X_k \in G(m_{k-1}, \lambda_k), \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

θ_k : the k -th-stage estimate of θ .

$\{a_k\}$: a sequence of positive numbers, with $a_k \leq 1 \forall k$.

f : $\mathbb{R}^N \rightarrow \mathbb{R}^q$ is a continuous mapping, defining an unbiased statistic for θ .

m_{k-1} : the $(k-1)$ -th stage GGA estimate of μ .

$$G(m_{k-1}, \lambda_k) = \{x: x \in \mathbb{R}^N, d_k(x, m_{k-1}) \leq \lambda_k\},$$

$$d_k^2(x, y) = (x - y)' A_k (x - y),$$

A_k : A symmetric, positive definite matrix, which may or may not be a function of X , and/or θ , $i = 1(1)k$.

λ_k : A positive number, prespecified.

In essence, this algorithm uses only those training samples for updating the estimate, which lie within the corresponding guard zone centred at the preceding estimate of the mean. Training samples which lie outside it are ignored and the estimate kept unchanged at the corresponding stages.

Special cases

(1) When $q = N$, $\theta = \mu$, i.e., only the mean vector is to be learned, we have $f(X) = X$, and the algorithm is as follows:

$$\theta_1 = m_1 = X_1,$$

and for $k > 1$, with $\theta_k = m_k$,

$$m_k = \begin{cases} m_{k-1} - a_k (m_{k-1} - X_k) & \text{if } d_k(m_{k-1}, X_k) \leq \lambda_k, \\ m_{k-1} & \text{otherwise.} \end{cases} \quad (3)$$

(2) For estimating σ_{ij} , $i = 1(1)N$, $j = 1(1)N$, recursively, there are two alternative procedures (we write $\hat{\sigma}_{ij(k)} = s_{ij}^{(k)}$):

(a) $s_{ij}^{(1)} = 0$, and for $k > 1$,

$$s_{ij}^{(k)} = \begin{cases} s_{ij}^{(k-1)} - a_k [s_{ij}^{(k-1)} - (X_{ki} - m_{(k-1)}) (X_{kj} - m_{(k-1)})] \\ \text{if } d_k(m_{n-1}, X_k) \leq \lambda_k, \\ s_{ij}^{(k-1)} \text{ otherwise.} \end{cases} \quad (4a)$$

(b) $s_{ij}^{(1)} = 0$, and for $k > 1$,

$$s_{ij}^{(k)} = c_{ij}^{(k)} - m_{(k-1)}, m_{(k-1)}, \quad (4b)$$

where

$$c_{ij}^{(k)} = \begin{cases} c_{ij}^{(k-1)} - a_k [c_{ij}^{(k-1)} - X_{ki} X_{kj}] \\ \text{if } d_k(m_{n-1}, X_k) \leq \lambda_k, \\ c_{ij}^{(k-1)} \text{ otherwise.} \end{cases} \quad (5)$$

(3) When $\theta = [\mu' : \sigma']'$, $\sigma' = [\sigma_{11} \sigma_{22} \dots \sigma_{NN}]$, $q = 2N$,

$$A_k = \begin{pmatrix} s_{11} & 0 & 0 & \dots & 0 \\ 0 & s_{22} & 0 & \dots & 0 \\ 0 & 0 & s_{33} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & s_{NN} \end{pmatrix}^{-1} \\ = [\text{Diag}(s_{11}, s_{22}, \dots, s_{NN})]^{-1}, \\ a_k = \frac{1}{k}, \quad \lambda_k = \frac{1}{\lambda},$$

then the GGA reduces to the algorithm of Pal et al. (1980).

(4) When $\theta = \mu$, $q = N$,

$A_k = K_1$, the dispersion in X due to ordinary measurement variation, i.e., type I noise (see Chien (1970)),

$\lambda_k = \theta_k$ (the threshold in Chien's algorithm),

$u_k = [(k-1) + \nu]^{-1}$ where ν is such that

$$\text{Disp}(\mu_0) = K_1 / \nu,$$

μ_0 = being the initial estimate, $\nu > 0$,

the GGA reduces to the non-linear learning algorithm of Chien (1970).

(5) As all λ_k 's decrease progressively, the system approaches the nonadaptive state. Clearly, this is because the dimensions of the guard zones decrease and hence the probability of a training sample to fall within the guard zone decreases with decrease in the corresponding λ_k -value; so the number of training samples getting selected for the updating process decreases.

(6) On the other hand, when the λ_k 's increase progressively, the system approaches the non-supervised state for, as the 'extent' of the guard zones increases, more and more training samples get selected for the updating of estimates; that is, the updating programme becomes less and less restrictive.

3. Convergence of the generalized guard zones algorithm

The convergence of a recursive discrete algorithm for estimating a parameter θ by θ_n , can be defined in various ways. For instance, we say that

(i) the sequence $\{\theta_n\}$ converges to θ with probability one or almost surely if

$$P \left[\lim_{n \rightarrow \infty} \|\theta_n - \theta\| = 0 \right] = 1,$$

P being the probability measure.

(ii) $\{\theta_n\}$ converges to θ in the mean-square sense if

$$\lim_{n \rightarrow \infty} E[\|\theta_n - \theta\|^2] = 0,$$

E being the expectation operator.

For proving certain results on the convergence of the GGA, we shall be making use of the following results:

Theorem 1 (Schmetterer, 1968). Let $\{a_n\}$ be a sequence of positive real numbers such that

$$(B1) \quad \sum_{n=1}^{\infty} a_n^2 < \infty.$$

Let x_n and y_n be K -dimensional random vectors which satisfy

$$(B2) \quad x_{n+1} = x_n - a_n y_n, \quad n \geq 1.$$

Let M_n be a measurable mapping from \mathbb{R}^K to \mathbb{R}^K , such that

$$(B3) \quad E(y_n | x_1, x_2, \dots, x_n) = M_n(x_n) \text{ a.e.}$$

Let a, b, c be nonnegative real numbers and let

$$(B4) \quad E(\|y_n\|^2 | x_1, x_2, \dots, x_n) \leq a + b|x_n| + c|x_n|^2 \text{ a.e.}$$

Also, for every $x \in \mathbb{R}^K$ and $n \geq 1$,

(B5) $x' M_n(x) \geq 0$.

If x_1 is chosen in such a way that

(B6) $E\{|x_1|^2\}$ exists,

then the sequence $\{x_n\}$ converges with probability 1 (that is, almost surely) and the sequence $\{E|x_n|^2\}$ converges also.

Theorem 2 (Schmetterer, 1968). Suppose that assumptions (B1)-(B6) hold. If there exists for every $\eta > 0$ a $\delta > 0$ such that for $n \geq 1$

(B7) $\inf_{n < |x| < n+1} |x' M_n(x)| \geq \delta$,

then $\{x_n\}$ converges to the k -dimensional null vector θ almost surely.

We shall now be proving the following:

Proposition 1. For the problem of estimating $\theta = \mu$ recursively, let $\hat{\theta}_k = m_k$ be the sequence of estimates, where m_k is given by equation (3).

If

(C1) $\sum_{n=1}^{\infty} a_n^2 < \infty$

and

(C2) $P\{d_k(m_{k-1}, X_k) \leq \lambda_k | m_{k-1}\} > \delta, \forall k$,
for some $\delta > 0$,

then

- (a) $\{m_k\}$ converges with probability 1 to μ as $k \rightarrow \infty$;
- (b) $\{E|m_k - \mu|^2\}$ converges as $k \rightarrow \infty$.

Proposition 2. Consider the problem of estimating μ and Σ completely.

Let

$\theta' = [\mu' ; \sigma_*^{(1)} ; \sigma_*^{(2)} ; \dots ; \sigma_*^{(N)}]$ (6a)

and

$\hat{\theta}'_k = [m'_k ; c_k^{(1)} ; c_k^{(2)} ; \dots ; c_k^{(N)}]$, (6b)

where

$q = \frac{N(N+1)}{2}$,

$\sigma_*^{(i)} = [\sigma_{ii}^* \sigma_{i(i+1)}^* \dots \sigma_{iN}^*]$,

$q_i = N - i + 1$,

$\sigma_{ij}^* = E(X_i X_j) = \sigma_{ij} + \mu_i \mu_j$,

and the elements of m_k and $c_k^{(i)} = [c_{ii}^{(i)} c_{i(i+1)}^{(i)} \dots c_{iN}^{(i)}]$ are given by equations (3) and (5) respectively.

If the conditions (C1) and (C2) hold, and besides, the following condition is also satisfied:

(C3) $\eta_{ij} = E(X_i^2 X_j^2)$ exists $\forall i, j = 1(1)N$ for the elements of the feature vector X ,

then

- (a) $\{\hat{\theta}_k\}$ (given by (6b)) converges with probability 1 to θ (given by (6a)) as $k \rightarrow \infty$;
- (b) $\{E|\hat{\theta}_k - \theta|^2\}$ converges as $k \rightarrow \infty$.

Proposition 3. Consider the problem of estimating μ and Σ completely.

Let

$\theta = [\mu' ; \sigma^{(1)} ; \dots ; \sigma^{(N)}]'$ (7a)

and

$\hat{\theta}_k = [m'_k ; s^{(1)} ; s^{(2)} ; \dots ; s^{(N)}]'$, (7b)

where

$\sigma^{(i)} = |\sigma_{ii} \sigma_{i(i+1)} \dots \sigma_{iN}|'$,

with $\sigma_{ij} = E(X_i X_j) - \mu_i \mu_j$, and

$s^{(i)} = [s_{ii}^{(i)} s_{i(i+1)}^{(i)} \dots s_{iN}^{(i)}]'$,

the s_{ij} 's being given by equation (4b).

If the conditions (C1), (C2) and (C3) hold, then $\hat{\theta}_k$ (given by eqn. (7b)) converges with probability 1 to θ (given by eqn. (7a)).

Proofs of Propositions

Before we give the formal proof of the propositions, we would like to make the following point

If we subtract θ from both sides of the equation (1) we get, writing

$\hat{\theta}_k^* = \hat{\theta}_k - \theta$,

$\hat{\theta}_k^* = \begin{cases} f(X_k) - \theta, & \text{for } k = 1, \\ \{\hat{\theta}_{k-1}^* - a_k Y_k^*, & \text{for } k > 1. \end{cases}$ (8)

where

$$Y_k^* = \begin{cases} \theta_{k-1}^* - (f(X_k) - \theta) \\ \text{if } X_k \in G(m_{k-1}, \lambda_k), \\ 0 \text{ otherwise.} \end{cases} \quad (9)$$

This is because $0 < a_k \leq 1 \forall k$ by choice, so that for $k > 1$,

$$\begin{aligned} \theta_k^* &= \theta_k - \theta = \theta_{k-1} - a_k Y_k - \theta \\ &= \begin{cases} \theta_{k-1} - a_k (\theta_{k-1} - f(X_k)) - \theta \\ \text{if } X_k \in G(m_{k-1}, \lambda_k), \\ \theta_{k-1} - \theta \text{ otherwise,} \end{cases} \\ &= \begin{cases} (1 - a_k) [\theta_{k-1} - \theta] + a_k [f(X_k) - \theta] \\ \text{if } X_k \in G(m_{k-1}, \lambda_k), \\ \theta_{k-1}^* \text{ otherwise,} \end{cases} \\ &\quad \text{writing } \theta = (1 - a_k)\theta + a_k\theta. \\ &= \begin{cases} \theta_{k-1}^* - a_k [\theta_{k-1}^* - \{f(X_k) - \theta\}] \\ \text{if } X_k \in G(m_{k-1}, \lambda_k), \\ \theta_{k-1}^* \text{ otherwise.} \end{cases} \end{aligned}$$

Thus Propositions 1 and 2 can be shown to be true if we can show that under the conditions assumed therein,

- (i) $\theta_k^* \rightarrow 0$ almost surely as $k \rightarrow \infty$;
 (ii) $E\{|\theta_k^*|^2\}$ converges as $k \rightarrow \infty$,

To establish these we shall apply Theorems 1 and 2 directly, by showing that conditions (B1)-(B7) are true for θ_k^* as defined by equations (8) and (9).

The conditions (B1) and (B2) are true because of our assumption (C1) and equations (8) and (9) respectively.

Proof of Proposition 1. Here,

$$\begin{aligned} E\{Y_k^* | \theta_1^*, \theta_2^*, \dots, \theta_k^*\} \\ &= p_{k+1} E\{\theta_k^* - (X_{k+1} - \theta) | \theta_1^*, \theta_2^*, \dots, \theta_k^*\} \\ &\quad (\text{as } f(x) = x \text{ here}) \\ &= p_{k+1} \{\theta_k^* - E(x_{k+1} - \theta)\} \\ &\quad (\text{as } X_{k+1} \text{ is independent of} \\ &\quad X_1, X_2, \dots, X_k \text{ and hence } \theta_1^*, \dots, \theta_k^*) \\ &= p_{k+1} \theta_k^* \\ &\quad (\text{as } \theta = \mu \text{ here and } E(X_{k+1}) = \mu). \end{aligned}$$

This verifies condition (B3), with

$$M_k(x) = p_{k+1}x, \quad x \in \mathbb{R}^N.$$

Further,

$$\begin{aligned} E\{|\theta_k^*|^2 | \theta_1^*, \theta_2^*, \dots, \theta_k^*\} \\ &= p_{k+1} E\{[\theta_k^* - (X_{k+1} - \theta)]^2 | \theta_1^*, \dots, \theta_k^*\} \\ &= p_{k+1} \{|\theta_k^*|^2 - 2\theta_k^* E(X_{k+1} - \theta) + E\{X_{k+1} - \theta\}^2\} \\ &\quad (\text{as } X_{k+1} \text{ is independent of } \theta_1^*, \dots, \theta_k^*) \\ &= p_{k+1} \left[|\theta_k^*|^2 + E \sum_{n=1}^N (X_{(k+1)_n} - \mu_n)^2 \right] \\ &\quad (\text{as } E(X_{k+1} - \theta) = 0) \\ &= p_{k+1} \left[|\theta_k^*|^2 + \sum_{n=1}^N \theta_n^2 \right] \\ &\leq |\theta_k^*|^2 + \sum_{n=1}^N \sigma_n^2 \quad (\text{as } p_{k+1} \leq 1 \text{ and } a_k^2 \leq a_k \leq 1), \end{aligned}$$

which means that condition (B4) holds with

$$a = \sum_{n=1}^N \sigma_n^2, \quad b = 0, \quad c = 1.$$

Also, $x'M_k(x) = p_{k+1}x'x \geq 0 \forall x \in \mathbb{R}^N$ which verifies condition (B5). That (B6) holds, is rather obvious, as

$$E\{f(X_1) - \theta\}^2 = E\{X_1 - \mu\}^2 = \sum_{n=1}^N \sigma_n^2 < \infty.$$

Finally, we see that by virtue of our assumption (C2), the condition (B7) holds, for

$$\begin{aligned} \inf_{\eta < |x| < \eta^{-1}} \{x'M_k(x)\} \\ &= \inf_{\eta < |x| < \eta^{-1}} \{p_{k+1}x'x\} > \delta\eta^2. \end{aligned}$$

Thus Theorems 1 and 2 hold for θ_k^* . Hence the proposition is proved.

Proof of Proposition 2. Here

$$f(x) = \begin{bmatrix} x^{(0)'} \\ x^{(1)'} \\ \vdots \\ x^{(N)'} \end{bmatrix}'$$

with

$$\begin{aligned} x^{(0)} &= x, \\ x_{q \times 1}^{(i)} &= [x_i^2, x_i x_{i+1}, \dots, x_i x_N]', \\ q_i &= N - i + 1. \end{aligned}$$

Thus

$$\begin{aligned} E\{Y_k^* | \theta_1^*, \theta_2^*, \dots, \theta_k^*\} &= E\{\theta_k^* - (f(X_{k,1}) - \theta^*) | \theta_1^*, \theta_2^*, \dots, \theta_k^* | \rho_{k,1}\} \\ &= \rho_{k,1} \{\theta_k^* - E\{f(X_{k,1}) - \theta\}\} \\ &\quad (\text{for the same reasons as before}) \\ &= \rho_{k,1} \theta_k^* \quad (\text{as } E\{f(X_{k,1})\} = \theta). \end{aligned}$$

This verifies (B3). Now,

$$\begin{aligned} E\{\|Y_k^*\|^2 | \theta_1^*, \theta_2^*, \dots, \theta_k^*\} &= \rho_{k,1} \{ \|\theta_k^* - (f(X_{k,1}) - \theta)\|^2 | \theta_1^*, \theta_2^*, \dots, \theta_k^* \} \\ &= \rho_{k,1} \{ \|\theta_k^*\|^2 - 2\theta^{*T} (E\{f(X_{k,1}) - \theta\}) \\ &\quad + E\{[f(X_{k,1}) - \theta]^2\} \} \quad (\text{as before}) \\ &= \rho_{k,1} \{ \|\theta_k^*\|^2 + E\{[f(X_{k,1}) - \theta]^2\} \} \\ &\quad (\text{as } E\{f(X_{k,1})\} = \theta). \end{aligned}$$

However, as

$$\begin{aligned} E\{\|f(X_{k,1})\|^2\} &= E\left\{ \sum_{i=0}^N \|X_{k,i}^{(i)}\|^2 \right\} = \sum_{n=0}^N \{E\|X_{k,i}^{(n)}\|^2\} \\ &= \sum_{n=1}^N E(X_n^2) + \sum_{j=1}^N \sum_{i=j}^N E(X_j^2 X_i^2) \\ &\quad \left(\text{since } E\{X_{k,i}^{(n)}\} = E\left\{ \sum_{n=1}^N X_n^2 \right\} \right) \\ &\quad \text{and } E\{X_{k,i}^{(i)}\} = E\left\{ \sum_{i=1}^N X_i^2 X_i^2 \right\} \\ &= \sum_{n=1}^N (\sigma_n^2 + \mu_n^2) + \sum_{j=1}^N \sum_{i=j}^N \eta_{ij} \\ &\quad (\text{by assumption (C3)}) \\ &\leq K, \end{aligned}$$

with K a finite positive constant independent of $\theta_1^*, \theta_2^*, \dots, \theta_k^*$, we must have

$$\begin{aligned} E\{\|Y_k^*\|^2 | \theta_1^*, \theta_2^*, \dots, \theta_k^*\} &\leq \rho_{k,1} \{ \|\theta_k^*\|^2 + K \}, \quad (\text{as } E\{\|f(X_{k,1}) - \theta\|^2\}) \\ &\leq \|\theta_k^*\|^2 + K \\ &= E\{[f(X_{k,1}) - \theta]^2\} - 2\theta^{*T} E\{f(X_{k,1}) - \theta\} + \|\theta\|^2 \\ &= E\{[f(X_{k,1}) - \theta]^2\} - \|\theta\|^2 \\ &\quad (\text{as } E\{f(X_{k,1})\} = \theta) \end{aligned}$$

$$\leq E\{f(X_{k,1})\}^2.$$

Thus condition (B4) is satisfied, with $a = K$, $b = 0$, $c = 1$. Further, conditions (B5), (B6) and (B7) also can be seen to be true, as

$$\begin{aligned} x^T M_k(x) &= \rho_{k,1} x^T x \geq 0, \quad \forall x \in \mathbb{R}^d, \\ E\{[f(X_{k,1}) - \theta]^2\} &\leq E\{f(X_{k,1})\}^2 < K < \infty, \end{aligned}$$

as seen above, and

$$\inf_{\eta < |x| < \eta} \{x^T M_k(x)\} > \delta \eta^2 > 0,$$

because of our assumption (C2).

Thus Theorems 1 and 2 hold for $\{\theta_k^*\}$. This completes the proof of Proposition 2.

Proof of Proposition 3. The proof follows directly from Proposition 2 and the following lemma, if we note that $s_{ij}^{(k)}$ as given by eqn. (4b), is a continuous function, say, g_{ij} of θ_{k-1} (given by eqn (6a)), where $g_{ij}(x)$, $x \in \mathbb{R}^d$ is defined as

$$g_{ij}(x) = x_i - x_i x_j$$

where

$$r_{ij} = \begin{cases} N + j & \text{if } i = 1, \\ N + \sum_{t=1}^{i-1} (N - t) + (j - i + 1) & \text{if } 1 < i \leq N. \end{cases}$$

and x_k denotes the k -th element of x , $k = 1(1)q$. Further, $g_{ij}(\theta) = \sigma_{ij}^* - \mu_i \mu_j = \sigma_{ij}$. Hence Proposition 3 follows.

Lemma. Let $\{X_n\}$ be a sequence of random variables taking values in \mathbb{R}^p and let $g: \mathbb{R}^p \rightarrow \mathbb{R}^q$ be a continuous map. Then $\{X_n\}$ converges with probability 1 to a ($a \in \mathbb{R}^p$) implies that $g(X_n)$ converge with probability 1 to $g(a)$ ($a \in \mathbb{R}^p$).

References

Andrews, D.F. et al. (1972). *K-Nearest Neighbors: Estimates of Location, Survey and Advances*. Princeton University Press, Princeton, NJ.
 Chien, Y.T. (1970). The threshold effect of a nonlinear learning algorithm for pattern recognition. *Information Science*: 351-358.
 Fu, K.S. (1968). *Sequential Methods in Pattern Recognition and Machine Learning*. Academic Press, New York.
 Huber, P.J. (1981). *Robust Statistics*. Wiley, New York.

Pal, S.K., A.K. Datta and D. Dutta Majumder (1980). A self-supervised and vowel recognition system. *Patt. Recogn.* 12, 27-34.

Schmetterer, L. (1968). Multidimensional stochastic approximation. In: P.R. Krishnaiah, Ed., *Multivariate Analysis-II:*

Proc. Second Int. Symp. Multivariate Anal. Dayton, Ohio. Academic Press, New York.

Tsybkin, Ya.Z. (1973). *Foundations of the Theory of Learning Systems.* Academic Press, New York.