

ON A SUCCESSIVE TRANSFORMATION OF PROBABILITY DISTRIBUTION AND ITS APPLICATION TO THE ANALYSIS OF THE OPTIMUM GRADIENT METHOD

BY HIROTUGU AKAIKE

(Received July 25, 1959)

§ 0. Introduction and summary.

In this paper we define a type of transformation of probability distribution and analyze the limiting behavior of the result of successive applications of the transformation to some initial probability distribution. By using the results of this analysis we can get a fairly general insight into the so-called optimum-gradient method in numerical analysis. We can prove the conjecture which was stated by Forsythe and Motzkin [7] and was used as the logical basis of an acceleration procedure for the optimum gradient method [4][5][6]. It was stated by Forsythe [4] that this conjecture seems to be hard to prove as the related transformation is rather complicated. But our present proof is rather simple. Further, we can see the relation between the condition-number of the related matrix and the convergence rate of the optimum gradient method. By using the relation which according to [5] is first proved by Kantorovich [8], we can say that when the matrix is ill-conditioned the convergence rate tends near to its worst possible value. Using the same data as those treated by Forsythes in paper [5], we give some numerical examples.

There are many important problems, not necessarily of linear type, where the gradient method is applicable [2], [3]. Even in these non-linear case we can expect that when the approximation proceeds the problem will be reduced essentially to the linear one. One of such examples was discussed in the former paper [1] of the present author. Thus, though the results of this paper are concerned with the solution of the simultaneous linear equation $Ax = b$, these results will be generally useful for the analysis of limiting behavior of the approximate solutions in the optimum gradient method.

§ 1. A successive transformation of probability distribution.

In this section we shall treat probability distributions over a set of

mutually different real numbers $\lambda_1, \lambda_2, \dots, \lambda_n$. First we shall give the definition of our transformation T . For a probability distribution $\mathbf{P} \equiv \{(\mathbf{P})_\nu; \nu=1, 2, \dots, n\}$ where $(\mathbf{P})_\nu$'s represent the probabilities attached to λ_ν 's respectively the image $T\mathbf{P} \equiv \{(T\mathbf{P})_\nu; \nu=1, 2, \dots, n\}$ of \mathbf{P} by T is given by

$$(T\mathbf{P})_\nu = \frac{(\mathbf{P})_\nu (\lambda_\nu - \bar{\lambda}(\mathbf{P}))^2}{\sum_{\mu=1}^n (\mathbf{P})_\mu (\lambda_\mu - \bar{\lambda}(\mathbf{P}))^2} \quad \nu=1, 2, \dots, n$$

where $\bar{\lambda}(\mathbf{P}) = \sum_{\mu=1}^n \lambda_\mu (\mathbf{P})_\mu$ and T has as its domain the set of \mathbf{P} 's with $\sum_{\mu=1}^n (\mathbf{P})_\mu (\lambda_\mu - \bar{\lambda}(\mathbf{P}))^2 > 0$. Here we shall give some preliminary lemmas. We shall hereafter use the notation $\bar{f}(\bar{\lambda})(\mathbf{P})$ for $\sum_{\mu=1}^n f(\lambda_\mu) (\mathbf{P})_\mu$.

LEMMA 1. *We have*

$$\bar{\lambda}(T\mathbf{P}) = \bar{\lambda}(\mathbf{P}) + \frac{\overline{(\lambda - \bar{\lambda}(\mathbf{P}))^2}(\mathbf{P})}{(\lambda - \bar{\lambda}(\mathbf{P}))^2(\mathbf{P})}.$$

PROOF.

$$\begin{aligned} \bar{\lambda}(T\mathbf{P}) &= \sum_{\nu=1}^n \lambda_\nu (T\mathbf{P})_\nu = \frac{\sum_{\nu=1}^n \lambda_\nu (\lambda_\nu - \bar{\lambda}(\mathbf{P}))^2 (\mathbf{P})_\nu}{(\lambda - \bar{\lambda}(\mathbf{P}))^2(\mathbf{P})} \\ &= \frac{\sum_{\nu=1}^n (\lambda_\nu - \bar{\lambda}(\mathbf{P}))^2 (\mathbf{P})_\nu + \bar{\lambda}(\mathbf{P}) \sum_{\nu=1}^n (\lambda_\nu - \bar{\lambda}(\mathbf{P}))^2 (\mathbf{P})_\nu}{(\lambda - \bar{\lambda}(\mathbf{P}))^2(\mathbf{P})} \\ &= \frac{\overline{(\lambda - \bar{\lambda}(\mathbf{P}))^2}(\mathbf{P})}{(\lambda - \bar{\lambda}(\mathbf{P}))^2(\mathbf{P})} + \bar{\lambda}(\mathbf{P}). \end{aligned}$$

LEMMA 2. *We have*

$$\overline{(\lambda - \bar{\lambda}(T\mathbf{P}))^2}(T\mathbf{P}) \geq \overline{(\lambda - \bar{\lambda}(\mathbf{P}))^2}(\mathbf{P})$$

where $=$ holds only when there are only two ν 's with $(\mathbf{P})_\nu > 0$.

Note; We are considering only \mathbf{P} 's with $\overline{(\lambda - \bar{\lambda})(\mathbf{P}))^2}(\mathbf{P}) > 0$, and the case where there is only one ν with $(\mathbf{P})_\nu > 0$ is out of our consideration. The lemma states that by transformation T the variance of the distribution increases except for the special type of distributions stated in the lemma. Thus for \mathbf{P} in the domain of T , $T\mathbf{P}$ is again in the domain of T .

PROOF.

By using the result of lemma 1 we have

$$\lambda_\nu - \bar{\lambda}(TP) = \lambda_\nu - \bar{\lambda}(P) + \{\bar{\lambda}(P) - \bar{\lambda}(TP)\} = \lambda_\nu - \bar{\lambda}(P) - \frac{(\lambda - \bar{\lambda}(P))^2(P)}{(\lambda - \bar{\lambda}(P))^2(P)}.$$

Thus we get

$$\begin{aligned} & \overline{(\lambda - \bar{\lambda}(TP))^2(TP)} - \overline{(\lambda - \bar{\lambda}(P))^2(P)} \\ &= \frac{\sum_{\nu=1}^n (\lambda_\nu - \bar{\lambda}(TP))^2 (\lambda_\nu - \bar{\lambda}(P))^2(P)_\nu - \{(\lambda - \bar{\lambda}(P))^2(P)\}^2}{(\lambda - \bar{\lambda}(P))^2(P)} \\ &= \frac{\overline{(\lambda - \bar{\lambda}(P))^4(P)} \overline{(\lambda - \bar{\lambda}(P))^2(P)} - \{(\lambda - \bar{\lambda}(P))^2(P)\}^2 - \overline{(\lambda - \bar{\lambda}(P))^2(P)}^3}{\{(\lambda - \bar{\lambda}(P))^2(P)\}^2}. \end{aligned}$$

For the sake of simplicity we shall use here the abbreviated notation $M_k(P)$ in place of $\overline{(\lambda - \bar{\lambda}(P))^k(P)}$. Then we have

$$\begin{aligned} & \overline{(\lambda - \bar{\lambda}(P))^4(P)} \overline{(\lambda - \bar{\lambda}(P))^2(P)} - \{(\lambda - \bar{\lambda}(P))^2(P)\}^2 - \overline{(\lambda - \bar{\lambda}(P))^2(P)}^3 \\ &= \begin{vmatrix} 1 & 0 & M_2(P) \\ 0 & M_2(P) & M_3(P) \\ M_2(P) & M_3(P) & M_4(P) \end{vmatrix}. \end{aligned}$$

If we represent by λ the random variable which follows the probability distribution P i. e. $\text{prob}\{\lambda = \lambda_\nu\} = (P)_\nu$, $\nu = 1, 2, \dots, n$, the above stated determinant is the determinant of the product moment matrix of the random variables $(\lambda - E(\lambda))^0 \equiv 1$, $\lambda - E(\lambda)$ and $(\lambda - E(\lambda))^2$. The product moment matrix is positive semi-definite and thus the value of the above determinant is non-negative. The value of the determinant is equal to zero if and only if there is non trivial linear relation between 1, $\lambda - E(\lambda)$ and $(\lambda - E(\lambda))^2$. Suppose that there exist constants $\alpha_0, \alpha_1, \alpha_2$ not all equal to zero satisfying the relation $\alpha_0 1 + \alpha_1(\lambda - E(\lambda)) + \alpha_2(\lambda - E(\lambda))^2 = 0$ with probability 1. Such a relation can hold only when there are not more than two values of ν for which $\text{prob}\{\lambda = \lambda_\nu\} = (P)_\nu > 0$. Taking into account the fact that we have excluded the case where λ is identically equal to a fixed constant, we get the proof of the final statement of the lemma.

Using the results of these lemmas we want to analyze the limiting behavior of $P^{(k)} \equiv T^k P^{(0)}$ as k tends to infinity where $T^k P^{(0)}$ denotes the result of k successive applications of the transformation T to $P^{(0)}$ and the initial probability distribution $P^{(0)}$ is supposed to satisfy the condition $M_2(P^{(0)}) > 0$.

It is obvious that any infinite subsequence of the sequence $\{P^{(k)}; k=1, 2, \dots\}$ contains its own convergent subsequence. For a convergent subsequence $\{P^{(\alpha_j)}; j=1, 2, \dots\}$ of original $\{P^{(k)}\}$ we shall represent by $P_\alpha^{(\infty)}$ its limiting distribution.

Now as $(\lambda - \bar{\lambda}(P^{(\alpha_j)}))^2(P^{(\alpha_j)}) = \sum_{\nu=1}^n (\lambda_\nu - \bar{\lambda}(P^{(\alpha_j)}))^2(P^{(\alpha_j)})$ and $\bar{\lambda}(P^{(\alpha_j)}) = \sum_{\nu=1}^n \lambda_\nu(P^{(\alpha_j)})$, it can be seen that

$$\lim_{j \rightarrow \infty} \bar{\lambda}(P^{(\alpha_j)}) = \bar{\lambda}(P_\alpha^{(\infty)})$$

$$\lim_{j \rightarrow \infty} \overline{(\lambda - \bar{\lambda}(P^{(\alpha_j)}))^2(P^{(\alpha_j)})} = \overline{(\lambda - \bar{\lambda}(P_\alpha^{(\infty)}))^2(P_\alpha^{(\infty)})}.$$

Further by taking into account the result of lemma 2 we can see that $\overline{(\lambda - \bar{\lambda}(P^{(\alpha_j)}))^2(P^{(\alpha_j)})}$ tends monotone-increasingly to $\overline{(\lambda - \bar{\lambda}(P_\alpha^{(\infty)}))^2(P_\alpha^{(\infty)})}$ and also that

$$\overline{(\lambda - \bar{\lambda}(P_\alpha^{(\infty)}))^2(P_\alpha^{(\infty)})} = \lim_{k \rightarrow \infty} \overline{(\lambda - \bar{\lambda}(P^{(k)}))^2(P^{(k)})} = \overline{(\lambda - \bar{\lambda}(TP_\alpha^{(\infty)}))^2(TP_\alpha^{(\infty)})}$$

holds.

Now as $M_2(P^{(0)}) > 0$, we have $M_2(P_\alpha^{(\infty)}) > 0$, and we have seen that the equality $M_2(TP_\alpha^{(\infty)}) = M_2(P_\alpha^{(\infty)})$ holds. Thus from the result of lemma 2 we can see that there are just two numbers $\nu(\alpha, 1)$ and $\nu(\alpha, 2)$ for which $(P_\alpha^{(\infty)})_{\nu(\alpha, 1)} > 0$, $(P_\alpha^{(\infty)})_{\nu(\alpha, 2)} > 0$ and $(P_\alpha^{(\infty)})_{\nu(\alpha, 1)} + (P_\alpha^{(\infty)})_{\nu(\alpha, 2)} = 1$ hold. Thus it has been proved that the limit points of the sequence $\{P^{(k)}\}$ are contained in the set of points or distributions of the type analogous to the $P_\alpha^{(\infty)}$ which we have obtained above.

Now we shall further prove that the sequence $\{P^{(k)}\}$ is itself oscillatorily convergent in the sense which will be described bellow.

Taking into account of the monotone-increasing property of the sequence $M_2(P^{(k)})$ we have seen that any $P_\alpha^{(\infty)}$ has one and the same variance $M_2(P_\alpha^{(\infty)}) = \lim_{k \rightarrow \infty} M_2(P^{(k)})$. Now the distribution $P_\alpha^{(\infty)}$ is characterized by the quantities $\nu(\alpha, 1)$, $\nu(\alpha, 2)$, $(P_\alpha^{(\infty)})_{\nu(\alpha, 1)}$, $(P_\alpha^{(\infty)})_{\nu(\alpha, 2)}$ and by the relation $M_2(P_\alpha^{(\infty)}) = (P_\alpha^{(\infty)})_{\nu(\alpha, 1)}(P_\alpha^{(\infty)})_{\nu(\alpha, 2)}(\lambda_{\nu(\alpha, 1)} - \lambda_{\nu(\alpha, 2)})^2 = \lim_{k \rightarrow \infty} M_2(P^{(k)})$.

Thus the set Π_∞ of all $P_\alpha^{(\infty)}$'s are composed of at most finitely many distributions. We shall denote the element of Π_∞ by P_i $i=1, 2, \dots, N$. P_i is characterized by the quantities $(P_i)_{\nu(i, 1)} > 0$ $(P_i)_{\nu(i, 2)} > 0$.

Now we can see that if there is some P_i with $(P_i)_{\nu(i, 1)} \neq (P_i)_{\nu(i, 2)}$ then the distribution P_i^* with $(P_i^*)_{\nu(i, 1)} = (P_i)_{\nu(i, 2)}$ and $(P_i^*)_{\nu(i, 2)} = (P_i)_{\nu(i, 1)}$ is also contained in Π_∞ . This is proved by taking into account of the fact that if $P^{(j)} \rightarrow P_i(j \rightarrow \infty)$ then $TP^{(j)} \rightarrow TP_i(j \rightarrow \infty)$ and the fact that $TP_i = P_i^*$.

It is obvious that if $(P_i)_{\nu(i, 1)} = (P_i)_{\nu(i, 2)} = 1/2$ we have $TP_i = P_i$. Here-

after, we shall sometimes use the notation $P(\nu(i, 1), \nu(i, 2))$ to represent the P_i to clarify the character of the distribution. We shall here observe the domain of T in its natural relative topology generated by the ordinary topology of n -dimensional Euclidean space. Now as is easily seen T is continuous in its domain, and for any neighborhood $U\{P(\nu(i, 1), \nu(i, 2))\}$ of $P(\nu(i, 1), \nu(i, 2))$ there exists a proper neighborhood $V\{P(\nu(i, 2), \nu(i, 1))\}$ of $P(\nu(i, 2), \nu(i, 1))$ with the property that for any P in $V\{P(\nu(i, 2), \nu(i, 1))\}$ its transform TP is contained in $U\{P(\nu(i, 1), \nu(i, 2))\}$. Here we take the neighborhoods $U\{P(\nu(i, 1), \nu(i, 2))\}$ of P_i 's so as to be mutually disjoint. We define the neighborhood $W\{P_i\}$ of P_i as the intersection of $U\{P_i\}$ and $V\{P_i\}$ where V is defined as above. Now there exists a number M such that $P^{(k)}$ with k greater than M are all contained in $\sum_{i=1}^N W\{P_i\}$, otherwise there must be some limit point outside Π_∞ .

We shall represent by $\{P^{(i,j)}, j=1, 2, \dots\}$ the subsequence of $\{P^{(k)}; k=M+1, M+2, \dots\}$ contained in $W\{P_i\}$. Then it is seen that $\{P^{(k)}; k=M+1, M+2, \dots\} = \sum_{i=1}^N \{P^{(i,j)}; j=1, 2, \dots\}$ holds as the relation between the sets of points in the sequences. Now from the definition of M and $P^{(i,j)}$ we can see that $TP^{(i,j)} \in \sum_{i=1}^N W\{P_i\}$ holds for any i and j and further by the definition of W and V we can see that $\{TP^{(i,j)}; j=1, 2, \dots\} \subset U\{TP_i\} = U\{P(\nu(i, 2), \nu(i, 1))\}$. As the $U\{P_i\}$'s are taken to be mutually disjoint, so $TP^{(i,j)} \notin W\{P_i\}$ for P_i different from TP_i . Thus we have $\{TP^{(i,j)}; j=1, 2, \dots\} \subset W\{TP_i\}$. From this last relation we can also see that $\{T^2P^{(i,j)}; j=1, 2, \dots\} \subset W\{P_i\}$, consequently, that $T^{2r+1}P^{(i,j)} \in W\{TP_i\}$ and $T^{2r+2}P^{(i,j)} \in W\{P_i\}$ hold for $r=0, 1, 2, \dots$ and that if $P^{(i,j)}$ is identical to some $P^{(k)}$ of the original sequence $\{P^{(k)}\}$, then $P^{(k+s)} \in W\{P_i\} \cup W\{TP_i\}$ holds for $s=0, 1, 2, \dots$. From these relations it follows that Π_∞ is composed of at most two points $P^{(\infty)}$ and $P^{*(\infty)}$ satisfying the relation $P^{*(\infty)} = TP^{(\infty)}$. When $TP^{(\infty)} = P^{(\infty)}$, Π_∞ is composed of only one point $P^{(\infty)} = P(\nu(\infty, 1), \nu(\infty, 2))$ with $(P^{(\infty)})_{\nu(\infty, 1)} = (P^{(\infty)})_{\nu(\infty, 2)} = 1/2$. Thus we have proved the following;

THEOREM 1: For any $P^{(0)}$ with $(\bar{\lambda} - \lambda(P^{(0)}))^2(P^{(0)}) > 0$ the sequences of distributions $\{T^{2r}P^{(0)}; r=1, 2, \dots\}$ and $\{T^{2r+1}P^{(0)}; r=0, 1, 2, \dots\}$ are convergent to some limiting distributions $P^{(\infty)}$ and $P^{*(\infty)}$, respectively. These distributions are characterized by $(P^{(\infty)})_{\nu(\infty, 1)} > 0$ and $(P^{(\infty)})_{\nu(\infty, 2)} > 0$ satisfying the relation $(P^{(\infty)})_{\nu(\infty, 1)} + (P^{(\infty)})_{\nu(\infty, 2)} = 1$ and $(P^{*(\infty)})_{\nu(\infty, 1)} =$

$(\mathbf{P}^{(\infty)})_{\nu(\infty,2)}, (\mathbf{P}^{*(\infty)})_{\nu(\infty,2)} = (\mathbf{P}^{(\infty)})_{\nu(\infty,1)}$. We have $T\mathbf{P}^{(\infty)} = \mathbf{P}^{*(\infty)}$ and $\mathbf{P}^{*(\infty)}$ is identical to $\mathbf{P}^{(\infty)}$ if and only if $(\mathbf{P}^{(\infty)})_{\nu(\infty,1)} = (\mathbf{P}^{(\infty)})_{\nu(\infty,2)} = 1/2$ holds.

Note; We can see that λ_i 's with $(\mathbf{P}^{(0)})_i = 0$ should have been entirely discarded at the very out-set of our whole discussion, and we shall hereafter assume that $(\mathbf{P}^{(0)})_i > 0$ for $i=1, 2, \dots, n$.

We shall here discuss some of the effects of the initial distribution $\mathbf{P}^{(0)}$ to the limiting distributions $\mathbf{P}^{(\infty)}$ and $\mathbf{P}^{*(\infty)}$.

THEOREM 2: If $\lambda_1 < \lambda_i < \lambda_n (i=2, 3, \dots, n-1)$, then we have

$$\{\nu(\infty, 1), \nu(\infty, 2)\} = \{1, n\}$$

that is, the limiting distributions have their total probability attached to both extremal points.

PROOF; It is obvious that $\lambda_1 < \bar{\lambda}(\mathbf{P}) < \lambda_n$ holds when $(\mathbf{P})_1 > 0$ and $(\mathbf{P})_n > 0$ hold. Thus from the definition of T it is seen that $(\mathbf{P}^{(k)})_1 > 0$ and $(\mathbf{P}^{(k)})_n > 0$ hold when $(\mathbf{P}^{(k-1)})_1 > 0$ and $(\mathbf{P}^{(k-1)})_n > 0$ hold. From this fact, under the condition of the theorem, we have $(\mathbf{P}^{(k)})_1 > 0$ and $(\mathbf{P}^{(k)})_n > 0$ for $k=0, 1, 2, \dots$. Suppose $\nu(\infty, 2) < n$ holds. Then we have $\lambda_{\nu(\infty,1)} < \bar{\lambda}(\mathbf{P}^{(\infty)}), \bar{\lambda}(\mathbf{P}^{*(\infty)}) < \lambda_{\nu(\infty,2)} < \lambda_n$ and we can find a K and $\bar{\lambda}, \underline{\lambda}$ such that, $\lambda_{\nu(\infty,1)} < \underline{\lambda} < \bar{\lambda}(\mathbf{P}^{(k)}) < \bar{\lambda} < \lambda_{\nu(\infty,2)}$ holds for all $k \geq K$. Now taking into account of the relation

$$\begin{aligned} & (\mathbf{P}^{(k+j)})_n / (\mathbf{P}^{(k+j)})_{\nu(\infty,2)} \\ &= \{(\mathbf{P}^{(k)})_n / (\mathbf{P}^{(k)})_{\nu(\infty,2)}\} \prod_{i=1}^j \{(\lambda_n - \bar{\lambda}(\mathbf{P}^{(k+i)}))^2 / (\lambda_{\nu(\infty,2)} - \bar{\lambda}(\mathbf{P}^{(k+i)}))^2\} \\ &\geq \{(\mathbf{P}^{(k)})_n / (\mathbf{P}^{(k)})_{\nu(\infty,2)}\} \prod_{i=1}^j \{(\lambda_n - \underline{\lambda})^2 / (\lambda_{\nu(\infty,2)} - \underline{\lambda})^2\}, \end{aligned}$$

we can see that $(\mathbf{P}^{(k)})_n$ cannot tend to 0, this contradicts to the definition of $\mathbf{P}^{(\infty)}$. The case where it is supposed that $1 < \nu(\infty, 1)$ holds is treated in the same manner.

THEOREM 3: Under the condition of theorem 2 if there is a λ_i ($i \neq 1, n$) for which $\bar{\lambda}(\mathbf{P}^{(k)}) \neq \lambda_i$ for all $k=0, 1, 2, \dots$ then the following inequality holds

$$\left(\frac{\lambda_n - \lambda_1}{2}\right)^2 + \left(\lambda_i - \frac{\lambda_n + \lambda_1}{2}\right)^2 \geq 2\left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2}\right)^2.$$

PROOF: We have

$$\begin{aligned}
 & (\mathbf{P}^{(k+2)})_i / (\mathbf{P}^{(k+2)})_1 \\
 &= \{(\mathbf{P}^{(k)})_i / (\mathbf{P}^{(k)})_1\} \{(\lambda_i - \bar{\lambda}(\mathbf{P}^{(k)}))^2 (\lambda_i - \bar{\lambda}(\mathbf{P}^{(k+1)}))^2\} \\
 & \quad / \{(\lambda_1 - \bar{\lambda}(\mathbf{P}^{(k)}))^2 (\lambda_1 - \bar{\lambda}(\mathbf{P}^{(k+1)}))^2\}.
 \end{aligned}$$

For any $\delta > 0$ there is a $K(\varepsilon)$ such that for $k \geq K(\varepsilon)$ we have

$$\begin{aligned}
 & \{(\lambda_i - \bar{\lambda}(\mathbf{P}^{(k)}))^2 (\lambda_i - \bar{\lambda}(\mathbf{P}^{(k+1)}))^2\} / \{(\lambda_1 - \bar{\lambda}(\mathbf{P}^{(k)}))^2 (\lambda_1 - \bar{\lambda}(\mathbf{P}^{(k+1)}))^2\} \\
 & > \{(\lambda_i - \bar{\lambda}(\mathbf{P}^{(\infty)}))^2 (\lambda_i - \bar{\lambda}(\mathbf{P}^{*(\infty)}))^2\} / \{(\lambda_1 - \bar{\lambda}(\mathbf{P}^{(\infty)}))^2 (\lambda_1 - \bar{\lambda}(\mathbf{P}^{*(\infty)}))^2\} - \varepsilon.
 \end{aligned}$$

From this relation we have

$$(\lambda_i - \bar{\lambda}(\mathbf{P}^{(\infty)}))^2 (\lambda_i - \bar{\lambda}(\mathbf{P}^{*(\infty)}))^2 / (\lambda_1 - \bar{\lambda}(\mathbf{P}^{(\infty)}))^2 (\lambda_1 - \bar{\lambda}(\mathbf{P}^{*(\infty)}))^2 \leq 1$$

and

$$\begin{aligned}
 & \left(\lambda_i - \frac{\lambda_n + \lambda_1}{2} - \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2 \left(\lambda_i - \frac{\lambda_n + \lambda_1}{2} - \left(\bar{\lambda}(\mathbf{P}^{*(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2 \\
 & \leq \left(\lambda_1 - \frac{\lambda_n + \lambda_1}{2} - \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2 \left(\lambda_1 - \frac{\lambda_n + \lambda_1}{2} - \left(\bar{\lambda}(\mathbf{P}^{*(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2.
 \end{aligned}$$

By using the relation $\bar{\lambda}(\mathbf{P}^{*(\infty)}) = \lambda_n + \lambda_1 - \bar{\lambda}(\mathbf{P}^{(\infty)})$ we have

$$\begin{aligned}
 & \left(\lambda_i - \frac{\lambda_n + \lambda_1}{2} - \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2 \left(\lambda_i - \frac{\lambda_n + \lambda_1}{2} + \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2 \\
 & \leq \left(\frac{\lambda_1 - \lambda_n}{2} - \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2 \left(\frac{\lambda_1 - \lambda_n}{2} + \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right) \right)^2
 \end{aligned}$$

or

$$\begin{aligned}
 & \left(\left(\lambda_i - \frac{\lambda_n + \lambda_1}{2} \right)^2 - \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right)^2 \right)^2 \\
 & \leq \left(\left(\frac{\lambda_1 - \lambda_n}{2} \right)^2 - \left(\bar{\lambda}(\mathbf{P}^{(\infty)}) - \frac{\lambda_n + \lambda_1}{2} \right)^2 \right)^2.
 \end{aligned}$$

From this last inequality we can get the desired result.

§ 2. Application to the optimum gradient method

In this section we shall discuss the application of the results of the former section to the so-called optimum gradient method. Let A be an n -by- n matrix, not singular, and let x , b denote the n -rowed column vectors. The optimum gradient method for the solution of simultaneous linear equation $Ax = b$ with respect to the metric $\| \cdot \|_P$ is defined as fol-

lows*); Take a positive definite symmetric matrix P . Then the gradient of the error function $\|Az-b\|_P^2 = (Az-b, P(Az-b))$ at z is given by $2(A'P(Az-b))$. Given the k -th approximate solution $x_k = x + \varepsilon_k$ where x represents the desired solution the optimum gradient method proceeds by finding the $(k+1)$ -th approximate solution $x_{k+1} = x_k - \gamma_k \zeta_k$ where

$$\begin{aligned} \text{Min}_{\gamma} \quad & \|A(x_k - \gamma \zeta_k) - b\|_P^2 = \|A(x_k - \gamma_k \zeta_k) - b\|_P^2 \\ \gamma_k = & (A\varepsilon_k, A\zeta_k)_P / \|A\zeta_k\|_P^2 \\ \zeta_k = & \frac{1}{2} \text{ gradient at } x_k = A'PA\varepsilon_k. \end{aligned}$$

We shall represent by $(0 <) \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ the eigenvalues of $A'PA$ and by $\xi_1, \xi_2, \dots, \xi_n$ the corresponding eigenvectors which are supposed to be orthonormal. We exclude the trivial case where $\lambda_1 = \lambda_2 = \dots = \lambda_n$. Then we have the following.

THEOREM 4; *In the optimum gradient method with respect to the metric $\|\cdot\|_P$*

- i) $\varepsilon_k (=x_k - x)$ tends to be approximated by a linear combination of two fixed eigenvectors of $A'PA$ with the eigenvalues equal to $\text{Max}(\lambda_i; (\varepsilon_0, \xi_i) \neq 0)$ and $\text{Min}(\lambda_i; (\varepsilon_0, \xi_i) \neq 0)$, respectively, and
- ii) ε_k alternates asymptotically in two fixed directions.

PROOF: Here we prove the theorem for the case where $\lambda_1 < \lambda_2 < \dots < \lambda_n$ and $(\varepsilon_0, \xi_1) \neq 0, (\varepsilon_0, \xi_n) \neq 0$ hold. Modifications necessary for the proofs of other cases are obvious and are omitted here.

Suppose that the error ε_k of the k -th approximate solution is represented as $\varepsilon_k = \sum_{i=1}^n \alpha_i^{(k)} \lambda_i^{-1} \xi_i$. Then ε_{k+1} is given by the following;

$$\begin{aligned} \varepsilon_k &= \sum_{i=1}^n \alpha_i^{(k+1)} \lambda_i^{-1} \xi_i = \varepsilon_k - \gamma_k \zeta_k = \varepsilon_k - \gamma_k A'PA\varepsilon_k \\ &= \sum_{i=1}^n \alpha_i^{(k)} (1 - \gamma_k \lambda_i) \lambda_i^{-1} \xi_i = \gamma_k \sum_{i=1}^n \alpha_i^{(k)} (\gamma_k^{-1} - \lambda_i) \lambda_i^{-1} \xi_i \end{aligned}$$

where $\gamma_k^{-1} = \|A\zeta_k\|_P^2 / (A\varepsilon_k, A\zeta_k)_P = \sum_{i=1}^n (\alpha_i^{(k)})^2 \lambda_i / \sum_{i=1}^n (\alpha_i^{(k)})^2 > 0$.

If we consider the set of values $((\alpha_i^{(k)})^2 / \sum_{i=1}^n (\alpha_i^{(k)})^2, (\alpha_2^{(k)})^2 / \sum_{i=1}^n (\alpha_i^{(k)})^2, \dots, (\alpha_n^{(k)})^2 / \sum_{i=1}^n (\alpha_i^{(k)})^2)$ as a probability distribution $P^{(k)}$ over $(\lambda_1, \lambda_2, \dots, \lambda_n)$, then

*) Here we use the notations $\|x\|_P^2$ and $(x, y)_P$ to represent the quantities (x, Px) and (x, Py) , respectively, where (x, y) denotes the ordinary inner product of vectors x and y . A' denotes the transposition of A .

$\mathbf{P}^{(k+1)} = ((\alpha_1^{(k+1)})^2 / \sum_{i=1}^n (\alpha_i^{(k+1)})^2, (\alpha_2^{(k+1)})^2 / \sum_{i=1}^n (\alpha_i^{(k+1)})^2, \dots, (\alpha_n^{(k+1)})^2 / \sum_{i=1}^n (\alpha_i^{(k+1)})^2)$ is represented as $\mathbf{P}^{(k+1)} = T\mathbf{P}^{(k)}$, transformation T being defined in § 1. Thus the results of the former section are applicable to the present sequence of $\mathbf{P}^{(k)}$'s and we have for some non-zero c

$$\lim_{h \rightarrow \infty} \mathbf{P}^{(2h)} = \left(\frac{1}{1+c^2}, 0, \dots, 0, \frac{c^2}{1+c^2} \right) \equiv \mathbf{P}^{(\infty)}$$

$$\lim_{h \rightarrow \infty} \mathbf{P}^{(2h+1)} = \left(\frac{c^2}{1+c^2}, 0, \dots, 0, \frac{1}{1+c^2} \right) \equiv \mathbf{P}^{*(\infty)}$$

and the limit points of the sequence of direction-cosines

$$\left(\frac{\alpha_1^{(2h)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(2h)})^2}}, \frac{\alpha_2^{(2h)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(2h)})^2}}, \dots, \frac{\alpha_n^{(2h)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(2h)})^2}} \right)$$

are limited to the set of points of the type

$$\left(\pm \frac{1}{\sqrt{1+c^2}}, 0, \dots, 0, \pm \frac{c}{\sqrt{1+c^2}} \right).$$

Now consider the transformation \hat{T} defined over the set of points $\{(\delta_1, \delta_2, \dots, \delta_n); \sum_{i=1}^n \delta_i^2 = 1, \delta_i^2 < 1 (i=1, 2, \dots, n)\}$ by the following;

$$\hat{T}(\delta_1, \delta_2, \dots, \delta_n)$$

$$= \left(\frac{\delta_1(\bar{\lambda} - \lambda_1)}{\sqrt{\sum_{i=1}^n \delta_i^2(\bar{\lambda} - \lambda_i)^2}}, \frac{\delta_2(\bar{\lambda} - \lambda_2)}{\sqrt{\sum_{i=1}^n \delta_i^2(\bar{\lambda} - \lambda_i)^2}}, \dots, \frac{\delta_n(\bar{\lambda} - \lambda_n)}{\sqrt{\sum_{i=1}^n \delta_i^2(\bar{\lambda} - \lambda_i)^2}} \right)$$

where $\bar{\lambda} = \sum_{i=1}^n \delta_i^2 \lambda_i$. Then we have

$$\left(\frac{\alpha_1^{(k+1)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(k+1)})^2}}, \frac{\alpha_2^{(k+1)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(k+1)})^2}}, \dots, \frac{\alpha_n^{(k+1)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(k+1)})^2}} \right)$$

$$= \hat{T} \left(\frac{\alpha_1^{(k)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(k)})^2}}, \frac{\alpha_2^{(k)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(k)})^2}}, \dots, \frac{\alpha_n^{(k)}}{\sqrt{\sum_{i=1}^n (\alpha_i^{(k)})^2}} \right)$$

and relations such as

$$\hat{T} \left(\frac{1}{\sqrt{1+c^2}}, 0, \dots, 0, \frac{c}{\sqrt{1+c^2}} \right) = \left(\frac{|c|}{\sqrt{1+c^2}}, 0, \dots, 0, \frac{-c}{|c|\sqrt{1+c^2}} \right)$$

$$\hat{T} \left(\frac{|c|}{\sqrt{1+c^2}}, 0, \dots, 0, \frac{-c}{|c|\sqrt{1+c^2}} \right) = \left(\frac{1}{\sqrt{1+c^2}}, 0, \dots, 0, \frac{c}{\sqrt{1+c^2}} \right).$$

Obviously \hat{T} is continuous in its domain and by applying the entirely similar arguments in the proof of theorem 1 of § 1, which were used to show the oscillatory convergence of $P^{(k)}$, we can get the conclusions i) and ii) of the present theorem.

Note: It may be of some use to remark here the relations $(\zeta_k, \zeta_{k+1})=0$,

$$\frac{\|\zeta_{k+1}\|_I^2}{\|\zeta_k\|_I^2} = \frac{(\lambda - \bar{\lambda}(P^{(k)}))^2(P^{(k)})}{\{\bar{\lambda}(P^{(k)})\}^2} = (\text{coefficient of variation of } P^{(k)})^2.$$

Thus taking into account the fact that $\gamma_k^{-1} = \bar{\lambda}(P^{(k)})$ and lemma 2 of § 1, we can see that the ratio $\|\zeta_{k+1}\|_I^2 / \|\gamma_k \zeta_k\|_I^2 = (\lambda - \bar{\lambda}(P^{(k)}))^2(P^{(k)})$ tends monotone-increasingly to the limit $(\lambda_n - \lambda_1)^2 c^2 (1 + c^2)^{-2}$.*)

This theorem provides the theoretical foundations of the following acceleration procedure which is a slightly modified version of the acceleration procedure proposed by Forsythe and Motzkin [5].

Acceleration procedure for the optimum gradient method:

In the optimum gradient method when the direction of vector ζ_k is nearly the same as that of ζ_{k-2} , it is recommendable to insert the step defined by the following;

$$x_{k+1} = x_k - \hat{\gamma}_k(x_{k-2} - x_k)$$

where

$$\hat{\gamma}_k = (A\varepsilon_k, A(x_{k-2} - x_k))_P / \|A(x_{k-2} - x_k)\|_P^2.$$

The rationale for this procedure is as follows; The fact that the direction of ζ_k is nearly the same of that of ζ_{k-2} means that $\zeta_{k-2} \doteq \|\zeta_{k-2}\|_I \|\zeta_k\|_I^{-1} \zeta_k$ where the sign \doteq is used in place of the description "is approximately equal to". Then we can see that

$$\begin{aligned} x_{k-2} - x_k &= \varepsilon_{k-2} - \varepsilon_k = (A'PA)^{-1}(\zeta_{k-2} - \zeta_k) \\ &\doteq (\|\zeta_{k-2}\|_I \|\zeta_k\|_I^{-1} - 1)(A'PA)^{-1}\zeta_k \\ &= (\|\zeta_{k-2}\|_I \|\zeta_k\|_I^{-1} - 1)\varepsilon_k \end{aligned}$$

holds and it is obvious that in this case $x_{k-2} - x_k$ will be a good candidate for the correcting term of the $x_k = x + \varepsilon_k$. Now if we can suppose that $\varepsilon_{k-2} \doteq \beta(\lambda_1^{-1}\xi_1 + c\lambda_n^{-1}\xi_n)$ then we have

*) It is supposed here that $(P^{(0)})_1 > 0$ and $(P^{(0)})_n > 0$ hold.

$$\begin{aligned}
 \zeta_{k-2} &\doteq \beta(\xi_1 + c\xi_n) \\
 \gamma_{k-2} &= (A\varepsilon_{k-2}, A\zeta_{k-2})_P / \|A\zeta_{k-2}\|_P^2 \doteq (1+c^2)/(\lambda_1 + c^2\lambda_n) \\
 \varepsilon_{k-1} &= \varepsilon_{k-2} - \gamma_{k-2}\zeta_{k-2} \doteq \beta(\lambda_1 + c^2\lambda_n)^{-1}(\lambda_n - \lambda_1)c^2\{\lambda_1^{-1}\xi_1 - c^{-1}\lambda_n^{-1}\xi_n\} \\
 \|A\varepsilon_{k-2}\|_P^2 &= (\varepsilon_{k-2}, \zeta_{k-2}) \doteq \beta^2(\lambda_1^{-1} + c^2\lambda_n^{-1}) \\
 \|A\varepsilon_{k-1}\|_P^2 &= (\varepsilon_{k-1}, \zeta_{k-1}) \doteq \beta^2(\lambda_1 + c^2\lambda_n)^{-2}(\lambda_n - \lambda_1)^2 c^4(\lambda_1^{-1} + c^{-2}\lambda_n^{-1}) \\
 \varepsilon_k &\doteq \beta(\lambda_n - \lambda_1)^2(\lambda_1 + c^2\lambda_n)^{-1}(\lambda_1 + c^{-2}\lambda_n)^{-1}\{\lambda_1^{-1}\xi_1 + c\lambda_n^{-1}\xi_n\} \\
 &\doteq \|A\varepsilon_{k-1}\|_P^2 / \|A\varepsilon_{k-2}\|_P^2 \varepsilon_{k-2}.
 \end{aligned}$$

It is obvious that in this case the direction of ζ_k is nearly the same as that of ζ_{k-2} . This fact and the result of theorem 1 which assures that we can expect that ε_{k-2} takes the form just stated when k becomes large, give the theoretical basis of the recommendation of the present acceleration procedure.

The rate of convergence of the optimum gradient method :

By using the results of calculations in the former paragraph we can see that when the distribution $P^{(k)}$ corresponding to ε_k tends alternately to the limiting distributions $P^{(\infty)} = (1/(1+c^2), 0, \dots, 0, c^2/(1+c^2))$ and $P^{*(\infty)} = (c^2/(1+c^2), 0, \dots, 0, 1/(1+c^2))$ i.e. ε_k tends to be approximated by $\beta_k\{\lambda_1^{-1}\xi_1 + c\lambda_n^{-1}\xi_n\}$ with some scalar β_k and c , the rate of convergence $\|A\varepsilon_{k+1}\|_P^2 / \|A\varepsilon_k\|_P^2$ tends to the value $(\lambda_n - \lambda_1)^2\{(\lambda_1 + \lambda_n)^2 + (c - c^{-1})^2\lambda_1\lambda_n\}^{-1}$. Thus the rate of convergence of the optimum gradient method is eventually determined by the value of c^2 which is inherited from the initial vector x_0 or $P^{(0)}$. The values $(\lambda_n - \lambda_1)^2\{(\lambda_1 + \lambda_n)^2 + (c - c^{-1})^2\lambda_1\lambda_n\}^{-1}$ attains its maximum value $(\lambda_n - \lambda_1)^2(\lambda_1 + \lambda_n)^{-2}$ when $c^2 = 1$ holds. Now $(\lambda_n - \lambda_1)^2 / (\lambda_n + \lambda_1)^2 = (t-1)^2 / (t+1)^2$ where $t \equiv \lambda_n / \lambda_1$ is a so-called condition-number of the matrix of $A'PA$.

We shall here make a slight digression for investigation of the meaning of the condition-number. Now define the function $f(\hat{x}, c) \equiv \|A(\hat{x} - c\hat{\zeta}) - b\|_P^2 / \|A\hat{x} - b\|_P^2$ of a vector $\hat{x} (\neq x)$ and a scalar c where $\hat{\zeta} = A'PA\hat{\varepsilon}$, $\hat{\varepsilon} = \hat{x} - x$ (x is the solution of $Ax = b$). If $\hat{\varepsilon} = \sum_{i=1}^n \beta_i \xi_i$, we have $f(\hat{x}, c) = \sum_{i=1}^n \beta_i^2 \lambda_i (1 - c\lambda_i)^2 / \sum_{i=1}^n \beta_i^2 \lambda_i$ and thus $\text{Min}_c \text{Max}_{\hat{x}} f(\hat{x}, c) = \text{Min}_c (\text{Max}_{\hat{x}} ((1 - c\lambda_1)^2, (1 - c\lambda_n)^2)) = (1 - 2\lambda_1 / (\lambda_1 + \lambda_n))^2 = (\lambda_n - \lambda_1)^2 / (\lambda_n + \lambda_1)^2$ for $c = 2 / (\lambda_1 + \lambda_n)$. Obviously $\text{Min}_c \text{Max}_{\hat{x}} f(\hat{x}, c) \geq \text{Max}_{\hat{x}} \text{Min}_c f(\hat{x}, c)$ holds and $\text{Min}_c f(\hat{x}, c) = f(\hat{x}, \hat{\gamma})$ where $\hat{\gamma} = (A\hat{\varepsilon}, A\hat{\zeta})_P / \|A\hat{\zeta}\|_P^2$ is given by the optimum gradient method. Thus we get $f(\hat{x}, \hat{\gamma}) \leq (\lambda_n - \lambda_1)^2 / (\lambda_n + \lambda_1)^2$. We have already seen that $f(\hat{x}, \hat{\gamma}) = (\lambda_n - \lambda_1)^2 / (\lambda_n + \lambda_1)^2$ holds for \hat{x} with corresponding probability

distribution $\hat{P} = (1/2, 0, \dots, 0, 1/2)$ and we can see that the value $(\lambda_n - \lambda_1)^2/(\lambda_n + \lambda_1)^2 = (t - 1)^2/(t + 1)^2$ gives the least upper bound of the convergence rate $f(\hat{x}, \hat{y})$ of the optimum gradient method. It is stated in [5] that a proof of this fact was given by Kantrovitch [8] but as the Kantrovitch's paper was not available, we gave a proof to it for the sake of completeness of the following discussion.

Here we use the result of theorem 3 of § 1 to see why the optimum gradient method often converges with the convergence rate nearly equal to its worst possible value $(\lambda_n - \lambda_1)^2/(\lambda_n + \lambda_1)^2$. This fact was also noticed by Forsythes [5]. For $P^{(\infty)} = (1/(1+c^2), 0, \dots, 0, c^2/(1+c^2))$ and $P^{*(\infty)} = (c^2/(1+c^2), 0, \dots, 0, 1/(1+c^2))$ we have $\bar{\lambda}(P^{(\infty)}) = (\lambda_1 + c^2\lambda_n)/(1+c^2)$ and $\bar{\lambda}(P^{*(\infty)}) = (c^2\lambda_1 + \lambda_n)/(1+c^2)$ and thus we have $(\bar{\lambda}(P^{(\infty)}) - (\lambda_n + \lambda_1)/2)^2 = (\bar{\lambda}(P^{*(\infty)}) - (\lambda_n + \lambda_1)/2)^2 = (\lambda_n - \lambda_1)^2(1-c^2)^2/4(1+c^2)^2$. Thus under the condition which assumes that the point λ_i is not discarded during the course of approximation procedure, we have from the result of theorem 3 of § 1

$$\left(\frac{\lambda_n - \lambda_1}{2}\right)^2 + \left(\lambda_i - \frac{\lambda_n + \lambda_1}{2}\right)^2 \geq \frac{(1-c^2)^2}{2(1+c^2)^2} (\lambda_n - \lambda_1)^2.$$

By putting $\delta_i \equiv (\lambda_i - (\lambda_n + \lambda_1)/2)/((\lambda_n - \lambda_1)/2)$ we get from the above inequality the following

$$4\left\{\frac{1+\delta_i^2}{1-\delta_i^2}\right\} \geq (c-c^{-1})^2.$$

Thus, for example, when there exists some λ_i which satisfies the condition of theorem 3 and with $|\delta_i| \ll 1$ we can expect that $(c-c^{-1})^2$ is near or less than 4. Now the convergence rate at the point corresponding to the $P^{(\infty)} = (1/(1+c^2), 0, \dots, 0, c^2/(1+c^2))$ is given by $(\lambda_n - \lambda_1)^2(\lambda_n + \lambda_1)^{-2} \{1 + (c-c^{-1})^2(\lambda_n\lambda_1^{-1} + \lambda_1\lambda_n^{-1} + 2)^{-1}\}^{-1}$ and by putting $\varepsilon = 1 - (\lambda_n - \lambda_1)/(\lambda_n + \lambda_1)$ this is represented as $(\lambda_n - \lambda_1)^2(\lambda_n + \lambda_1)^{-2} \{1 + (c-c^{-1})^2(2\varepsilon^{-1} + 1 + \varepsilon(2-\varepsilon)^{-1})^{-1}\}^{-1}$. Thus we can see that the convergence rate tends to be greater than $(\lambda_n - \lambda_1)^2(\lambda_n + \lambda_1)^{-2} \{1 + (\varepsilon/2)(c-c^{-1})^2\}^{-1}$ or a fortiori than $(\lambda_n - \lambda_1)^2(\lambda_n + \lambda_1)^{-2} \{1 + 2\varepsilon(1+\delta_i^2)(1-\delta_i^2)^{-1}\}^{-1}$. From this last relation we can see that when we use the optimum gradient method for ill-conditioned ($\varepsilon \ll 1$) matrix $A'PA$, then the rate of convergence tends near to its worst possible value, especially when there is some λ_i near the value $(\lambda_n + \lambda_1)/2$ i.e. $|\delta_i| \ll 1$.

As to one of the numerical examples described by Forsythes [5], and which will be treated fully in the following paragraph of this section, we have $\epsilon=0.01073$ and

$$\begin{aligned} & (\lambda_n - \lambda_1)^2 (\lambda_n + \lambda_1)^{-2} \{1 + 2\epsilon(1 + \delta_1^2)(1 - \delta_1^2)^{-1}\}^{-1} \\ &= (\lambda_n - \lambda_1)^2 (\lambda_n + \lambda_1)^{-2} \times 0.9789 = 0.9580 \quad \text{for } \delta_1^2 = 0.001321 \\ &= (\lambda_n - \lambda_1)^2 (\lambda_n + \lambda_1)^{-2} \times 0.7171 = 0.7018 \quad \text{for } \delta_2^2 = 0.896854 . \end{aligned}$$

We can clearly see in this example the effect of the small value of ϵ making the convergence rate tend near to its worst possible value. Thus, it seems that our present analysis gives fairly general theoretical explanation to the fact that the optimum gradient method converges slowly when the matrix is ill-conditioned.

Numerical example

Here we shall present some numerical results obtained by using the same A and b as those treated by Forsythes [5] and by putting $P=I$ (identity matrix). These results are obtained by using a FACOM-128 relay computer of our institute. We have used the computation scheme of optimum gradient method where the acceleration step is inserted automatically when the condition $(\zeta_{k-2}, \zeta_k)_I / (\|\zeta_{k-2}\|_I \|\zeta_k\|_I) > \delta$ is satisfied for some preassigned value of $\delta (1 > \delta > 0)$.

We have

$$A = \begin{bmatrix} \sqrt{0.00268704} & & & & & \\ & \sqrt{0.01581310} & & & & 0 \\ & & \sqrt{0.08234830} & & & \\ & & & \sqrt{0.17590130} & & \\ 0 & & & & \sqrt{0.25946632} & \\ & & & & & \sqrt{0.49823436} \end{bmatrix}$$

$$b = [0, \quad 0, \quad 0, \quad 0, \quad 0, \quad 0] \quad \text{and}$$

$$A'PA = \begin{bmatrix} 0.00268704 (= \lambda_1) & & & & & \\ & 0.01581310 (= \lambda_2) & & & & 0 \\ & & 0.08234830 (= \lambda_3) & & & \\ & & & 0.17590130 (= \lambda_4) & & \\ 0 & & & & 0.25946632 (= \lambda_5) & \\ & & & & & 0.49823436 (= \lambda_6) \end{bmatrix} .$$

We have used three values 0.99, 0.999 and 0.9999 as δ . In fig's. 1, 2 and

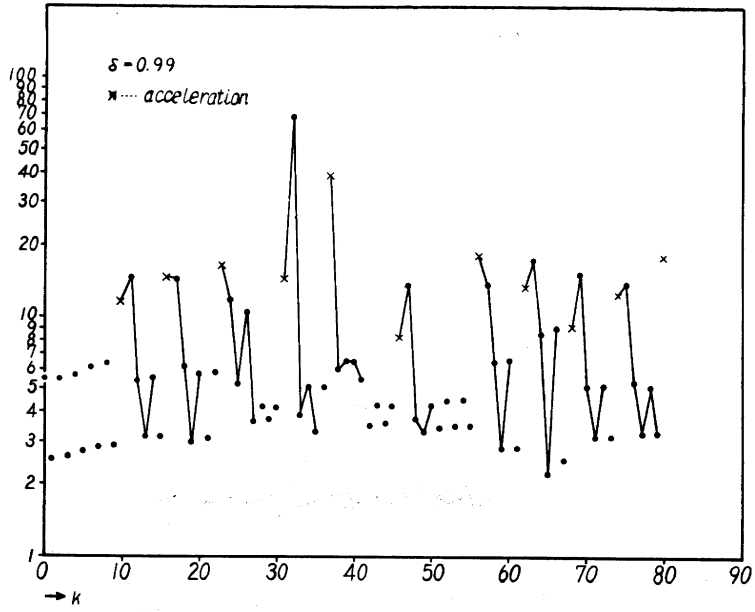


Fig. 1

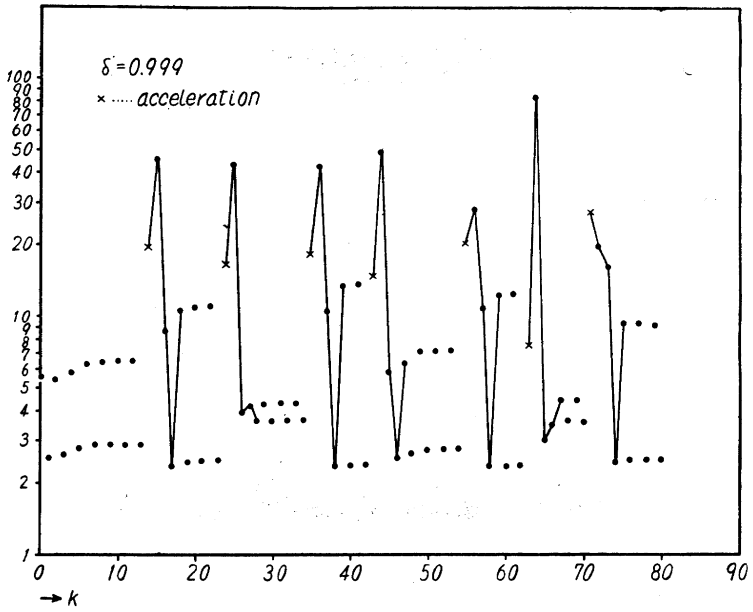


Fig. 2

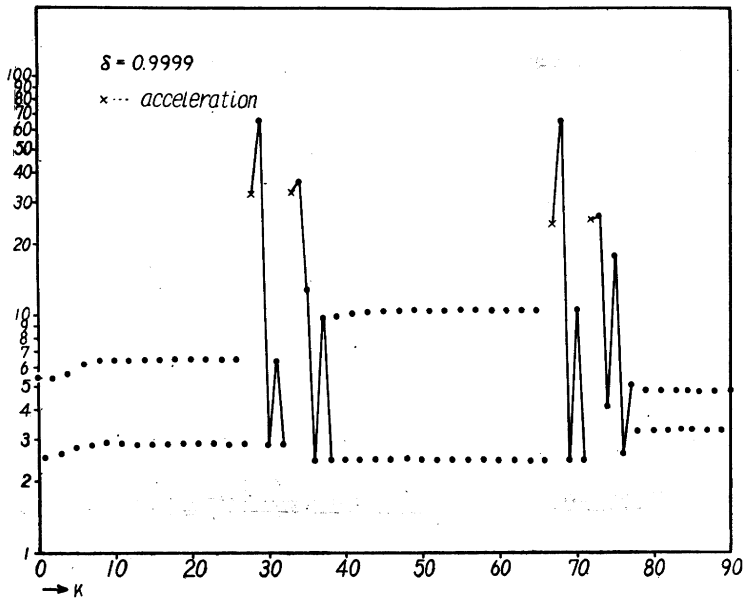


Fig. 3

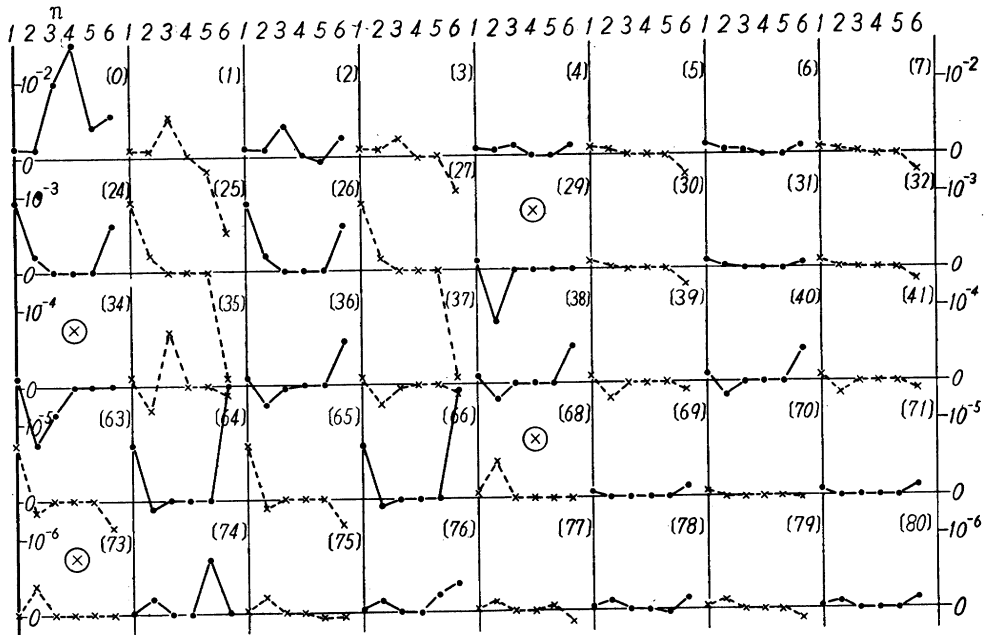


Fig. 4

3 the sequence of γ_k 's are represented by \bullet and $\hat{\gamma}_k$'s (acceleration step) are represented by \times . These three figures correspond to the above mentioned three values of δ . In fig. 4 part of the values of ζ_k 's are illustrated for the case $\delta=0.9999$. The numbers in the bracket in fig. 4 correspond to the k 's of fig. 3 and the symbols \otimes signify the result of acceleration procedure.

We can see in these examples a good agreement between the results of our theoretical analysis made in this paper and the results of practical computations.

Acknowledgement

I wish to express my thanks to Mr. K. Isii for reading the manuscript of §1 of this paper. Thanks are also due to Miss Y. Saigusa who performed the programmings and operations of FACOM-128 automatic relay computer to prepare the numerical results illustrated in §2.

THE INSTITUTE OF STATISTICAL MATHEMATICS

REFERENCES

- [1] Akaike, H. "On a computation method for eigenvalue problems and its application to statistical analysis," *Annals of the Institute of Statistical Mathematics* 10, 1-20 (1958).
- [2] Chernoff, H. and Divinsky, N. "The computation of maximum-likelihood estimates of linear structural equations," Chapter 10 in *Studies in Econometric Method*, Cowles Commission Monograph 14, Wm. C. Hood and T. C. Koopmans editors, John Wiley and Sons, Inc. New York, (1953).
- [3] Crockett, J. B. and Chernoff, H. "Gradient methods of maximization," *Pacific J. Math.*, 5, 33-50 (1955).
- [4] Forsythe, G. E. "Solving linear algebraic equations can be interesting," *Bull. Amer. Math. Soc.*, 59, 299-329 (1953).
- [5] Forsythe, A. I. and Forsythe, G. E. "Punched-card experiments with accelerated gradient methods for linear equations," *Contributions to the Solution of systems of Linear Equations and the Determination of Eigenvalues*, N.B.S. Applied Mathematics Series 39, Olga Taussky editor, 55-69 (1954).
- [6] Forsythe, G. E. and Motzkin, T. S. "Asymptotic properties of the optimum gradient method," *Bul. Am. Math. Soc.*, 57, 183 (1951) (Abstract).
- [7] Forsythe, G. E. and Motzkin, T. S. "Acceleration of the optimum gradient method," Preliminary report, *Bul. Am. Math. Soc.*, 57, 304-305 (1951) (Abstract).
- [8] Kantrovich, L. V. "Functional analysis and applied mathematics," *Uspekhi Matematicheskikh Nauk* (NS), 3, No. 6, 89-185 (1951).