

1

FROM "REACTION CONCEPT" TO "CONJUGATE GRADIENT": HAVE WE MADE ANY PROGRESS?

T. K. Sarkar

- 1.1 Introduction
 - 1.2 Application of the Reaction Concept to Scattering Problems
 - 1.3 The Method of Moments
 - 1.4 The Conjugate Gradient Method for Solution of Operation Equations
 - 1.4a. The difference between application of the conjugate gradient method (CG) to solve a matrix equation as opposed to direct application of CG to solution of an operator equation
 - 1.5 Development of a Numerical Method
 - 1.6 Numerical Considerations
 - 1.7 A Note on Computational Efficiency
 - 1.8 Conclusion
- References

1.1 Introduction

In recent times there seems to be an attempt to interpret every numerical method as a generalization of the moment methods [1]. Even though such an attempt is extremely useful, the fact of the matter is that there are some fundamental philosophical differences among various techniques which are difficult to reconcile. In this chapter an attempt has been made to describe each of the methods historically, by tracing them back to their roots. By doing this, several philosophical concepts become quite clear, and the rationale for utilizing a particular technique becomes apparent. We first start with the reaction concept,

and follow how this technique has been utilized effectively in solving scattering problems. In the third section we present the method of moments. Section 1.4 describes the conjugate gradient method, and the various philosophical differences are outlined. The rate of convergence of each technique is also outlined. The unity in diversity amongst the various techniques is presented in Section 1.6 from a purely numerical standpoint.

1.2 Application of the Reaction Concept to Scattering Problems

Historically, Rumsey [2] utilized the "reaction concept", developed by Schelkunoff and Friis [3], to solve electromagnetic scattering problems, utilizing surface currents. This work was later extended by Cohen [4] for the analysis of electromagnetic scattering utilizing volume currents.

The application of the reaction concept to electromagnetics is natural, as the "reaction" is a basic physical observable [2]. Also, according to Cohen, "Any microwave measurement consists of measuring, not the electric field at a point, but, rather, the signal at the terminations of the antenna. Thus two antennas, for transmission and reception, are inherent in a measurement." [4] The reaction, between two sources a and b consisting of electric currents $\mathbf{J}(a)$ and $\mathbf{J}(b)$ and magnetic currents $\mathbf{M}(a)$ and $\mathbf{M}(b)$, is defined by

$$\langle a ; b \rangle_s = \iiint_{V_i} [\mathbf{E}(a) \bullet \mathbf{J}(b) - \mathbf{H}(a) \bullet \mathbf{M}(b)] dV \quad (1)$$

where the electric and the magnetic fields radiated by the sources a and b are $\mathbf{E}(a)$, $\mathbf{H}(a)$ and $\mathbf{E}(b)$, $\mathbf{H}(b)$, respectively. The integration volume, V_b , contains $\mathbf{J}(b)$ and $\mathbf{M}(b)$. The use of $\langle \bullet ; \bullet \rangle$ for the reaction is rather unfortunate, as the same symbol is used for the definition of the inner product. We will use the subscript 'S' to represent a symmetric product, namely reaction. The subscript 'H' will be used to represent the actual inner product, namely, the Hilbert inner product. Therefore,

$$\langle V ; I \rangle_H = \int V(z) \bar{I}(z) dz \quad (2)$$

$$\langle V; I \rangle_S = \int V(z) I(z) dz \quad (3)$$

where the overbar represents the complex conjugate. From a purely circuits point of view, the symmetric product arises from reciprocity, whereas the Hilbert inner product describes power [5,6]. Therefore, in this chapter, reciprocity will be treated as a totally different entity than power, as they are physically different properties. The significance of the use of the different inner products will be made clear when one talks about convergence.

Equation (2) is the definition for power, and conservation of power always holds, irrespective of what type of circuit one is dealing with. Whereas reciprocity holds only under certain circumstances and, hence, is not very general, in a philosophical development. Reciprocity is usually defined utilizing the definition described by [3].

It is important to reemphasize that the symmetric product in (1) does not contain the conjugate. Hence, it is reaction, and not power. We will often dwell on this difference.

Under suitable conditions, the reciprocity theorem holds [4], and we have

$$\langle a; b \rangle_S = \langle b; a \rangle_S \quad (4)$$

Analysis of electromagnetic scattering from perfectly conducting structures is now developed utilizing this reaction technique, described by (4) [7]. For the problem of scattering from a perfectly conducting body, the impressed electric and magnetic currents ($\mathbf{J}_i; \mathbf{M}_i$) generate the electric and magnetic field intensities ($\mathbf{E}^i, \mathbf{H}^i$), in the absence of the conductor. We assume that the exterior medium is free space. From the surface equivalence theorem of Schelkunoff [3], the interior field will vanish (without disturbing the exterior field) if we introduce surface electric currents

$$\mathbf{J}_s = \mathbf{n} \times \mathbf{H}^+ \quad (5)$$

on the closed surface S_c of the scatterer. The unit vector, \mathbf{n} , is directed outward and \mathbf{H}^+ is the total magnetic field just outside S_c , the conductor surface. In this situation, the scatterer is replaced by free space without disturbing the field anywhere, since the internal fields are zero. We now place an electric test source \mathbf{J}_t in this internal region, and find from the reciprocity theorem that

$$\iint_{S_c} \mathbf{J}_s \cdot \mathbf{E}_t \, ds + \iint_{S_i} \mathbf{J}_i \cdot \mathbf{E}_t \, ds = \iint_T \mathbf{J}_t \cdot \mathbf{E} \, ds \quad (6)$$

where \mathbf{E}_t is the free space field of the test source, and \mathbf{E} is the total field produced by \mathbf{J}_s and \mathbf{J}_i inside S_c , where the test source is placed. T is the integration region of the test source. Since in the interior, the test source has zero reaction with the other sources, it follows from using reciprocity on \mathbf{J}_t and \mathbf{J}_s , and from (6), that

$$\iint_T \mathbf{E}_s \bullet \mathbf{J}_t \, ds = \iint_{S_c} \mathbf{J}_s \bullet \mathbf{E}_t \, ds = - \iint_{S_i} \mathbf{J}_i \bullet \mathbf{E}_t \, ds = - \iint_T \mathbf{E}_t \bullet \mathbf{J}_t \, ds \quad (7)$$

The above equation is obtained from direct application of reciprocity, as the right hand side of (6) is zero. This is because there cannot be any fields internal to the conductor. One can now see that the first and the fourth integrals in (7) constitute an integral equation for the scattering problem. The objective is to use this equation to determine the unknown \mathbf{J}_s in a finite series of known expansion functions \mathbf{F}_n , with unknown expansion coefficients a_n , to be determined. Therefore,

$$\mathbf{J}_s \cong \mathbf{J}_s^N = \sum_{n=1}^N a_n \mathbf{F}_n \quad (8)$$

It is often "assumed" that the test sources, \mathbf{J}_m , are of the same size, shape and functional form as the expansion functions, \mathbf{F}_n , for computational simplicity. Then from (7) and (8) we get the following matrix equation [7,8]

$$\sum_{n=1}^N a_n C_{mn} = B_m, m = 1, 2, \dots, N \quad (9)$$

where

$$C_{mn} = \iint_T \mathbf{E}_n \bullet \mathbf{J}_m \, ds \quad (10)$$

$$B_m = \iint_T \mathbf{E}_m \bullet \mathbf{J}_m \, ds \quad (11)$$

and \mathbf{E}_n is the field due to \mathbf{F}_n .

Even though the reaction concept leads one to arrive at a stationary form more naturally and pictorially than the conventional variational methods, the two techniques yield identical analysis equations. However, "the variational methods start with a functional and the

motivation of the manipulation to arrive at a stationary form is often obscure." In addition, there are several important questions that need to be addressed:

- A) What are the best location, size and shape for the test sources?
- B) How many test sources should one choose to get an accuracy of (say) 1% in the solution?
- C) Does this procedure guarantee monotonic convergence of the sequence J^N as the number of expansion functions and test sources is increased?

From a theoretical point of view, all these questions remain unanswered for this method. What the reaction theorem guarantees is only a distributional convergence or weak convergence. Distributional convergence implies that a sequence of solutions $\{J^N\}$ converges distributionally to the distribution $\{J\}$ as $N \rightarrow \infty$, and we write [p. 107, ref. 18]

$$\lim_{N \rightarrow \infty} \langle J^N; \phi \rangle_H = \langle J; \phi \rangle_H \quad (12)$$

where ϕ is called a test function. A test function is infinitely differentiable and vanishes outside some bounded region. [The infinitely differentiable condition can be relaxed depending on the nature of the solution J (namely differentiability conditions)]. It is therefore observed that the distributional convergence no longer is a convergence of functions, but rather a convergence of numbers which may or may not be related to the norm of the square of the residuals, or to the norm of the solution error. This will be explained in the next section.

However, from a practical point of view, distributional convergence is a direct outcome of reciprocity, as we are not interested in the currents at any point, but the signal field at a point. So the problem is, what signal is produced by the currents at the required point? This information is extremely useful from a practical standpoint, but it does not provide a clue as to how many expansion functions are needed in (8) to reduce the error in the solution of (7).

For example, in an electrostatic problem, one is interested in computing the charge distribution on a structure for a given applied potential to evaluate its capacitance. As the number of expansion functions is increased, the convergence of the charge distribution on the structure is of great interest, as it has a direct relevance to the solution of the integral equation. The convergence of the value of the capacitance as a function of the number of expansion functions thus yields a distributional convergence for the problem of interest, as we are talking

about the integral of the charge. From the convergence of the value of capacitance it is difficult to extrapolate as to the exact nature of the charge distribution. However, from the behaviour of the charge distribution it is possible to accurately predict the nature of convergence of the capacitance. Therefore, the convergence of the charge distribution is defined as "strong convergence", and the convergence of the capacitance is defined as "weak convergence", or a distributional convergence.

1.3 The Method of Moments

The method of moments generalized the reaction concept, and introduced more flexible terminologies and test sources [9].

In the method of moments one starts with the operator equation, which describes the fundamental equation. We consider the same example that we utilized in the last section, i.e., scattering from a perfectly conducting body. The starting equation states that the total tangential electric field is zero on the surface of the conductor, i.e.,

$$\mathbf{E}^{tan} = \mathbf{E}_i^{tan} + \mathbf{E}_s^{tan} = 0 \quad (13)$$

In other words,

$$\mathbf{E}_s^{tan} = \mathbf{L}[\mathbf{J}_s] = -\mathbf{E}_i^{tan} \quad (14)$$

where \mathbf{L} is a linear operator. Operating on \mathbf{J}_s , it produces the electric field. Equivalently, in terms of an operator equation, (14) can be rewritten as

$$\mathbf{A}\mathbf{J} = \mathbf{Y} \quad (15)$$

We expand \mathbf{J}_s in terms of known expansion functions \mathbf{F}_n , with unknown coefficients a_n , so that

$$\mathbf{J}_s \cong \mathbf{J}_s^N = \sum_{n=1}^N a_n \mathbf{F}_n \quad (16)$$

We now substitute \mathbf{J}_s^N in the original operator equation and form the residual, \mathbf{R} , as

$$\mathbf{R} = \sum a_n \mathbf{L}[\mathbf{F}_n] + \mathbf{E}_i^{tan} = \mathbf{A}\mathbf{J}_s^N - \mathbf{Y} \quad (17)$$

Since the objective is to make $\sum a_n L[F_n]$ approximate $-E_i^{tan}$ in some sense, it then becomes obvious that the set $L[F_n]$ should be linearly independent and must be complete, so that E_i^{tan} can be approximated to any arbitrary degree of accuracy. It really is immaterial whether the expansion functions form an orthonormal set or not [10,11]! The residual, R in (17), is now weighted to zero with respect to some weighting functions, W_j , such that

$$\langle R; W_j \rangle_S = \iint R W_j ds = 0 \text{ for } j = 1, 2, \dots, N \quad (18)$$

The subscript S under the inner product defines a symmetric product, and not the conventional Hilbert inner product. In the conventional Hilbert inner product, the multiplication is defined with respect to the complex conjugate, i.e.

$$\langle R; W_j \rangle_H = \iint R \overline{W_j} ds = 0 \text{ for } j = 1, 2, \dots, N \quad (19)$$

By defining a symmetric product, the method of moments provides a generalization of the reciprocity principles.

If the weighting functions, W_j , are real, then the symmetric inner product and the Hilbert inner product become the same. However, problems arise when the weighting functions are complex. This may occur when one wants to utilize, for example, the method of least squares, where the weighting functions are chosen as

$$W_j = L[F_j] \quad (20)$$

for the solution of the operator equation (15). So, a true least squares problem cannot fit under the basic framework of the method of moments. This is because for the true least squares case we have to use the Hilbert inner product - the power concept - and not reciprocity. One might argue, however, that the weighting functions in (19) could be redefined as the conjugate of the weighting functions in (18). This is fine, but one should realize that one is now talking about power and not reciprocity, and this is the distinction we want to emphasize.

This philosophical dilemma also surfaces for another class of problems: analysis of electromagnetic scattering from bodies of revolution. There, it is often assumed that the circumferential variation of the

current is of the form $\exp(jm\phi)$, where m is an integer, $0, 1, \dots, M$, describing the order of the variation. The transverse variation of the current is approximated by pulses or triangles. For example, on an ogive, the current basis functions may be chosen as

$$J_s(\phi; t) = \sum_m \sum_n A_{mn} \exp(jm\phi) J_n(t) \quad (21)$$

It has been observed that for this type of expansion function, Galerkin's method, as defined in the method of moments context, is not applicable. In the methods of moments context, Galerkin's method requires taking the expansion functions to be identical to the weighting functions. However, it has been observed [12] that the choice $e^{-jm\phi}$ as weighting functions leads to stable results. Therefore the weighting functions are the conjugate of the basis functions [9, p.86].

The point here is that if one is going to choose the weighting functions as the conjugate of the expansion functions in the method of moment context, then it is natural to talk about power and not reciprocity. Hence, all the confusions of using a symmetric inner product would be cleared up. The symmetric inner product can only be utilized for real weighting functions, whereas the Hilbert inner product must be used when the weighting functions are complex. Therefore, why not start the solution procedure with only the Hilbert inner product, and not introduce the symmetric product at all?

In the solution of operator equations, "power conservation" is more important than satisfaction of reciprocity. In other words, the Hilbert inner product *must* be chosen, instead of the symmetric product. This is because the Hilbert inner product defines a norm, which can be used as a criterion for the measurement of error (for example, the error can be characterized by $\langle R; R \rangle_H = \|R\|^2$, whereas the symmetric product does not define a norm and, hence, no "a priori" measure of accuracy is possible.

Next, the nature of convergence of the method of moments is addressed. Like the "reaction concept", the "method of moments" guarantees only distributional convergence, or weak convergence, of the residuals. Neither the strong convergence of the residuals nor the convergence of the solution is guaranteed.

In the distributional convergence of the residuals, $\langle R; W_j \rangle_s = 0$ holds. Use of this error criterion can sometimes be disturbing. Take, for example, the situation when R is an oscillatory function and W_j

is a constant weight function. By enforcing (18), the integral will be zero even though \mathbf{R} at each point may be extremely large. This is why in all optimization schemes one attempts to minimize a quantity which is always positive. Hence, it would make sense to minimize the integral of $|\mathbf{R}^2|$, which is done in the method of least squares.

However, note that when the \mathbf{W}_j exist over a finite support, then the integral of the residual over that small region is equated to zero to solve for the unknown coefficients of the current. This may yield reasonable results, as long as the residual is of the same sign in that interval over which \mathbf{W}_j is not zero. Even though this makes sense from a practical point of view, (16) may not provide a better solution of (15) as N increases!

1.4 The Conjugate Gradient Method for Solution of Operator Equations

A significant difference between the application of the conjugate gradient method to the solution of operator equations, and the "reaction concept" and the method of moments, is that the conjugate gradient method guarantees strong convergence of the residuals, like the classical least squares approach. Strong convergence of the residuals implies mean square convergence, as opposed to weak convergence or distributional convergence. In this technique, the solution is upgraded by minimizing a functional, F , defined by

$$\begin{aligned} F(\mathbf{J}_a) &= \langle \mathbf{Q}\mathbf{R}; \mathbf{R} \rangle = \langle \mathbf{Q}\{\mathbf{L}[\mathbf{J}_a] + \mathbf{E}_i\}; \{\mathbf{L}[\mathbf{J}_a] + \mathbf{E}_i\} \rangle \\ &= \langle \mathbf{Q}(\mathbf{A}\mathbf{J} - \mathbf{Y}); (\mathbf{A}\mathbf{J} - \mathbf{Y}) \rangle_{\mathbf{H}} \end{aligned} \quad (22)$$

Here, \mathbf{Q} is a positive definite operator, which is assumed to be known. The approximate solution, \mathbf{J}_a is sought in the form

$$\mathbf{J}_a = \sum_n \mathbf{a}_n (\mathbf{A}^H \mathbf{A})^{n-1} \mathbf{A}^H \mathbf{Y} \quad (23)$$

where \mathbf{A}^H is the adjoint operator, and the unknown coefficients, \mathbf{a}_n , are selected in such a way that the positive functional $F(\mathbf{J}_a)$ is minimized. The advantage of choosing the expansion functions in the form defined by (23) is that the unknown coefficients, \mathbf{a}_n , can be computed recursively, instead of solving a matrix equation [13].

Observe that when $Q = I$, an identity operator, one is using the classical least squares techniques, and hence strong convergence of the residuals in the mean square sense is guaranteed, as we have

$$F_R(\mathbf{J}_s) = \langle \mathbf{R}; \mathbf{R} \rangle_H = \|\mathbf{R}\|^2 = \int R(z) \bar{R}(z) dz \quad (24)$$

It is interesting to note that when Q is chosen as $(\mathbf{A}^H \mathbf{A})^{-1}$, then the functional that is minimized is the following [13, case B, p. 1062]:

$$\begin{aligned} F_s(\mathbf{J}_s) &= \langle (\mathbf{J}_{\text{exact}} - \mathbf{J}_s); (\mathbf{J}_{\text{exact}} - \mathbf{J}_s) \rangle_H \\ &= \|\mathbf{J}_{\text{exact}} - \mathbf{J}_s\|^2 \end{aligned} \quad (25)$$

This means that the unknown coefficients, a_n , in (23) are chosen such that the error between the exact solution and the approximate solution is minimized at each iteration. From a philosophical point of view, this is very attractive, as this technique is theoretically guaranteed to give the exact solution as $N \rightarrow \infty$.

What the term $\mathbf{J}_{\text{exact}}$ means from a numerical standpoint will be explained in Section 1.5. The basic difference between the terms "theoretical" and "numerical" is that for the former we are utilizing infinite precision arithmetic, whereas for the latter we are using finite precision arithmetic, in the computations. Ill-posed problems often arise because of using finite precision arithmetic. Some of the consequences are outlined in [14].

1.4a. The difference between application of the conjugate gradient method (CG) to solve a matrix equation as opposed to direct application of CG to solution of an operator equation.

Recently, many communications have been devoted to this interesting topic. From the previous sections, it is clear that there are fundamental philosophical differences between applying the conjugate gradient method to solve the moment matrix equation of (18), as opposed to applying CG to minimize (24) recursively. The difference is that the latter guarantees strong convergence of the residuals and the solution (depending on case A or case B; [13]), whereas the former guarantees weak convergence. The question now arises, does that difference still exist when we look at the problem from a computational point of view? The answer is YES, and it is demonstrated by a simple example.

To illustrate this point, we start with an operator equation, $\mathbf{A}\mathbf{J} = \mathbf{Y}$, and convert it to a matrix equation utilizing the method of moments principles.

$$\mathbf{A}_M \mathbf{J}_M = \mathbf{Y}_M \quad (26)$$

where \mathbf{A}_M , \mathbf{J}_M , and \mathbf{Y}_M are the matrix version of the continuous operator equation. If we now solve the matrix version in (26) by the conjugate gradient method, then we are essentially solving the following equations, known as the normal equations:

$$\mathbf{A}_M^H \mathbf{A}_M \mathbf{J}_M = \mathbf{A}_M^H \mathbf{Y}_M \quad (27)$$

where \mathbf{A}_M^H is the adjoint matrix of \mathbf{A}_M (which is simply the conjugate transpose). Now, if we apply the conjugate gradient method directly to the solution of the operator equation we solve

$$\mathbf{A}^H \mathbf{A}\mathbf{J} = \mathbf{A}^H \mathbf{Y} \quad (28)$$

Here \mathbf{A}^H is the adjoint operator. For numerical computation, (28) has to be discretized, and we obtain

$$\mathbf{A}_D^H \mathbf{A}_D \mathbf{J}_D = \mathbf{A}_D^H \mathbf{Y}_D \quad (29)$$

Since $\mathbf{A}_D = \mathbf{A}_M$, $\mathbf{J}_D = \mathbf{J}_M$ and $\mathbf{Y}_D = \mathbf{Y}_M$. Therefore, during numerical implementation, the basic difference between (27) and (29) will be as to how the continuous adjoint operator \mathbf{A}_H has been discretized to \mathbf{A}_D^H and whether the "matricised" adjoint operator \mathbf{A}_M^H is the conjugate transpose of the discretized original operator, \mathbf{A}_D . If the two matrices (namely, the discretized adjoint operator \mathbf{A}_D^H and the conjugate transpose of the "matricised" operator \mathbf{A}_M^H) are not identical, then there will be a difference between the application of the CG method to the solution of a matrix equation, as opposed to the solution of the operator equation. For many problems, (27) and (29) are not identical.

As an example, consider the convolution equation [15]

$$\int \mathbf{x}(t-u) \mathbf{v}(u) du = \mathbf{y}(t), \text{ for } 0 \leq t < \infty \quad (30)$$

where $\mathbf{x}(t)$ and $\mathbf{y}(t)$ are assumed to be known. We redefine (30) as $\mathbf{A}\mathbf{V} = \mathbf{Y}$. This operator equation can now be written in the matrix form, utilizing the method of moments concept of pulse basis

functions for the unknown \mathbf{v} , and with impulse weighting, to give $\mathbf{A}_M \mathbf{V}_M = \mathbf{Y}_M$, or as

$$\begin{bmatrix} x_1 & & & & \\ & x_1 & & & \\ & \bullet & & & \\ & \bullet & & & \\ x_N & x_{N-1} & \bullet & x_{N-M} & \end{bmatrix}_{N \times M} \begin{bmatrix} v_1 \\ v_2 \\ \bullet \\ \bullet \\ v_M \end{bmatrix}_{M \times 1} = \begin{bmatrix} y_1 \\ y_2 \\ \bullet \\ \bullet \\ y_N \end{bmatrix}_{N \times 1} \quad (31)$$

If we now apply the conjugate gradient method to solve this matrix equation we will be solving $\mathbf{A}_M^H \mathbf{A}_M \mathbf{V}_M = \mathbf{A}_M^H \mathbf{Y}_M$. The adjoint matrix \mathbf{A}_M^H in this case is

$$\mathbf{A}_M^H = \begin{bmatrix} x_1 & x_2 & \bullet & \bullet & x_N \\ & x_1 & \bullet & \bullet & x_{N-1} \\ & & & & \bullet \\ & & & & \bullet \\ & & & x_1 & x_{N-M} \end{bmatrix} \quad (32)$$

We get a solution for \mathbf{v} (in this case) by solving (27) by CG. Now let us apply the conjugate gradient method directly to the operator equation of (30). We consider the adjoint operator for the integral equation of (30). The adjoint operator, \mathbf{A}^H , of \mathbf{A} , is defined by

$$\langle \mathbf{A}\mathbf{V}; \mathbf{Z} \rangle_H = \langle \mathbf{V}; \mathbf{A}^H \mathbf{Z} \rangle_H \quad (33)$$

or equivalently,

$$\int \mathbf{z}(t) dt \int \mathbf{x}(t-u) \mathbf{v}(u) du = \int \mathbf{v}(u) du \int \mathbf{z}(t) \mathbf{x}(t-u) dt \quad (34)$$

So the adjoint operator in this case is the advance convolution operator. By comparing (32) and (34) it is apparent that \mathbf{A}_M^H is a restricted version of the actual adjoint operator in (34). This is because the adjoint operator is

$$\mathbf{A} \star \mathbf{z} = \int \mathbf{z}(t) \mathbf{x}(t-u) dt \quad (35)$$

This continuous operator can only take the form of (32) under the assumption that $z(n)$ is identically zero beyond n , and this may not be true, in general. Therefore, A_D^H and A_M^H are not identical for this problem unless some additional assumptions are made! In the signal processing literature, a clear distinction has been made between (27) and (29). In (27), $A_M^H \bullet A_M$ is called the covariance matrix of the data, whereas $A_D^H \bullet A_D$ in (29) is called the autocorrelation matrix of the data. It is well known that (27) and (29) yield quite different results. In the electromagnetics literature, several techniques have been developed to solve (29). These techniques depend on the assumptions that need to be made regarding the behaviour of x and z for $t \geq n$. One possible assumption may be that the waveforms have become almost zero for $t \geq n$. This has been implemented by Tjihuis [16]. An alternate numerical implementation has been considered by Tseng and Sarkar [15], where instead of assuming the nature of decay for x and z for $t \geq n$, a weighted inner product has been defined to minimize the error introduced in discretizing the continuous adjoint operator.

The adjoint operator provides physical insight into the system. The operator A tells us how the system will behave for a given external source. The adjoint operator, on the other hand, tells us how the system responds to sources in general. In short, for causal systems (i.e. $x(n) = 0$ for $n < 0$), the application of the conjugate gradient method to the direct solution of operator equations is quite different from the application of the conjugate gradient method to the solution of matrix equations. However, when the operator A exists ($x(n)$, in this case) from $-\infty$ to $+\infty$, and is even and symmetric about the origin, then the two techniques may yield similar results.

However, for nonequally sampled data, no generalization can be made.

1.5 Development of a Numerical Method

So far the techniques have been presented from a heuristic point of view. Next we look at a systematic development of the boundary value problem.

Suppose one is interested in the solution of the scattering of electromagnetic fields from a conducting body, when it is illuminated by the incident field. This is the same problem that was addressed in [4].

It is now desirable to delve more deeply into the philosophy of setting up the problem. In the solution procedure, one replaces the conducting structure by an equivalent electric current in free space. The current is located at the same position as that of the conductor. In the next step, the scattered fields are computed from this equivalent electric current in free space by the operator $L(J)$, where L is the operator acting on the electric current which produces the scattered electric field, E^{sc} . One then applies the continuity of the total tangential electric field on the conductor surface, namely $E^{sc} + E^{inc} = 0$. The rationale for matching the tangential fields on the conductor surface is given by the Uniqueness Theorem: "A field in a lossy region is uniquely specified by the source within the region plus the tangential components of E over the boundary ... Note that our uniqueness breaks down for dissipationless media. To obtain uniqueness in this case, we consider the field in a dissipationless medium to be the limit of the corresponding field in a lossy medium as the dissipation goes to zero..." [11, p.102]. Therefore, the electric field everywhere in space can be obtained from [the current] J , once we match the tangential components of the fields and solve the boundary value problem $L(J) = -E^{inc}$ on the surface of the conductor. Therefore, if one's objective is to find a solution to this problem, then one must find a method to match $L(J)$ equal to $-E^{inc}$ at all possible locations. Therefore the ideal situation is to seek a solution J_∞ such that the error

$$E_\infty = \max |L(J_\infty) + E^{inc}| \quad (36)$$

is minimized. The error E_∞ in the solution is defined as follows: One computes $L(J)$ at all positions z on the surface of the conductor, and compares it with $-E^{inc}$. If the computed field $L(J_\infty)$ does not match $-E^{inc}$ at all positions z then an error is generated. We now choose the maximum absolute value of this error, which occurs perhaps at the point z , or may even occur at several points. So one seeks a solution J_∞ so the E_∞ is 0.001 (say, for example). Even though this is the ideal solution, practically, there is no solution technique that solves this nonlinear minimization problem in an efficient way. So for this reason, an alternate error criterion is sought. The next best choice for the error is

$$E_1 = \int_z |L(J_1) + E^{inc}| dz \quad (37)$$

In this, the error is defined by the integral of the absolute value of the error over the entire surface. So the objective in this case would be to

find a solution \mathbf{J}_1 such that \mathbf{E}_1 is 0.001 (say). Clearly, this is a much more relaxed error criterion than \mathbf{E}_∞ . This is because in \mathbf{E}_1 , the actual error (the difference in the tangential field), may be large. So even though a large difference may exist over a small region the integral in \mathbf{E}_1 is small. Then the solution \mathbf{J}_1 is acceptable. This also gives rise to a nonlinear solution procedure. Several well-known techniques exist to solve these types of minimization problems, of which the algorithm by Nedler and Mead [12] is notable. Recently, Karmarker claims to have developed an efficient algorithm to solve this class of problems.

An alternate error criterion, which leads to a linear minimization problem, is the classical "least squares" method. In this case, the error is defined as

$$\mathbf{E}_2 = \int_z |\mathbf{L}(\mathbf{J}_2) + \mathbf{E}^{inc}|^2 dz \quad (38)$$

So a solution \mathbf{J}_2 is sought which minimizes the least square error. Since this error criterion leads to a linear matrix problem, it has been very popular in science and engineering. One possible way to obtain a solution \mathbf{J}_2 is to assume that the solution \mathbf{J}_2 is of the form of (16) where \mathbf{F}_i are certain expansion functions which are known, and a_i are the unknown weights to be solved for. Now \mathbf{J}_2^N is substituted in \mathbf{E}_2 and we form

$$E_2 = \int_z \left| \sum_{i=1}^N a_i \mathbf{L}\mathbf{F}_i + \mathbf{E}^{inc} \right|^2 dz \quad (39)$$

We next minimize E_2 to solve for a_i . This is accomplished by

$$\frac{\partial E_2}{\partial a_i} = 0 \quad (40)$$

This leads to a set of equations

$$\sum_{i=1}^N a_i \langle \mathbf{L}\mathbf{F}_i; \mathbf{L}\mathbf{F}_i \rangle_H = - \langle \mathbf{E}^{inc}; \mathbf{L}\mathbf{F}_i \rangle_H \quad (41)$$

from which the a_i 's are computed. This raises an interesting question as how to choose the expansion functions \mathbf{F}_i . Since the objective is to minimize the error $(\sum_{i=1}^N a_i \mathbf{L}\mathbf{F}_i + \mathbf{E}^{inc})$, therefore it is essential that $\mathbf{L}\mathbf{F}_i$ must be linearly independent. Please note that orthogonality of

the basis functions has nothing to do with the actual solution to the problem [10,11]. The emphasis is on the completeness of \mathbf{LF}_i , as this is the requirement to solve the problem. The least squares solution also stipulates that the error $(\sum_{i=1}^N a_i \mathbf{LF}_i + \mathbf{E}^{inc})$ will be orthogonal to $\sum_{i=1}^N a_i \mathbf{LF}_i$. An alternative way of viewing this problem is to rewrite (7) as

$$\begin{aligned} \sum_{i=1}^N a_i < \mathbf{LF}_i; \mathbf{LF}_i >_H + < \mathbf{E}^{inc}; \mathbf{LF}_i >_H = \\ < \sum_{i=1}^N a_i \mathbf{LF}_i + \mathbf{E}^{inc}; \mathbf{LF}_i >_H = \\ < -\mathbf{E}_N^{inc} + \mathbf{E}^{inc}; \mathbf{LF}_i >_H = \\ < \text{error}; \mathbf{LF}_i >_H = 0 \end{aligned} \quad (42)$$

Therefore, the functions \mathbf{LF}_i are orthogonal to the error space $\mathbf{E}^{inc} - \mathbf{E}_N^{inc}$. For orthogonality of two complex vectors, it is necessary to use the Hilbert Inner Product, thus the nature of the inner product becomes obvious. One can never use the symmetric product, as it cannot define orthogonality between two functions. One has to use the classical Hilbert product. Why? Because this is the methodology to define a stable numerical procedure. The definition of the inner product is secondary, the primary objective being to solve $\mathbf{L}(\mathbf{J}) + \mathbf{E}^{inc} = 0$.

The classical least squares approach can be quite time consuming as, from a purely computational point of view, evaluation of $< \mathbf{LF}_i; \mathbf{LF}_i >$ is very elaborate. The question now arises, can the computational efficiency be improved without sacrificing the scientific methodology; i.e., as one increases the summation of \mathbf{J}^N in (8) from N to $N + 1$, one is indeed guaranteed to have a better solution. In other words, the error is orthogonal to $-\mathbf{E}_N$. This is where the concept of weighting functions \mathbf{W}_i comes in. And so the weighting functions are defined by $\mathbf{LF}_i = \mathbf{W}_i$. The job of the weighting functions is:

(I)

$$S(\mathbf{W}_i) = S(\mathbf{LF}_i) \perp S(\text{error}) \quad (43)$$

i.e., the space of weighting functions is orthogonal to the error space for best approximation, and

(II)

$$-\mathbf{E}^{inc} \in S(\mathbf{W}_i) \quad (44)$$

In the limit, the excitation is an element in the space of weighting functions.

So now one introduces a class of weighting functions satisfying the above two rules so that

$$\sum_{i=1}^N a_i \langle \mathbf{L}\mathbf{F}_i; \mathbf{W}_i \rangle_H = - \langle \mathbf{E}^{inc}; \mathbf{W}_i \rangle_H \quad (45)$$

The solution given by the above equation will be identical to the least squares solution from a computation point of view if the weighting functions have the two required properties and the Hilbert inner product is used. So now one could choose some simple, known, analytic functions for \mathbf{W}_i such that $S(\mathbf{W}_i) = S(\mathbf{L}\mathbf{F}_i)$. Thereby, one gets the least squares solution with less computation. From a theoretical point of view it may be quite difficult to select \mathbf{W}_i such that $S(\mathbf{W}_i) = S(\mathbf{L}\mathbf{F}_i)$; for $i = 1, \dots, \infty$. However, from a computational point of view, any set of N linearly independent functions would suffice so that $S(\mathbf{L}\mathbf{F}_i) = S(\mathbf{W}_i)$ for $i = 1, 2, \dots, N$. Under the above assumptions, equation (45) can be interpreted as an integral of the error. However, to obtain a meaningful solution, certain constraints are necessary, as outlined above. Unfortunately, one often starts with equation (45) and interprets it utilizing reaction technique and symmetric products. However, neither the reaction technique nor the symmetric products have anything to do with the solution of the scattering problem defined in (45). The equation (45) represents that the integral of the error with respect to the weighting function is zero. If one interprets equation (45) as the integral of the error and one integrates an error function which has large oscillations, then even though the integral is zero, the actual errors may be quite large.

Therefore from a computational viewpoint, the recipe is quite different from that of a purely philosophical analysis presented in the previous sections.

1.6 Numerical Considerations

In the last sections we have observed that, from a philosophical standpoint, there is a significant difference between the reaction concept, the method of moments, and the conjugate gradient method. And yet, numerical computations in many cases (except body of revolution type problems, where the weighting function is complex) have

demonstrated that the numbers yielded by the three techniques are often identical to several decimal places. This raises the question as to why they yield similar results, even though each technique starts from a totally different philosophical origin. In the solution of $\mathbf{AJ} = \mathbf{Y}$, the unknown \mathbf{J} is approximated by

$$\mathbf{J}_a = \sum_n a_n \mathbf{F}_n$$

and the residual is formed as

$$\mathbf{R} = \mathbf{AJ} - \mathbf{Y} = \sum_{n=1}^N a_n \mathbf{AF}_n - \mathbf{Y} \quad (46)$$

The residual is weighted to zero by some weighting functions \mathbf{W}_j , so that

$$\langle \mathbf{R}; \mathbf{W}_j \rangle_H = \sum_{n=1}^N a_n \langle \mathbf{AF}_n; \mathbf{W}_j \rangle_H - \langle \mathbf{Y}; \mathbf{W}_j \rangle_H = 0$$

for $j = 1, \dots, M$ (47)

Here we assume the inner products to be the Hilbert inner product. However, the final conclusion will be independent of the choice of the inner product!

Generally, we have to utilize some sort of quadrature formula to evaluate the inner products in (47), as the inner products in many cases cannot be evaluated analytically. In that case, we have

$$\langle Q; W_j \rangle_H = \sum_k c_k Q(z_k) \overline{W_j}(z_k) \quad (48)$$

where the z_k are the points at which the functions in the inner products are evaluated, and the c_k are the quadrature weights. Now (47) can be rewritten in a factored matrix form [17]:

$$\begin{bmatrix} \overline{w}_1(z_1) & \overline{w}_1(z_2) & \cdots & \overline{w}_1(z_k) \\ \overline{w}_2(z_1) & \overline{w}_2(z_2) & \cdots & \overline{w}_2(z_k) \\ \vdots & \vdots & & \vdots \\ \overline{w}_M(z_1) & \overline{w}_M(z_2) & \cdots & \overline{w}_M(z_k) \end{bmatrix} \bullet \begin{bmatrix} c_1 & & & \\ & c_2 & & \\ & & \ddots & \\ & & & c_k \end{bmatrix} \bullet$$

$$\begin{bmatrix} AF_1(z_1) & AF_2(z_1) \cdots & AF_N(z_1) \\ AF_2(z_2) & AF_2(z_2) \cdots & AF_N(z_2) \\ \vdots & \vdots & \vdots \\ AF_1(z_k) & AF_2(z_k) \cdots & AF_N(z_k) \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} = \\
 \begin{bmatrix} \bar{w}_1(z_1) & \bar{w}_1(z_2) \cdots & \bar{w}_1(z_k) \\ \bar{w}_2(z_1) & \bar{w}_2(z_2) \cdots & \bar{w}_2(z_k) \\ \vdots & \vdots & \vdots \\ \bar{w}_M(z_1) & \bar{w}_M(z_2) \cdots & \bar{w}_M(z_k) \end{bmatrix} \cdot \begin{bmatrix} c_1 & & \\ & c_2 & \\ & & \ddots \\ & & & c_k \end{bmatrix} \cdot \begin{bmatrix} Y(z_1) \\ Y(z_2) \\ \vdots \\ Y(z_k) \end{bmatrix} \quad (49)$$

Equivalently,

$$\begin{bmatrix} \bar{w}_1(z_1) & \bar{w}_1(z_2) \cdots & \bar{w}_1(z_k) \\ \bar{w}_2(z_1) & \bar{w}_2(z_2) \cdots & \bar{w}_2(z_k) \\ \vdots & \vdots & \vdots \\ \bar{w}_M(z_1) & \bar{w}_M(z_2) \cdots & \bar{w}_M(z_k) \end{bmatrix} \cdot \begin{bmatrix} c_1 & & \\ & c_2 & \\ & & \ddots \\ & & & c_k \end{bmatrix} \cdot \\
 \left\{ \begin{bmatrix} AF_1(z_1) & AF_2(z_1) \cdots & AF_N(z_1) \\ AF_2(z_2) & AF_2(z_2) \cdots & AF_N(z_2) \\ \vdots & \vdots & \vdots \\ AF_1(z_k) & AF_2(z_k) \cdots & AF_N(z_k) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} - \begin{bmatrix} Y(z_1) \\ Y(z_2) \\ \vdots \\ Y(z_k) \end{bmatrix} \right\} = 0 \quad (50)$$

If the square matrices $[C]_{K \times K}$ and $[W]_{K \times K}$ are not singular (for a general quadrature rule $[C]$ cannot be singular, and since the weighting functions are different $[W]$ cannot be singular), then one is essentially solving a set of point-matched equations, irrespective of what the starting point may have been: the "reaction concept", the "method of moments", or the "conjugate gradient". For the first two methods, the weighting functions are redefined in (50) when the symmetric products are utilized.

The interesting point in (50) is that the final equation is independent of the choice of weighting functions, provided the inner product in (47) is computed by a quadrature rule.

There are several comments that need to be made regarding (50).

- A. When one utilizes point matching in the conventional method of moments context, then $K = N$, and one is solving a square matrix equation. This is really the worst case.

- B. When one utilizes Galerkin's method in the method of moments, then K is larger than N , and one is then solving a set of weighted point-matched equations. Since $K > N$, the result will be better. This is the scientific reason why Galerkin's method always provides a better result than point matching, defined by case A.
- C. It is rather interesting that if one had utilized only point matching in the method of moments with $K > N$, and with the match points being placed at the identical points where the Galerkin's method of case B will sample the residual R , then one would obtain a rectangular matrix equation. If this rectangular matrix is solved in a least squares sense, one would obtain a better solution than Galerkin's method, and this solution would be the least squares solution.
- D. The conjugate gradient method and the method of moments will yield similar results, up to a few decimal places, for the case $K = N$. However, if in the conjugate gradient method, the inner products are evaluated at more points, with $K > N$, then the solution given by the conjugate gradient method will be identical to that of the method of least squares.
- E. Note that numerical results will be different from technique to technique depending on how AJ_n is evaluated.
- F. If analytical integration is utilized instead of numerical quadrature in the evaluation of the inner products in (47), then this discussion is no longer valid!
- G. Note that from a computational point of view, the final results obtained by the two versions of the conjugate gradient method are identical (namely, case A and case B [13]); however, the numerical procedures utilized to obtain the final result are quite different.

In numerical computation there are three variables that are at our disposal. These are K , the number of points at which the residual is sampled; M , the number of weighting functions; and N , the number of expansion functions. Often, for simplicity, M is chosen to be equal to N . So, suppose we are given K and N . A question then arises as to how the solution behaves as a function of K and N . Given a fixed K , as we increase N , the sequence of solutions for the conjugate gradient method monotonically converges to the best solution that can be obtained with a fixed K . This solution is termed J_{exact} in (25). In the limit as K (the number of match points) approaches infinity, the sequence of solutions called J_{exact} will converge to the solution in the infinite-dimensional

space. So the definition of $\mathbf{J}_{\text{exact}}$ in (25) is a function of K , the number of points at which the residuals have been sampled. However, this convergence may not be monotonic, as the definition of the functional is changed for each K . Monotonic convergence for a fixed K , as N increases, is not guaranteed by either the "reaction concept", "point matching", or "Galerkin's method". Only the method of least squares and the conjugate gradient method guarantee monotonic convergence.

1.7 A Note on Computational Efficiency

Even though from a purely computational point of view, the theoretical advantages of the conjugate gradient method over conventional numerical techniques are lost, as all numerical techniques perform weighted point matching anyway, the conjugate gradient technique may still be quite efficient for solving certain class of matrix equations.

The class of matrix problems, where the application of the conjugate gradient method may lead to significant savings of CPU time, are the Hankel matrices of which Toplitz matrices is a subset. Most of the CPU time in the utilization of the conjugate gradient method is used in the computation of the operator/ matrix products $\mathbf{A}\mathbf{P}_n$ and $\mathbf{A}^H\mathbf{R}_n$ [13]. However, if the matrices \mathbf{A} belong to the Hankel system then FFT (Fast Fourier Transform) can be utilized in reducing the CPU time significantly. This is because the computation of $\mathbf{A}\mathbf{P}_n$ and $\mathbf{A}^H\mathbf{R}_n$ are typically convolutions. Hence the application of a FFT and an Inverse FFT would be more efficient in computing the matrix products. This is because the computation of $\mathbf{A}\mathbf{P}_N$ and $\mathbf{A}^H\mathbf{R}_n$ via the FFT route is $\Theta(2N\log N)$ process as compared to the usual $\Theta(N^2)$ matrix products.

Solution of Hankel Systems by Trench's algorithm typically is $\Theta(N^2)$ process. This means as the size of the matrix increases the CPU time increases by the square of the dimension of the matrix. This is in contrast to $\Theta(N^3)$ for arbitrary matrix equations. However, when FFT and conjugate gradient is utilized, it is seen that the CPU time increases as $\Theta(N)$ rather than as $\Theta(N^2)$ of contemporary techniques [19].

In Figure 1.1, the CPU time for solving the electromagnetic scattering from broadside incident of a 2.5λ antenna is presented. It is seen from Fig. 1.1, that for 100-900 unknowns, the CPU time increased linearly with the number of unknowns. [The CPU time represents the

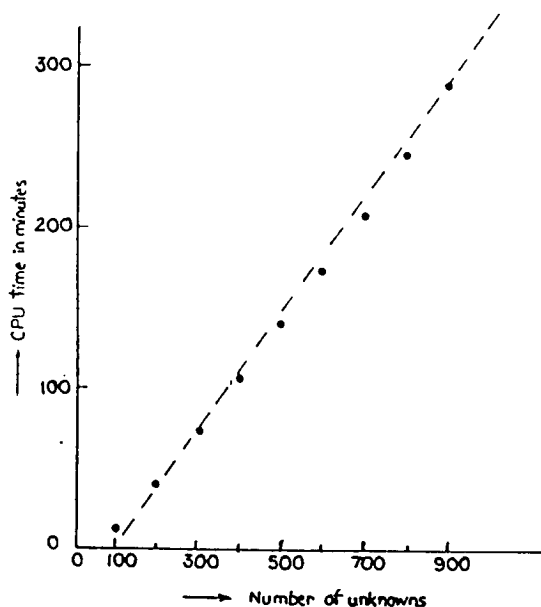


Figure 1.1 CPU time against the number of unknowns for the Ordinary Conjugate gradient method.

total time taken for the program to run on a Microvax 2 computer. A Microvax 2 is about 90% efficient of a VAX 11/780.] The interesting point is that the CPU time increases linearly with the number of knowns.

This linearity is still exhibited for 2D-problems as described in [19]. The same linearity has been observed for 3D problems, when “page faults” are not a major problem in the computation.

Fortran programs for the various versions of the conjugate gradient method are available in [20].

1.8 Conclusion

It has been shown that the “reaction concept” and the “method of moments” are derived from the principle of reciprocity, whereas the conjugate gradient method is derived from the concept of power conservation. Hence, the symmetric products are utilized in the method of moments, whereas the Hilbert inner product is utilized in the formula-

tion of the conjugate gradient techniques. For real weighting functions, it really does not make any difference as to whether the symmetric or the Hilbert product is chosen. However, for complex weighting functions, the Hilbert inner product *must* be chosen as the symmetric product does not define a norm. Hence, the conjugate gradient method provides strong convergence for the residuals, whereas the "reaction concept" and the method of moments provide weak convergence.

Depending on the Green's function, the application of the conjugate gradient method to the solution of a matrix equation can yield quite different results, as opposed to the application of the conjugate gradient method to the direct solution of the operator equation. The difference becomes quite obvious for the solution of integral equations in the time domain and for nonequally-sampled functionals.

However, from a purely computational point of view, all techniques essentially yield a solution corresponding to a weighted point-matched technique, if the inner products are evaluated using a numerical quadrature. If the inner products are evaluated analytically, then the above conclusion does not hold.

References

- [1] Wang, J. J. H., "Generalized moment methods in electromagnetics," *IEEE Proc.*, **137**, pt H, 2, 127-132, April 1990.
- [2] Rumsey, V. H., "Reaction Concept in Electromagnetic Theory," *Phys. Rev.*, **94**, 1483-1491, June 1954.
- [3] Schelkunoff, S. A., and H. Friis, *Antenna Theory and Practice*, New York, John Wiley and Sons, 1952.
- [4] Cohen, M. R., "Application of the reaction concept to scattering problems," *IRE Trans. Ant. Prop.*, 193-199, October 1955.
- [5] Carson, J. R., "Reciprocal theorems in radio communication," *Proc. IRE*, **17**, 952-956, June 1929.
- [6] Carter, P. S., "Circuit relations in radiating systems and applications to antenna problems," *Proc. IRE*, **20**, 1004-1041, June 1932.
- [7] Wang, N., J. H. Richmond, and M. C. Gilreath, "Sinusoidal re-

- action formulation for radiation and scattering from conducting surfaces," *IEEE Trans. Ant. Prop.*, AP-23, 376-382, May 1975.
- [8] Richmond, J., "A reaction theorem and its applications to antenna impedance calculations," *IRE Trans. Ant. and Prop.*, 515-520, November 1961.
- [9] Harrington, R. F., *Field Computation by Moment Methods*, Kreiger Publications, 1985.
- [10] Sarkar, T. K., "A note on the choice of weighting functions in the method of moments," *IEEE Trans. Ant. Prop.*, 436-441, April 1985.
- [11] Sarkar, T. K., A. R. Djordjevic, and E. Arvas, "On the choice of expansion and weighting functions in the numerical solution of operator equations," *IEEE Trans. Ant. Prop.*, 33, 988-996, September 1985.
- [12] Mautz, J. R., and R. F. Harrington, "H-field, E-field and combined field solutions for bodies of revolution," *AEU*, 32, 159-164, 1978.
- [13] Sarkar, T. K., and E. Arvas, "On the class of finite step iterative methods (Conjugate Directions) for the solution of an operator equation arising in electromagnetics," *IEEE Trans. on Ant. Prop.*, 33, 1058-1066, October 1985.
- [14] Sarkar, T. K., D. D. Weiner, and V. K. Jain, "Some mathematical considerations in dealing with the inverse problem," *IEEE Trans. Ant. Prop.*, 29, 373-379, March 1981.
- [15] Tseng, F. I., and T. K. Sarkar, "Deconvolution of impulse response of a conducting sphere by the conjugate gradient method," *IEEE Trans. Ant. Prop.*, 35, 105-110, January 1987.
- [16] Tijhuis, A. G., *Electromagnetic Inverse Profiling*, VNU Science Press, Utrecht, the Netherlands, 1987
- [17] Djordjevic, A. R., and T. K. Sarkar, "A theorem on the method of moments," *IEEE Trans. Ant. Prop.*, 35, 353-355, March 1987.
- [18] Stakgold, I., *Green's Functions and Boundary Value Problems*, J. Wiley and Sons, New York, 1979.
- [19] Sarkar, T., "On the application of the generalized biconjugate gradient method," *Journal of Electromagnetic Waves and Applications*, 1, No. 4, 325-345, 1987.

- [20] Sarkar, T., X. Yang and E. Arvas, "A limited survey of various conjugate gradient methods for solving complex matrix equations arising in electromagnetic wave interactions," *Wave Motion*, 10, 527–546, 1988.