

A Radically Different Method of Moments

William A. Gardner

Department of Electrical and Computer Engineering, University of
California Davis, Orcid number: 0000-0003-3840-7191

January 5, 2023

Abstract

A method of parameter estimation using only specified moments of the observed data is described. It is radically different from the classical method of moments (MoM) introduced at the end of the 19th Century and shows promise for being competitive. The alternative method uses estimates of posterior PDF values of the unknown parameters — estimates that are constrained to be linear combinations of specified nonlinear transformations of the observed data. These estimates are the solutions to linear equations specified in terms of first- and second-order moments from a probabilistic model of the nonlinearly transformed data. For polynomial nonlinearities up to the order n , these are equivalent to moments of the observed random variables up to the order $2n$, revealing that this general method includes an alternative MoM as a special case; however, in place of the sample moments of the data used along with the probabilistic moments in the Classical MoM, more general weighted averages of products of the data with itself are used, and the weighting functions are optimized according to a Bayesian minimum-risk criterion. The solution for the posterior PDF estimate is studied analytically. Results are encouraging.

Keywords: Parameter estimation, Multivariate Model Fitting, Methods of Moments, Bayesian Inference

1 Introduction

The traditional Method of Moments is said to have been introduced by Pearson (1936) and also by Chebyshev in 1887 (see Wikipedia (2022)). This method, when applied to either multivariate or time-series data, consists of equating sample moments measured from the data with theoretical moments obtained from a probabilistic model of the time series, and then solving for the unknown values of parameters in the theoretical moment expressions. The theoretical moments can be interpreted as unconditional moments depending on unknown parameters, or moments conditioned on unknown values of random parameters. The choice of interpretation has no impact on the method. However, the latter interpretation can be used to formulate a radically new approach to parameter estimation based on concepts from Bayesian Inference.

The classical MoM is a mainstay of parameter estimation for probabilistic models of data in econometrics and other fields for which knowledge of the likelihood function is often unavailable or complexity of the known likelihood function prevents its use for maximum-likelihood estimation.

In the Method of Moments, one can use the mean and centralized moments or the mean and non-centralized moments, and one can use as many moments as there are unknowns, and there are other variations that have been devised. One such variation uses the fact that the theoretical moments for an M -th order autoregressive time series model satisfy a set of $M + 1$ linear equations in $M + 1$ unknowns involving only 2nd-order moments, the autoregressive model coefficients, and these equations can be solved for these unknown coefficients. This method is very common in data modeling and time-series prediction and associated studies of causality. More generally, however, the Method of Moments requires the solution of nonlinear equations.

In the alternative approach, conditional moments are used and the objective is not to match theoretical moments to samples moments but rather to estimate the posterior PDF of the unknown parameters using the observed data, and then select the values of the conditioning parameters that maximize the estimated posterior PDF, thereby obtaining the maximally “probable” solution for the parameter values, where the quotation marks denote the fact that the PDF used is only an estimated PDF.

In this alternative method, one can use moments of a linear combination of any user-specified nonlinear functions of the observations; the equations to be solved are always linear, regardless of the particular functional dependence of the theoretical conditional moments on the parameters. But, when polynomial nonlinearities are used, the higher the order of the polynomials, the higher the order of the moments required. The posterior PDF estimator requires orders of moments up to twice the order of the polynomial.

The alternative method is optimal in the Bayesian sense that its estimate of the posterior PDF is a minimum mean-squared-error estimate subject to the selected constraint on the nonlinearities used.

This method requires calculating the sum of the moments, conditioned on the unknown parameter values, weighted by a prior PDF for those parameters. But one can always use a uniform PDF over a sufficiently large finite region of the domain if there is no knowledge of a prior PDF. This is tantamount to switching from a Bayesian approach to a Maximum-Likelihood (ML) approach, since the posterior PDF, as a function of the unknown parameters, is proportional to the Likelihood Function over the admissible region of the domain of the uniform prior PDF. However, the ML approach here is still only Maximum-“Likelihood” because the likelihood function used is only an estimated likelihood function.

This alternative method can use a single sample of a stationary (or cyclostationary) sequence of random variables, which favors applications to time-series analysis, or it can use multiple samples of one or more random variables.

To provide the reader with a look ahead at what this new method offers in comparison with the classical MoM, Table 1 summarizes the key characteristics of these two methods.

Expanding on Table 1, the Radically Different MoM produces more than just a single estimate for each parameter; it produces an estimate of the posterior PDF of the unknown parameters. This “PDF” can be used to calculate most “probable” estimates of the unknown parameters (their posterior “PDF” modes) which equal maximum-“Likelihood” estimates when the prior PDF is chosen to be uniform, or it can produce minimum-“mean”-squared-error estimates (their posterior “means”) or minimum-“mean”-absolute-value of error estimates (their posterior “median”) where, in all cases, the quotation marks denote the fact that the posterior PDF used in the estimates is the MMSE estimate of the true posterior PDF, subject to a user specified constraint on the structure (functional form) of that estimator’s dependence on the observed data.

Prior to recognizing the applicability of this method to multivariate statistics, it was devised and used to design linear and quadratic communications receivers for digitally modulated signals and was found to have strong resemblances to statistically optimum receivers that are linear or quadratic under the simplifying assumption (for the optimum receiver) of additive Gaussian noise Gardner (1973)-Gardner (1976*a*).

The purpose of this article is to show that this method is promising for not only statistical inference based on single samples of time series data but also for multivariate statistical inference based on multiple samples in general. The particular applications studied in the original papers Gardner (1973)-Gardner (1976*a*), which focused on data communications

Table 1: Advantages of the Radically Different MoM

#	Comparison Basis	Classical Method of Moments (MoM)	Radically Different MoM
1.	Functionals of Data Used	Uses sample moments	Uses linear combinations of any specified functionals of samples, including for example optimally weighted sample moments
2a.	Model Used	Probabilistic conditional (on the parameters) moments of the data	Probabilistic conditional 1st and 2nd order moments of specified functionals of the data
2b.	Model Restrictions	Must use only moments for which the equations are solvable	Requires a data model in which the unknown parameters appear explicitly; the equations will always be solvable
3.	Nature of Equations to be Solved	Generally nonlinear, except for Auto Regressive models	Always linear
4.	Breadth of Optimality Criteria	Produces a single solution with no optimality properties in general	Produces different solutions for different choices of optimality criteria
5.	Convergence to ML or Min-Risk Estimate	Has no general relationship to ML or Min-Risk estimates, except asymptotically as the amount of data grows without bound	Converges to Bayesian Min-Risk estimate (when prior PDF is known) or ML estimate of parameters as order of polynomial estimator increases for any fixed finite amount of data.
6.	Use of Prior Information	Does not use prior information	Uses prior information in an optimal manner when available
7.	Number of Samples Used for Each Random Variable	Typically, as many as possible	One (e.g., for long time series) or many
8.	Philosophy of Approach	Purely ad hoc	Disciplined application of Bayesian methodology
9.	Ability to Address Dynamic as Well as Static Models	Is not convenient for tracking rapid changes in parameters of interest	Is inherently amenable to tracking rapidly changing parameters

system design, demonstrated that the new method is analytically tractable and can indeed produce useful parameter estimates. But, the applicability to multivariate statistical inference using multiple samples has not been recognized or pursued. The theoretical advantages of the alternative method identified in Table 1 provides strong motivation for showing how to apply the new method to multivariate statistics in general.

The Radically Different MoM is no less different from the much newer Generalized MoM introduced by L. P. Hanson in 1982 Wikipedia (2022) than it is from the Classical MoM from a century earlier, and the Generalized MoM is equivalent to several other methods introduced 20-to-30 years earlier Wikipedia (2022).

There is a limitation to the applicability of this new method. As stated in row 2b of Table 1, the user must be able to calculate the moments specified in row 2a of this table, as explicit functions of the unknown parameters. This typically requires a model in which the unknown parameters appear explicitly in the data model. For example, if the unknown parameter is the variance of one random variable for which multiple samples are available, the expected values of any nonlinear functionals to be used for estimation (e.g. the squaring function), conditioned on knowledge of the variance, are not defined; e.g., the fourth moment conditioned on knowledge of the variance is undefined, except in very unique cases like jointly Gaussian variables.

Because there is much distracting detail in the following presentation of the derivation of the new MoM, a streamlined summary of this derivation is provided in the Appendix. Readers may prefer to read the Appendix first in order to know in advance where the derivation is heading as it proceeds through the following sections of this paper.

2 Classical method of moments

Assume we have R observations (samples) $\{x_{k,r} : k = 1, 2, \dots, K, r = 1, 2, \dots, R\}$ of

K random variables $\{X_k\}$ and a model of the functional dependence of these random variables on Q unknown parameters $\boldsymbol{\theta} = \{\theta_q: q = 1, 2, \dots, Q\}$ and L random variables $\mathbf{Z} = \{Z_l: l = 1, 2, \dots, L\}$

$$X_k = f_k(\boldsymbol{\theta}; \mathbf{Z}) \quad (1)$$

We briefly consider three alternative assumptions and then down-select to one:

1. the joint PDF of $\{X_k\}$ is known, or
2. the joint PDF of $\{Z_l\}$ is known and this enables calculation of the joint PDF of $\{X_k\}$,
or
3. a formulaic probabilistic model of \mathbf{X} is available and enables the calculation of the joint moments of $\{X_k\}$

Cases 1) and 2) are quickly dispensed with here because resorting to ad hoc methods in these cases is generally not necessary unless issues of complexity arise. To be more specific, we consider the well-known relationship among prior (before data observation) and posterior (after data observation) probabilities

$$p(\boldsymbol{\theta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{x})} \quad (2)$$

where $p(\boldsymbol{\theta})$ is the prior PDF of the parameters, $p(\boldsymbol{\theta}|\mathbf{x})$ is the posterior PDF, $p(\mathbf{x}|\boldsymbol{\theta})$ is the Likelihood Function and $p(\mathbf{x})$ is the unconditional data PDF, which can be decomposed into conditional PDFs (likelihood functions) as follows:

$$p(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \quad (3)$$

or, for discrete-valued parameters,

$$p(\mathbf{x}) = \sum_{\boldsymbol{\theta}} p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

In the abbreviated notation used here, the particular PDF function is denoted by the

symbol used for its arguments.

Given knowledge of the functions $p(\boldsymbol{\theta})$ and $p(\mathbf{x}|\boldsymbol{\theta})$, the other two functions in (2) can be calculated, and one can choose to use an ML estimate or any Minimum-Bayes'-Risk estimate of the parameter vector $\boldsymbol{\theta}$.

Consequently, resorting to the ad hoc MoM is generally not necessary for parameter estimation unless these functions are not known or are exceedingly difficult to calculate, particularly the Likelihood Function.

For case 3), the classical MoM for estimating the values of $\{\theta_q\}$ is to:

- (1) equate $M \geq Q$ calculated joint probabilistic moments of $\{X_k\}$ to the M corresponding sample moments of $\{x_{k,r} : k = 1, 2, \dots, K; r = 1, 2, \dots, R\}$. For example, some subset M of the $(K^2 + 1)/2$ unique moments from the set of K^2 2nd-order moments can be used:

$$\mathbb{E}\{X_j X_k | \boldsymbol{\theta}\} = \frac{1}{R} \sum_{r=1}^R x_{j,k} x_{k,r} \quad \text{for } j, k = 1, 2, \dots, K \quad (4)$$

Then,

- (2) try to solve this set of simultaneous equations.

3 Radically different method of moments

In preparation for introducing the alternative MoM, we briefly expand the above discussion of Cases 1) and 2). It follows from (2) that any difference between the ML estimate and the MAP estimate is completely determined by the prior PDF. In the event that the prior PDF is uniform over the region where the likelihood function reaches its maximum value, then the ML and MAP estimates are equal. In situations where knowledge of a non-uniform prior PDF is not available, it is common to assume it is uniform over a sufficiently large finite region A in the prior-PDF domain, Q -dimensional Euclidean space:

$$p(\boldsymbol{\theta}) = \begin{cases} \frac{1}{|A|}, & \boldsymbol{\theta} \in A \\ 0, & \boldsymbol{\theta} \notin A \end{cases}$$

where $|A|$ denotes the volume of A . In this case, (2) reduces to

$$p(\boldsymbol{\theta}|\mathbf{x}) = \begin{cases} \frac{p(\mathbf{x}; \boldsymbol{\theta})}{p(\boldsymbol{\theta})}, & \boldsymbol{\theta} \in A \\ 0, & \boldsymbol{\theta} \notin A \end{cases}$$

and (3) reduces to

$$p(\mathbf{x}) = \frac{1}{|A|} \int_A p(\mathbf{x}|\boldsymbol{\theta}) d\boldsymbol{\theta}$$

It follows that the maxima of the likelihood function and the posterior PDF coincide, and these two alternative methods become one and the same:

$$\operatorname{argmax}_{\boldsymbol{\theta} \in A} \{p(\boldsymbol{\theta}|\mathbf{x})\} = \operatorname{argmax}_{\boldsymbol{\theta} \in A} \{p(\mathbf{x}|\boldsymbol{\theta})\}$$

Unfortunately, whether or not the prior PDF is known, if either the likelihood function or the posterior probability is unknown, neither the ML nor Min-Risk methods can be used, and the MoM is likely to be resorted to. However, the interpretation of $\boldsymbol{\theta}$ as a vector of random variables instead of unknown constants enables a radically different alternative to the MoM to be derived. In the original formulation of this method for time series analysis, the name *Structurally Constrained Bayesian Methodology* (SCBM) Gardner (1976*b*) was introduced. This descriptive name is also appropriate in the more general setting of multivariate statistics based on multiple samples; however, it has the disadvantage of suggesting that a full probabilistic model for the data is available, as it must be in classical Bayesian statistics. For this reason, the alternative name *Structurally Constrained Bayesian Method of Moments* is suggested here. The same acronym can be used. The substantive advantages of the SCBM over the classical MoM are described in Table 1 (entries 1, 3-6, 8, 9 in Table 1).

In the SCBM, knowledge of the posterior PDF required by the MAP and other Bayesian methods is replaced with the requirement of knowledge of moments of \mathbf{X} conditioned on $\boldsymbol{\theta}$, as in the classical MoM. Such moments can often, in practice, be calculated from the model (1), even when the posterior PDF and the likelihood function cannot be calculated.

The SCBM specifies (1) a constraint that the estimator be confined to some linear space derived from the observations \mathbf{X} , and (2) a performance criterion for optimizing an estimate of the posterior PDF.

The following discussion explains the extension of the original work on the SCBM to multivariate statistics for which the classical MoM was devised. This discussion is not a necessary part of the Radically Different MoM when applied to the classical MoM data model. But it does lead to an understanding of why the SCBM is able to address dynamic and as well as static MOM problems and produce tracking parameter estimates.

When more than one sample of the vector of observed random variables is available, say R as in (1)), each sample can be interpreted as originating from a distinct K -dimensional vector \mathbf{X}_r , all R of which are identically distributed and are concatenated to form the composite RK -dimensional column-vector of observations $\mathbf{X} = [\mathbf{X}_1^T \mathbf{X}_2^T \dots \mathbf{X}_R^T]$. In this case, each \mathbf{X}_r can contain as few as $K = 1$ random variable, X_r . This enables the SCBM to accommodate scalar-valued time-series of observations — for which the r -th observed random variable X_r is the r -th time sample of a scalar-valued stochastic process — as well as the classical MoM setup involving multiple statistical samples $\{\mathbf{x}_r\}$ of a single vector \mathbf{x}_r of observations, each sample vector $\boldsymbol{\theta}$ depending on the same unknown parameter vector in which time may play no role. The SCBM was originally proposed for time-series analysis for communi-

cations systems, i.e., for statistical signal processing. But, with this simple device of reindexing and re-interpreting, it becomes apparent that the methodology applies as well to the classic multivariate statistics problem for which the MoM was created. The stochastic process model with a single scalar-valued sample path of length KR that is equivalent to the K -variate model with R sample vectors has a special temporal structure because the sequence of RK times samples has a block structure in which the joint PDF of any subset of time samples depends periodically on the time shift parameter with a period of K . That is, the process is Cyclostationary Gardner (2022), Gardner et al. (2006). And this is precisely the type of stochastic process for which the models originally addressed with the SCBM, the sequence of R random Q -vector parameters $\{\Theta_r\}$ is stationary whereas, for the classical MoM model, reinterpreted as a scalar-valued time series model, the Q -vector is fixed from one period to another $\{\theta_r = \theta\}$ and is not treated as a realization of a random vector. Therefore, instead of estimating a time-sequence of Q -vectors, as in the typical communications system application, there is only one Q -vector for all blocks of K time samples.

However, this reveals that the SCBM allows for the MoM data model to be generalized to allow for some evolution of parameter values as more samples are collected. This would be done by allowing the parameter vector to become dependent on the data-block index r so that the otherwise Static SCBM-based MoM becomes Dynamic and tracks evolving parameters. The case in which such changes occur slowly is accommodated by including high correlation in the stationary sequence of random vectors. For the other extreme of maximally rapid changes in the parameter vector, the parameter sequence can be modeled as independent and identically distributed. For this Dynamic MoM

problem, the original time-series formulation Gardner (1973) and Gardner (1976b) is preferred to the classical MoM problem formulation.

3.1 Structural Constraint

The estimator must be some linear functional of some set of specified nonlinear functions $\{g_j(\mathbf{X})\}$ of the random variables \mathbf{X} modeling the observations. A sufficiently general linear functional for many applications has the form

$$\hat{p}(\boldsymbol{\theta}|\mathbf{X}) = \sum_j H_j(\boldsymbol{\theta}) \cdot [g_j(\mathbf{X})] \quad (5)$$

The set of values of $\hat{p}(\boldsymbol{\theta}|\mathbf{X})$ for each specified value of $\boldsymbol{\theta}$ generated by all component linear functionals $\{H_j(\boldsymbol{\theta})\}$ of the specified set of nonlinear functions $\{g_j(\cdot)\}$ of \mathbf{X} is a linear vector space Λ of random variables. For example, one can choose

$$\begin{aligned} g_1(\mathbf{X}) &= \mathbf{X} \\ g_2(\mathbf{X}) &= \mathbf{X}\mathbf{X}^T \end{aligned} \quad (6)$$

in which case (5) reduces to

$$\hat{p}(\boldsymbol{\theta}|\mathbf{X}) = \sum_k h_k(\boldsymbol{\theta})X_k + \sum_{k,l} h_{k,l}(\boldsymbol{\theta})X_kX_l \quad (7)$$

which is a multivariate polynomial of order 2. Of course, this linear-plus-quadratic form is easily generalized to higher order polynomials. In (7), $\{h_k(\boldsymbol{\theta})\}$ is a representation (kernel) of the functional $H_1(\boldsymbol{\theta})$ and $\{h_{k,l}(\boldsymbol{\theta})\}$ is a representation of the functional $H_2(\boldsymbol{\theta})$. In general, the functions $\{g_j(\mathbf{X})\}$ are each tensors of various dimensions, as illustrated in (6), and the functionals $\{H_j(\boldsymbol{\theta})\}$ each map these tensors into scalars.

For applications in which the observed data is a continuous-time stochastic process, $\{X(t) : t \in [a, b] \subset (-\infty, +\infty)\}$, (7) becomes

$$\widehat{p}(\boldsymbol{\theta}|\mathbf{X}) = \int_a^b h(t; \boldsymbol{\theta})X(t)dt + \int_a^b \int_a^b h(t, u; \boldsymbol{\theta})X(t)X(u)dtdu \quad (8)$$

which is a 2nd-order Volterra-like counterpart of a 2nd-order polynomial. In this case, $\{g_j(\mathbf{X})\}$ and $\{H_j(\boldsymbol{\theta})\}$ are continuous counterparts of the discrete tensors and functionals described above.

This choice to constrain the estimator to be in a specified linear vector space facilitates the analytical optimization of the estimator.

In the above example, the solution uses non-centralized moments. A recommended alternative is to replace \mathbf{X} with $\overline{\mathbf{X}} = \mathbf{X} - \mathbb{E}\{\mathbf{X}\}$. Also, the term $g_0(\overline{\mathbf{X}}) = 1$ for which $H_0(\boldsymbol{\theta}) \cdot g_0(\overline{\mathbf{X}}) = h(\boldsymbol{\theta})$, a constant vector, can be included in (6). This adds to the RHS of (7) and (8) the constant term $h_0(\boldsymbol{\theta})$. These modifications are illustrated in Gardner (1976a), and they are made in some of the examples below. In addition, the quantity $p(\boldsymbol{\theta}|\overline{\mathbf{X}})$ can be replaced with $p(\boldsymbol{\theta}|\overline{\mathbf{X}}) - \mathbb{E}\{p(\boldsymbol{\theta}|\overline{\mathbf{X}})\} = p(\boldsymbol{\theta}|\overline{\mathbf{X}}) - p(\boldsymbol{\theta})$, which is done in the examples below. Observe that $p(\boldsymbol{\theta}|\overline{\mathbf{X}}) = p(\boldsymbol{\theta}|\mathbf{X})$.

This completes the explanation of the types of structural constraints imposed by the SCBM method. We now move on to a description of the optimality criterion.

3.2 Optimality Criterion

The performance criterion for optimizing the structurally constrained estimator of the posterior PDF arises from selecting squared error as a cost function. Then the Bayes Risk to be minimized, which is the expected value of the cost, is the Mean Squared Error (MSE):

$$\text{MSE} = \mathbb{E}\{[\widehat{p}(\boldsymbol{\theta}|\mathbf{X}) - p(\boldsymbol{\theta}|\mathbf{X})]^2\} \quad (9)$$

This may seem strange at first glance because probabilities are not random variables. However, when a random variable \mathbf{X} is substituted in place of a sample observation \mathbf{x} , inside the function $p(\boldsymbol{\theta}|\cdot)$, the function value becomes a random variable. For the parameter estimation problems of interest here, we have a set of multiple random variables indexed by the parameter vector $\boldsymbol{\theta}$.

The optimization problem before us is to find the estimator $\hat{p}(\boldsymbol{\theta}|\mathbf{X})$ for each value of $\boldsymbol{\theta}$ that minimizes the above MSE subject to the constraint that the random variable $\hat{p}(\boldsymbol{\theta}|\mathbf{X})$ is contained in the specified linear vector space Λ of all admissible estimates, which we denote by $\tilde{p}(\boldsymbol{\theta}|\mathbf{X})$. The solution to this optimization problem is well-known to be the orthogonal projection of the vector $p(\boldsymbol{\theta}|\mathbf{X})$, generally outside of Λ , onto the hyperplane Λ contained in the linear space of (loosely speaking) *all* functions of the observables.

A technical detail here is that, in order to apply the classical orthogonal projection theorem, the linear space must be an inner-product space, and this in turn requires that the vectors in the space all have finite norms; in this application, this means the probability model of the nonlinearly transformed observed random variables $\{g_j(\mathbf{X})\}$ must have finite-mean-squared values. This puts constraints on both the nonlinearities used, $\{g_j(\cdot)\}$, and the probabilistic model of the original observations $\{p(\mathbf{X}|\boldsymbol{\theta}) : \boldsymbol{\theta} \in A\}$. These constraints are to be expected: one cannot use mean-squared error if random variables of interest do not have finite mean-squared values. Nevertheless, even if \mathbf{X} does not have finite-mean-squared values, the nonlinearities $\{g_j(\cdot)\}$ can be chosen such that $\{g_j(\mathbf{X})\}$ do have finite mean-squared values.

3.3 Solution for Optimum Posterior PDF Estimate

The necessary and sufficient condition that characterizes the orthogonal projection solution described above is the following Orthogonality Condition:

$$\mathbb{E}\{[\hat{p}(\boldsymbol{\theta}|\mathbf{X}) - p(\boldsymbol{\theta}|\mathbf{X})]\tilde{p}(\boldsymbol{\theta}|\mathbf{X})\} = 0 \quad \forall \tilde{p}(\boldsymbol{\theta}|\mathbf{X}) \in \Lambda \quad (10)$$

By using the estimator characterization (5), this condition can be re-expressed as

$$\mathbb{E}\left\{\left[\sum_j H_j(\boldsymbol{\theta}) \cdot [g_j(\mathbf{X})] - p(\boldsymbol{\theta}|\mathbf{X})\right] g_k(\mathbf{X})\right\} = 0 \quad \forall \{k\} \quad (11)$$

which is equivalent to

$$\sum_j H_j(\boldsymbol{\theta}) \cdot \mathbb{E}\{[g_j(\mathbf{X})]g_k(\mathbf{X})\} = \mathbb{E}\{p(\boldsymbol{\theta}|\mathbf{X})g_k(\mathbf{X})\} \quad \forall \{k\} \quad (12)$$

where (because $\mathbb{E}\{\cdot\}$ is linear) $H_j(\boldsymbol{\theta})$ operates on the quantity in square brackets *after* the expectation is executed.

As a final step in simplifying these equations, we use the *magic relationship*:

$$\begin{aligned} \mathbb{E}\{p(\boldsymbol{\theta}|\mathbf{X})g_k(\mathbf{X})\} &= \mathbb{E}\left\{\frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})g_k(\mathbf{X})}{p(\mathbf{X})}\right\} \\ &= \int \frac{p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})g_k(\mathbf{x})}{p(\mathbf{x})} p(\mathbf{x}) d\mathbf{x} \\ &= \int p(\mathbf{x}|\boldsymbol{\theta})g_k(\mathbf{x}) d\mathbf{x} p(\boldsymbol{\theta}) \\ &= \mathbb{E}\{g_k(\mathbf{X}|\boldsymbol{\theta})\} p(\boldsymbol{\theta}) \end{aligned} \quad (13)$$

in which the unknown posterior PDF vanishes and the assumed-known prior PDF (possibly a uniform PDF when it is not known) appears. Substituting (13) into (12) produces

$$\sum_j H_j(\boldsymbol{\theta}) \cdot \mathbb{E}\{[g_j(\mathbf{X})]g_k(\mathbf{X})\} = \mathbb{E}\{g_k(\mathbf{X})|\boldsymbol{\theta}\} p(\boldsymbol{\theta}) \quad \forall \{k\} \quad (14)$$

The unconditional moments in the left member of this set of linear equations can be re-expressed in terms of conditional moments as follows:

$$\sum_j H_j(\boldsymbol{\theta}) \cdot \int \mathbb{E}\{[g_j(\mathbf{X})]g_k(\mathbf{X})|\tilde{\boldsymbol{\theta}}\} p(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{\theta}} = \mathbb{E}\{g_k(\mathbf{X})|\boldsymbol{\theta}\} p(\boldsymbol{\theta}) \quad \forall \{k\} \quad (15)$$

This is a set of linear equations in the unknown linear functionals $\{H_j(\boldsymbol{\theta})\}$. Thus, regardless of the nonlinear functions (tensors) selected in the structural constraint, the

equations to be solved are always linear. In addition, when $\{g_j(\mathbf{X})\}$ are comprised of homogeneous polynomials, as in the examples above, the linear equations are fully specified by moments of \mathbf{X} conditioned on the parameters $\boldsymbol{\theta}$.

This latter observation reveals that the SCBM is a method of moments in the special case for which polynomial nonlinearities $\{g_j(\cdot)\}$ are selected, and the radical difference between the details of the SCBM and those of the classical MoM explains why this method is called a radically different MoM.

Another interesting observation that can be made from (15) is the fact that by using the SCBM, the otherwise required knowledge of the likelihood function—the data PDF conditioned on the parameter values—is replaced with the required knowledge of the 1st and 2nd order moments of prescribed nonlinear functions of the data which, for up-to- n th-order polynomial functions of the data, are 1st through 2 n th order moments of the data. So, the required knowledge of likelihood functions—the data PDFs conditioned on parameter values—is replaced with the required knowledge of a finite set of data moments conditioned on parameter values.

As mentioned in a previous section, when the prior PDF is known, this is additional information the SCBM uses, which the classical MoM does not use. And, in addition, when the prior PDF is not known, it can be assumed to be uniform over a user specified region of parameter space which the user can specify according to any relevant prior information.

To illustrate the design equation whose solution fully specifies the posterior PDF estimate for each set of parameter values $\boldsymbol{\theta}$ of interest, we consider here the example (6), modified by inclusion of the $k = 0$ term and replacement of \mathbf{X} by $\overline{\mathbf{X}}$ as discussed in section 1. Using (12), modified by replacement of $p(\boldsymbol{\theta}|\mathbf{X})$ with $p(\boldsymbol{\theta}|\mathbf{X}) - p(\boldsymbol{\theta})$, we obtain

To illustrate the design equation whose solution fully specifies the posterior PDF estimate for each set of parameter values $\boldsymbol{\theta}$ of interest, we consider here the example (6), modified by inclusion of the $k = 0$ term and replacement of \mathbf{X} by $\overline{\mathbf{X}}$ as discussed in section 1. Using (12), modified by replacement of $p(\boldsymbol{\theta}|\mathbf{X})$ with $p(\boldsymbol{\theta}|\mathbf{X}) - p(\boldsymbol{\theta})$, we obtain

$$\begin{aligned} H_0(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [g_0(\overline{\mathbf{X}})] g_k(\overline{\mathbf{X}}) \} + H_1(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [g_1(\overline{\mathbf{X}})] g_k(\overline{\mathbf{X}}) \} + H_2(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [g_2(\overline{\mathbf{X}})] g_k(\overline{\mathbf{X}}) \} \\ = (\mathbb{E} \{ g_k(\overline{\mathbf{X}}) | \boldsymbol{\theta} \} - \mathbb{E} \{ g_k(\overline{\mathbf{X}}) \}) p(\boldsymbol{\theta}) \text{ for } k = 0, 1, 2 \end{aligned} \quad (16)$$

which, using modified (6), is equivalent to

$$\begin{aligned} H_0(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [1] g_k(\overline{\mathbf{X}}) \} + H_1(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}}] g_k(\overline{\mathbf{X}}) \} + H_2(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}} \overline{\mathbf{X}}^T] g_k(\overline{\mathbf{X}}) \} \\ = (\mathbb{E} \{ g_k(\overline{\mathbf{X}}) | \boldsymbol{\theta} \} - \mathbb{E} \{ g_k(\overline{\mathbf{X}}) \}) p(\boldsymbol{\theta}) \text{ for } k = 0, 1, 2 \end{aligned} \quad (17)$$

which can be more explicitly expressed as

$$\begin{aligned} H_0(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [1] \} + H_1(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}}] \} + H_2(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}} \overline{\mathbf{X}}^T] \} &= (\mathbb{E} \{ 1 | \boldsymbol{\theta} \} - \mathbb{E} \{ 1 \}) p(\boldsymbol{\theta}) \\ H_0(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [1] \overline{\mathbf{X}} \} + H_1(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}}] \overline{\mathbf{X}} \} + H_2(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}} \overline{\mathbf{X}}^T] \overline{\mathbf{X}} \} \\ &= (\mathbb{E} \{ \overline{\mathbf{X}} | \boldsymbol{\theta} \} - \mathbb{E} \{ \overline{\mathbf{X}} \}) p(\boldsymbol{\theta}) \\ H_0(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [1] \overline{\mathbf{X}} \overline{\mathbf{X}}^T \} + H_1(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}}] \overline{\mathbf{X}} \overline{\mathbf{X}}^T \} + H_2(\boldsymbol{\theta}) \cdot \mathbb{E} \{ [\overline{\mathbf{X}} \overline{\mathbf{X}}^T] \overline{\mathbf{X}} \overline{\mathbf{X}}^T \} \\ &= (\mathbb{E} \{ \overline{\mathbf{X}} \overline{\mathbf{X}}^T | \boldsymbol{\theta} \} - \mathbb{E} \{ \overline{\mathbf{X}} \overline{\mathbf{X}}^T \}) p(\boldsymbol{\theta}) \end{aligned} \quad (18)$$

Using (7), we can now re-express the above set of linear equations more explicitly as follows:

$$\begin{aligned} h(\boldsymbol{\theta}) + \sum_{k,l} h_{k,l}(\boldsymbol{\theta}) \mathbb{E} \{ \overline{X}_k \overline{X}_l \} &= 0 \\ \sum_k h_k(\boldsymbol{\theta}) \mathbb{E} \{ \overline{X}_k \overline{X}_j \} + \sum_{k,l} h_{k,l}(\boldsymbol{\theta}) \mathbb{E} \{ \overline{X}_k \overline{X}_l \overline{X}_j \} &= \mathbb{E} \{ \overline{X}_j | \boldsymbol{\theta} \} p(\boldsymbol{\theta}) \quad \forall j \\ h(\boldsymbol{\theta}) \mathbb{E} \{ \overline{X}_j \overline{X}_i \} + \sum_k h_k(\boldsymbol{\theta}) \mathbb{E} \{ \overline{X}_k \overline{X}_j \overline{X}_i \} + \sum_{k,l} h_{k,l}(\boldsymbol{\theta}) \mathbb{E} \{ \overline{X}_k \overline{X}_l \overline{X}_j \overline{X}_i \} \\ &= (\mathbb{E} \{ \overline{X}_j \overline{X}_i | \boldsymbol{\theta} \} - \mathbb{E} \{ \overline{X}_j \overline{X}_i \}) p(\boldsymbol{\theta}) \quad \forall i, j \end{aligned} \quad (19)$$

The unconditional moments in (19) can be characterized in terms of conditional moments using (3) as follows for example:

$$\mathbb{E}\{X_k X_l X_j X_i\} = \int \mathbb{E}\{X_k X_l X_j X_i | \boldsymbol{\theta}\} p(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (20)$$

If no prior information is available for specifying a prior PDF, a uniform PDF can be used to obtain

$$\int \mathbb{E}\{X_k X_l X_j X_i | \boldsymbol{\theta}\} p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{|A|} \int_A \mathbb{E}\{X_k X_l X_j X_i | \boldsymbol{\theta}\} d\boldsymbol{\theta} \quad (21)$$

The solutions to (19) are used in the estimator formula (7), modified by replacement of $\hat{p}(\boldsymbol{\theta} | \mathbf{X})$ with $\hat{p}(\boldsymbol{\theta} | \mathbf{X}) - p(\boldsymbol{\theta})$.

Example 1: Linear Estimator In the case of a constant-plus-linearly-constrained estimator of the posterior PDF, the design equation (18) reduces to

$$\begin{aligned} h(\boldsymbol{\theta}) + \sum_k h_k(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k\} &= 0 \\ h(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_j\} + \sum_k h_k(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k \bar{X}_j\} &= \mathbb{E}\{\bar{X}_j | \boldsymbol{\theta}\} p(\boldsymbol{\theta}) \quad \forall j \end{aligned} \quad (22)$$

which has the explicit solution

$$\begin{aligned} h(\boldsymbol{\theta}) &= 0 \\ \mathbf{h}(\boldsymbol{\theta}) &= \left[\mathbb{E}\{\bar{\mathbf{X}} \bar{\mathbf{X}}^T\} \right]^{-1} \mathbb{E}\{\bar{\mathbf{X}} | \boldsymbol{\theta}\} p(\boldsymbol{\theta}) \end{aligned} \quad (23)$$

Consequently, the posterior PDF estimator, given by the modified version of (7), reduces to

$$\hat{p}(\boldsymbol{\theta} | \bar{\mathbf{X}}) = (1 + \mathbf{h}^T(\boldsymbol{\theta}) \bar{\mathbf{X}}) p(\boldsymbol{\theta}) \quad (24)$$

Substituting (23) into (24) yields

$$\hat{p}(\boldsymbol{\theta} | \mathbf{X}) = p(\boldsymbol{\theta}) \left(1 + \mathbb{E}\{\bar{\mathbf{X}}^T | \boldsymbol{\theta}\} \left[\mathbb{E}\{\bar{\mathbf{X}} \bar{\mathbf{X}}^T\} \right]^{-1} \bar{\mathbf{X}} \right) \quad (25)$$

Denoting the square root of the inverse of the covariance matrix in (25) by \mathbf{W} , and denoting the decorrelated vector of observations by $\mathbf{Y} = \mathbf{W}\overline{\mathbf{X}}$, we can re-express (25) for a particular sample of data \mathbf{x} as

$$\widehat{p}(\boldsymbol{\theta}|\mathbf{x}) = p(\boldsymbol{\theta}) (1 + \mathbb{E}\{\mathbf{Y}^T|\boldsymbol{\theta}\}\mathbf{y}) \quad (26)$$

In words, the constant plus linear estimator probabilistically centers the data and probabilistically decorrelates it and then empirically correlates it with its probabilistic mean conditioned on the parameter vector.

In the special case for which the data consists of a known function of unknown parameters (call it a signal) in additive self-correlated zero-mean noise, $\mathbf{X} = \mathbf{s}(\boldsymbol{\Theta}) + \mathbf{N}$, with noise covariance \mathbf{R}_N , that is uncorrelated with the signal whose covariance is $\mathbf{R}_{s(\boldsymbol{\Theta})}$, we obtain $\mathbf{W} = [\mathbf{R}_N + \mathbf{R}_{s(\boldsymbol{\Theta})}]^{-1/2}$ and $\overline{\mathbf{X}} = \mathbf{X} - \mathbb{E}\{\mathbf{s}(\boldsymbol{\Theta})\}$. One cannot say much more about the posterior PDF estimate (26) without focusing on a particular function of the parameters $\mathbf{s}(\boldsymbol{\theta})$. For example, for the special case of a sinusoidal signal with unknown phase, θ , the mode of the estimated posterior PDF can be analytically shown to converge to the true value of the phase as the SNR increases, for any fixed data set (Gardner 1976b, p. 590).

For a pseudo-MAP estimator of $\boldsymbol{\theta}$, (25) yields

$$\widehat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \left\{ p(\boldsymbol{\theta}) \left(1 + \mathbb{E}\{\overline{\mathbf{X}}^T|\boldsymbol{\theta}\} \left[\mathbb{E}\{\overline{\mathbf{X}}\overline{\mathbf{X}}^T\} \right]^{-1} \overline{\mathbf{x}} \right) \right\} \quad (27)$$

which can be re-expressed using (26) as

$$\widehat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \left\{ p(\boldsymbol{\theta}) (1 + \mathbb{E}\{\mathbf{Y}^T|\boldsymbol{\theta}\}\mathbf{y}) \right\} \quad (28)$$

Similarly, the pseudo-MMSE estimator, which is the pseudo-posterior mean, is given by

$$\widehat{\boldsymbol{\theta}}_{\text{MMSE}} = \int p(\boldsymbol{\theta}) (1 + \mathbb{E}\{\mathbf{Y}^T|\boldsymbol{\theta}\}\mathbf{y}) \boldsymbol{\theta} d\boldsymbol{\theta} \quad (29)$$

where $d\boldsymbol{\theta} = d\theta_1 d\theta_2 \dots d\theta_Q$.

One way to investigate the utility of any of these estimators is to seek to determine the conditions under which the parameter estimate equals the true value of the parameter when the additive noise in the data is zero. In the case of zero noise, $\mathbf{N} = \mathbf{0}$ and therefore $\mathbf{R}_N = \mathbf{0}$. Assuming also that the signal covariance matrix has full rank, the posterior PDF estimate (26) reduces to

$$\hat{p}(\boldsymbol{\theta}|\mathbf{x}) = p(\boldsymbol{\theta}) \left(1 + \mathbb{E}\{\mathbf{Y}^T|\boldsymbol{\theta}\}\mathbf{y}\right) = p(\boldsymbol{\theta}) \left(1 + [\mathbf{s}(\boldsymbol{\theta}) - \boldsymbol{\mu}_s]^T \mathbf{R}_s^{-1} [\mathbf{s}(\boldsymbol{\theta}_o) - \boldsymbol{\mu}_s]\right) \quad (30)$$

One example, for which the behavior of this zero-noise PDF estimate is transparent is that for which the norm of the deviation $\mathbf{s}(\boldsymbol{\theta}) - \boldsymbol{\mu}_s$ of the signal from its mean is independent of the value of the parameter vector $\boldsymbol{\theta}$. In this case, the maximum of the factor multiplying the prior $p(\boldsymbol{\theta})$ occurs at the true value $\boldsymbol{\theta} = \boldsymbol{\theta}_o$. Nevertheless, if the prior PDF is non-uniform, it is possible for that factor in (30) to shift the peak in $\boldsymbol{\theta}$ away from the true value. Also, if the norm is strongly dependent on the parameter value, the peak might be shifted away from the true value and toward the value for which the norm is relatively large. Previous studies have shown that the pseudo posterior mean can outperform the pseudo posterior mode in some such cases.

Example 2: Linear Plus Quadratic Estimator In the case of a constant-plus-linearly-plus-quadratically constrained estimator of the posterior PDF, for the special case in which the odd-order unconditional moments of the observed data are zero, the design equations (19) reduce to the following equations:

$$\begin{aligned} h(\boldsymbol{\theta}) + \sum_{k,l} h_{k,l}(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k \bar{X}_l\} &= 0 \\ \sum_k h_k(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k \bar{X}_j\} &= \mathbb{E}\{\bar{X}_j|\boldsymbol{\theta}\} p(\boldsymbol{\theta}) \quad \forall j \\ h(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_j \bar{X}_i\} + \sum_{k,l} h_{k,l}(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k \bar{X}_l \bar{X}_j \bar{X}_i\} &= (\mathbb{E}\{\bar{X}_j \bar{X}_i|\boldsymbol{\theta}\} - \mathbb{E}\{\bar{X}_j \bar{X}_i\}) p(\boldsymbol{\theta}) \quad \forall i, j \end{aligned} \quad (31)$$

which can be solved to obtain

$$\begin{aligned}
h(\boldsymbol{\theta}) &= -\sum_{k,l} h_{k,l}(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k \bar{X}_l\} \\
h_k(\boldsymbol{\theta}) &= \sum_j [\mathbb{E}\{\bar{X}_k \bar{X}_j\}]^{-1} \mathbb{E}\{\bar{X}_j | \boldsymbol{\theta}\} p(\boldsymbol{\theta}) \quad \forall k \\
h_{k,l}(\boldsymbol{\theta}) &= \sum_{j,i} [\mathbb{E}\{\bar{X}_k \bar{X}_l \bar{X}_j \bar{X}_i\}]^{-1} \cdot [(\mathbb{E}\{\bar{X}_j \bar{X}_i | \boldsymbol{\theta}\} - \mathbb{E}\{\bar{X}_j \bar{X}_i\}) p(\boldsymbol{\theta}) - h(\boldsymbol{\theta}) \mathbb{E}\{\bar{X}_k \bar{X}_l\}] \quad \forall k, l
\end{aligned} \tag{32}$$

Finally, the third equation in (32) can be substituted into the first equation in (32) and the scalar $h(\boldsymbol{\theta})$ can be solved for and substituted back into the third equation to obtain the desired 3 explicit solutions for the unknown scalar, vector, and matrix defining the estimator. As can be seen, the third equation above requires the inversion of a rank-4 tensor. A standard approach to doing this is to represent the tensor in terms of matrices and use existing software to invert the matrices, and then convert those back to the desired inverse tensor. See, for example, the article Bu et al. (2014), and references therein, and Kisil et al. (2022).

Example 3: Higher-Order Polynomial Estimators Observe from (19) that the representations of the linear functionals $\{H_j(\boldsymbol{\theta})\}$ for homogeneous polynomial nonlinearities $\{g_j(\mathbf{X})\}$ are rank-1 tensors (vectors) in one linear design equation for 1st order polynomials, then rank-1 and rank-2 tensors (vectors and matrices) in two simultaneous linear design equations for 2nd order polynomials, then rank-1, rank-2, and (by extrapolating) rank-3 tensors in three linear design equations for 3rd order polynomials, etc; and the conditional moments of the modeled data defining these linear equations are rank-1 and rank-2 tensors for 1st-order polynomials, then rank-1 through rank-4 tensors for 2nd order polynomials, and then rank-1 through rank-6 tensors for 3rd order polynomials, etc.

This pattern enables one to simply write down the tensor design equations for any order polynomial estimator of the posterior PDF. All the analytical work has been done here, leaving for the user only the computational challenge of inverting tensors or otherwise solving explicit linear tensor equations.

It is worthy of note that the quadratic-plus-linear-plus-constant estimator of the posterior PDF obtained from the SCBM, specified by (19), is consistent with the linear-plus-constant estimator obtained from the SCBM, specified by (22), in the sense that the latter is included as a special case of the more general former. That is, the solution for the constant and linear parts of the former obtained by equating to zero the quadratic kernel and eliminating the highest-order (3rd) equation is identical to the solution for the constant-plus-linear estimator from the latter.

In conclusion, there is one unique SCBM solution formula for the structurally constrained posterior PDF estimate comprised of a linear combination of any set of finite-mean square non-linear functions of the data.

Moreover, for each natural number n , there is one single algorithm that solves all SCBM problems for polynomial-type nonlinearities of any order $m = 1$, or 2, or 3, or \dots , or n up to the point of specifying the particular algorithms to be used to invert the tensors of various ranks each corresponding to one of the various homogeneous polynomial terms used.

In this conclusion above, the various functional kernels need not be of different dimensions. For example, if the nonlinearities of interest produce the data $\{X_j\}$ and $\{X_j^2\}$ to be linearly combined, the two corresponding kernels are both of dimension 1. The only requirement is that the set of nonlinearly transformed data sets be linearly independent, which guarantees that the specified tensor inverses exist. Other simplifications of the general

solution can be obtained by orthogonalizing the nonlinearly transformed data sets prior to forming their linear combinations. This does not change the linear space Λ , but it does affect the form of the solution. In particular, it renders the off-diagonal ($j \neq k$) terms, in the summary in Section 4 below, zero.

4 Summing up the Radically Different MoM

The SCBM can be summed up as follows:

- The Pseudo Min-Risk Estimate of a parameter vector is calculated from the structurally constrained Min-MSE estimate of the random posterior PDF in the same manner that the true Min-Risk parameter estimate would be computed from the true posterior PDF, were it available.
- The Min-MSE Posterior PDF Estimate is calculated from the structurally constrained formula

$$\hat{p}(\boldsymbol{\theta}|\mathbf{x}) = \sum_j H_j [g_j(\mathbf{x})]$$

in which the nonlinear functions (tensors) $\{g_j(\mathbf{x})\}$ are specified by the user (e.g., (7) or (8)).

- The linear functionals $\{H_j\}$ in this formula are the solutions to the set of simultaneous linear equations

$$\sum_j H_j \cdot \int \mathbb{E} \left\{ [g_j(\mathbf{X})] g_k(\mathbf{X}) | \tilde{\boldsymbol{\theta}} \right\} p(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{\theta}} = \mathbb{E} \{ g_k(\mathbf{X}) | \boldsymbol{\theta} \} p(\boldsymbol{\theta}) \quad \forall \{k\}$$

(e.g., (19)). If the prior PDF is unknown, it is approximated with a uniform PDF over a user specified region of the parameter space (e.g., (21)).

- If the user specified nonlinear functions are multivariate polynomials, then all expected values in these linear equations are conditional moments obtained from a probabilistic model of the observations, justifying this as a *method of moments* (e.g., (19) - (21)).
- Moreover, for homogeneous polynomial nonlinearities, the linear design equations can be explicitly written down in terms of linear tensor equations, knowing nothing more than the specified order of the polynomial to be used. Similarly, the estimator formula can be explicitly written down as a polynomial in the observed data. The only work a user needs to do is solve the known simultaneous linear tensor equations and implement the polynomial posterior PDF estimator.
- As explained below in Section 7, the estimated posterior PDFs satisfy all but one of the 3 traditional axioms of probability

5 Options for SCBM Solutions for Parameter Estimates

Once we have the optimum estimate of the posterior PDF, we can proceed to choose a particular Bayesian Minimum-Risk performance criterion for estimating the parameters θ . For example, we can choose the posterior mode (MAP) criterion described above which, for the assumption of uniform prior PDF, is equivalent to ML; or we can choose the posterior median, which derives from using the absolute value of the error in each element of the estimate of the vector θ for the risk function. We also can use the posterior mean, which results from using the squared error of each element of the estimate of the vector θ . Some comparisons have been made between the pseudo posterior mode and pseudo posterior

mean estimates in Gardner (1973), Gardner (1976*b*), and especially Gardner (1981). The results of these comparisons depend on the particular structural constraints chosen. Consequently, there may be low likelihood of obtaining any general comparative results on performance dependence on the selected type of risk. Nevertheless, the results in Gardner (1981) establish some conditions under which the estimated posterior mean is superior to the estimated posterior mode for the decision problem of classifying observed data into one of a finite number of specified classes. This is interesting since the mode seems like a more natural choice and actually is when the posterior probability is not just an estimate.

6 Application of SCBM to Decision Making

The Bayesian approach to minimum-risk decision making uses the same performance criterion as that it uses for parameter estimation. The primary difference is that the parameters for decision making are discrete-valued, and each discrete value corresponds to a particular hypothesis. The hypothesis that is decided to be the correct one minimizes the risk, given the particular observed data. Consequently, the SCBM described in this paper applies as well to decision making as it does to parameter estimation. This has been pursued in the early work reported in Gardner (1973), Gardner (1976*b*), Gardner (1981). The Author does not know of any formalism that has been formulated for a decision-making counterpart to the classical MoM formulated for parameter estimation. (However, one would expect that some work on this concept has been done.) Consequently, no complement to Table 1 that applies to decision making is included herein. Nevertheless, it seems likely that Table 1 applies, as is, to both parameter estimation and decision making.

7 Properties of the SCBM Posterior PDF Estimator

It is shown in the original contribution Gardner (1976*b*) that the posterior PDF (and discrete probability mass function) estimates provided by the SCBM satisfy the traditional axioms of probability, regardless of the specific structural constraints chosen by the user, except for the positivity axiom. Another property of interest is revealed by the general solution (26) for a constant-plus-linear constraint, and this is that the posterior PDF estimate is explicitly specified in terms of the prior PDF and the conditional mean of the centered and decorrelated data. In all cases of essentially arbitrary nonlinearities in the structural constraints, the solution is fully specified in terms of the prior PDF and conditional first- and second-order moments of the nonlinearly transformed data. And for polynomial nonlinearities, these are equivalent to higher-order conditional moments of the model for the original random data, guaranteeing this is indeed a method of moments; however, in place of the sample moments of the data used in the Classical MoM, more general weighted averages of the data and products of the data with itself are used, and the weighting functions are optimized according to a Bayesian minimum-risk criterion.

8 Applications

To illustrate a nontraditional type of application of this alternative MoM, previously published work is referred to here. In Gardner (1973), Gardner (1976*b*) the problem of optimizing a digital communications system receiver is addressed. One of the models used for this is a continuous-time cyclostationary process defined for all time, and the unknown parameters in this process comprise an infinite sequence of discrete values from a finite alphabet of encoded symbols representing the information-bearing data being transmitted on a stream of pulses. Thus, this is an ongoing decision problem in which a decision as to

which symbol was transmitted is made every symbol interval (after some delay required to process data following each symbol interval) . The data received for each symbol extends over multiple symbol intervals, creating what is called inter-symbol interference. As shown in Gardner (1973), Gardner (1976*b*), the solution for a constant-plus-linearly-constrained receiver has much in common with the min-risk receiver for additive Gaussian noise: It is comprised of a parallel bank of matched filters, each filter matched to one of the finite set of transmitted pulse shapes, followed by a symbol-rate time sampler and a multi-input/multi-output sampled-data filter which produces SCBM estimates of the posterior probabilities of the transmitted symbols. This portion of the receiver structure that follows the bank of matched filters is known as a Fractionally Spaced Equalizer, which attempts to remove the intersymbol interference; however, its function is seen here to be much more than a traditional channel equalizer. In fact, it is more akin to a discrete-time Wiener filter. These probability estimates can be used for making decisions on which of the symbols from the finite alphabet were transmitted or for estimating symbol values or estimating the entire transmitted signal.

Another application, addressed in Gardner (1976*a*), considers parameter estimation and decision making for marked and filtered Poisson processes, used to model optical communications signals transmitted over optical fibers. Results obtained for a linearly constrained receiver strongly paralleling those obtained in Gardner (1973), Gardner (1976*b*).

Yet another application to communications receiver design is addressed in Gardner (1976*b*), where a linear-plus-quadratically constrained receiver for noncoherent decision making for sinewave-carrier modulated signals is considered. Again, results obtained are similar to optimum receivers for signals in Gaussian noise.

9 Reflection

Some of the concepts used to formulate the SCBM parameter estimation method could be said to be twisted—they are quite unconventional. Seeking a new MoM within the Bayesian framework seems unmotivated and, at first glance, unlikely to succeed. Yet the Bayesian formulation is logical, and it leads to a tractable genuine MoM for two reasons:

1. The infrequently used concept that the posterior probability, with the conditioning quantity — which is normally a sample of a set of observed random variables — replaced with the observable random variables (not their samples), is itself a random variable and can be subjected to classical random variable estimation theory; though, it is uncommon to apply such theory to the problem of estimating an unknown deterministic function $u(\mathbf{X})$ of the observations, which is exactly what the posterior probability is. In fact, such a problem is generally unsolvable because it generally requires knowledge of the unknown function, even when the estimates are constrained to belong to a linear space derived from the observations, such as Λ herein. It appears, at first glance, by comparing (12) and (13), to be solvable under only one condition and this is that $u(\mathbf{x})$ is proportional to the ratio $p(\mathbf{x}|\boldsymbol{\theta})/p(\mathbf{x})$ of the likelihood function to the unconditional PDF of the data, an example of which is the posterior PDF in which case the proportionality factor is the prior PDF $p(\boldsymbol{\theta})$. This condition is responsible for the disappearance of the unknown function $u(\mathbf{X}) = p(\boldsymbol{\theta}|\mathbf{X})$ in the RHS of the design equation (12) as per (13). However, a deeper look reveals that $u(\mathbf{X})$ and $p(\boldsymbol{\theta}|\mathbf{X})a$ for any scalar a can differ by any random variable that is orthogonal to $g_k(\mathbf{X})$ for all k . A good example is $u(\mathbf{X})$ equal to the event indicator function, $u(\mathbf{X}) = 1$ for all samples $\mathbf{X} = \mathbf{x}$ for which the event $\Theta = \boldsymbol{\theta}$ occurs and $u(\mathbf{X}) = 0$ for all other $\mathbf{X} = \mathbf{x}$. It is easily shown that (12) reduces to (14) with this choice for

$u(\mathbf{X})$. The reason for this is that $p(\boldsymbol{\theta}|\mathbf{X})$ is the orthogonal projection of this indicator function onto the space of all finite mean-square functions of \mathbf{X} (see (Gardner 1989, pp. 427-428)). Therefore, the orthogonal projection of this indicator function onto the linear sub-space Λ is identical to the orthogonal projection of $p(\boldsymbol{\theta}|\mathbf{X})$ onto Λ .

2. The adoption of minimum-mean-squared error as an optimality criterion for estimating the function $u(\mathbf{X})$, together with the constraint on the estimator to a hyperplane in the space of all functions $g(\mathbf{X})$ of the data. These two choices of formulation are responsible for the design equation (12) being a set of linear equations.

The observation above reveals that this alternative MoM could have been formulated in terms of estimating either any scaled version of the event indicator function or the ratio $p(\mathbf{x}|\boldsymbol{\theta})/p(\mathbf{x})$ instead of the posterior PDF $p(\boldsymbol{\theta}|\mathbf{x})$. In these cases, the prior PDF $p(\boldsymbol{\theta})$ disappears (with the appropriate scalar a) from the RHS of the general design equation (15), but not the LHS.

Because this methodology is so highly structured in terms of the algorithms required for implementation, namely linear equation solvers and multivariate polynomial functionals of the observations, it should be highly amenable to efficient algorithmic implementations in terms of either software computer applications or special purpose digital signal processing hardware.

As a final remark, it is mentioned that, unlike the Radically Different MoM, the Classical MoM does not appear to be nearly as convenient a starting point for developing a tracking parameter estimator, regardless of how the memory of the sample moments calculator is adjusted, because every change in the sample moments requires the solution of a new set of generally nonlinear equations.

APPENDIX: Outline of Derivation of New MoM

Background

- The Method of Moments (MoM) is a classical statistical technique for estimating the parameters of a probabilistic data model
- The MoM was introduced just prior to the turn of the 19th Century by K. Pearson and P. Chebyshev, independently
- It is designed for statistical inference where the available data consists of multiple samples of a set of random variables, with a partially specified probabilistic model
- **The partial model needed is a set of joint moments of various orders for the random variables, showing explicit dependence on unknown parameters**

The Classical Method of Moments

- The number of moments M needed is equal to the number of unknown parameters a in these moment models (formulas); e.g.,

$$M_{12} = \mathbb{E}\{X_1 X_2\} = f(a_1, a_2, a_3) \quad a_1, a_2 = \text{variances}, \quad a_3 = \text{covariance}$$

$$M_2 = \mathbb{E}\{(X_2)^2\} = g(a_2)$$

$$M_1 = \mathbb{E}\{(X_1)^2\} = h(a_1)$$

for which f, g, h are known functions

- The statistics that are computed from the data consist of the sample moments corresponding to the theoretical moment models, e.g.,

$$m_{12} = \frac{1}{n} \sum_{j=1}^n x_1^j x_2^j$$

$$m_2 = \frac{1}{n} \sum_{j=1}^n (x_2^j)^2$$

$$m_1 = \frac{1}{n} \sum_{j=1}^n (x_1^j)^2$$

- The inference procedure is to equate the computed sample moments to the theoretical moment formulas and attempt to solve these equations

$$m_{12} = f(a_1, a_2, a_3)$$

$$m_2 = g(a_2)$$

$$m_1 = h(a_1)$$

- The tractability of this MoM depends on the particular nonlinear equations

An Alternative Approach

- I recently observed that every multivariate statistical inference problem based on multiple samples can be *reformulated* as a problem of statistical inference for a single times series of data based on one sample path of the time series, consisting of concatenated time-series segments equal to a first sample of the ordered set of random variables, followed by a 2nd sample of the same random variables, and so on until all samples have been included, e.g.,

$$\{y_k\}_1^{16} = \{x_1^1, x_2^1, x_1^2, x_2^2, x_1^3, x_2^3, x_1^4, x_2^4, x_1^5, x_2^5, x_1^6, x_2^6, x_1^7, x_2^7, x_1^8, x_2^8\}$$

- The theoretical model for this time series is a single sample path of a cyclostationary stochastic process $\{Y_k\}$, with period equal to the number of random variables and with the time sequence of this set of random variables being i.i.d. from one period to the next: e.g., $\{x_1^1, x_2^1\}$ and $\{x_1^2, x_2^2\}$ are i.i.d.
- This is a special cyclostationary process because it contains the same unknown parameters in every period
- I generalized this model to allow the parameter values to change from one period to the next and modeled them as samples of a stationary sequence of random variables, which preserves the cyclostationarity

- Then I invoked an unusual methodology I had introduced in the early 1970s for this type of cyclostationary process model which I used for commonly encountered digital pulse-modulated signals used in communications transmission systems
- The unusual methodology uses Bayesian concepts to formulate the problem of estimating the parameter values (transmitted digits $\{a_i\}$) in terms of the sequence of posterior probabilities, which can be used to compute various minimum-risk parameter estimates, such as maximum-posterior-probability estimates and minimum-mean-squared-error estimates, e.g.,

$$\hat{a}_i = \max_{a_i} P(a_i|\{y_k\})$$

- Finally, I formulated an inference problem for estimating these posterior probabilities using structurally constrained minimum-MSE estimators: optimum linear combinations of any appropriate specified nonlinear transformation of the data samples

$$\hat{\hat{a}}_i = \max_{a_i} \hat{P}(a_i|\{y_k\})$$

- This particular formulation ensures the posterior probability estimates are always the solutions to sets of *simultaneous linear equations*
- By choosing polynomial nonlinearities, the equations are fully *specified by weighted sample moments* of the data; **this makes it a MoM**
- The weights are optimal in the sense of producing structurally constrained minimum-MSE estimates of the posterior probabilities
- In actuality, the reformulation process described above was performed in reverse order for the purpose of showing that **the original work on time series was equivalent to a radically new MoM.**

Summary

- We now have two radically different *Methods of Moments*
- The **numerous advantages** of the new method are fully described in the Table 1 in Section 1
- The utility of the new method was studied back in the 1970s for estimating digital symbols in digital transmission systems developed by Bell Telephone Labs
- **But more diverse applications to various specific multivariate parameter estimation problems, and comparison with the classical MoM, has not yet been pursued**

What's Unusual About this Application of Bayes Minimum Risk Methodology?

- The quantities to be estimated, the posterior probabilities of parameters, are **deterministic functions of the observed data**.
- So, why do we need to estimate them?
- For the same reason we would choose to use the MoM: we do not know the complete probabilistic model for the data
- The particular way I set up the problem for estimating the unknown function

$$P(a|\{y_k\})$$

of the known data **requires knowledge of only moments of orders determined by the orders of the polynomial nonlinearities selected for the structural constraint**

- This was not foreseen, but rather was discovered during my open-ended investigation as a young naïve investigator in my first year as an assistant professor

References

- Bu, C., Zhang, X., Zhou, J., Wang, W. & Wei, Y. (2014), ‘The inverse, rank and product of tensors’, *Linear Algebra and Its Applications* **446**, 269–280.
- Gardner, W. A. (1973), ‘The structure of least-mean-square linear estimators for synchronous M-ary signals (corresp.)’, *IEEE Transactions on Information Theory* **19**(2), 240–243.
- Gardner, W. A. (1976a), ‘An equivalent linear model for marked and filtered doubly stochastic Poisson processes with application to MMSE linear estimation for synchronous m-ary optical data signals’, *IEEE Transactions on Communications* **24**(8), 917–921.
- Gardner, W. A. (1976b), ‘Structurally constrained receivers for signal detection and estimation’, *IEEE Transactions on Communications* **24**(6), 578–592.
- Gardner, W. A. (1981), ‘Design of nearest prototype signal classifiers (corresp.)’, *IEEE Transactions on Information Theory* **27**(3), 368–372.
- Gardner, W. A. (1989), *Introduction to random processes with applications to signals & systems*, McGraw-Hill.
- Gardner, W. A. (2022), ‘Cyclostationarity: educational website’, www.cyclostationarity.com .

Gardner, W. A., Napolitano, A. & Paura, L. (2006), ‘Cyclostationarity: Half a century of research’, *Signal processing* **86**(4), 639–697.

Kisil, I., Calvi, G. G., Konstantinidis, K., Xu, Y. L. & Mandic, D. P. (2022), ‘Accelerating tensor contraction products via tensor-train decomposition [tips & tricks]’, *IEEE Signal Processing Magazine* **39**(5), 63–70.

Pearson, K. (1936), ‘Method of moments and method of maximum likelihood’, *Biometrika* **28**(1/2), 34–59.

Wikipedia (2022), ‘Generalized method of moments’,
https://en.wikipedia.org/wiki/Generalized_method_of_moments .