

Optimum Estimation via Gradients of Partition Functions and Information Measures: A Statistical–Mechanical Perspective

Neri Merhav

Department of Electrical Engineering
Technion – Israel Institute of Technology
Technion City, Haifa 32000, Israel
`merhav@ee.technion.ac.il`

Abstract

In continuation to a recent work on the statistical–mechanical analysis of minimum mean square error (MMSE) estimation in Gaussian noise via its relation to the mutual information (the I–MMSE relation), here we propose a simple and more direct relationship between optimum estimation and certain information measures (e.g., the information density and the Fisher information), which can be viewed as partition functions and hence are amenable to analysis using statistical–mechanical techniques. The proposed approach has several advantages, most notably, its applicability to general sources and channels, as opposed to the I–MMSE relation and its variants which hold only for certain classes of channels (e.g., additive white Gaussian noise channels). We then demonstrate the derivation of the conditional mean estimator and the MMSE in a few examples. Two of these examples turn out to be generalizable to a fairly wide class of sources and channels. For this class, the proposed approach is shown to yield an approximate conditional mean estimator and an MMSE formula that has the flavor of a single–letter expression. We also show how our approach can easily be generalized to situations of mismatched estimation.

Index Terms: Conditional mean estimation, minimum mean squared error, partition function, statistical mechanics, Fisher information.

1 Introduction

Relationships between signal estimation, signal detection, and information measures, both in discrete time and continuous time, have been known for decades [1],[3],[8] and have gained a remarkable degree of revived interest and research activity in the last several years, see, e.g., [4], [5], [6], [7], [11], [12], [13] and references therein.

In particular, in [5], Guo, Shamai and Verdú have derived a relation between the mutual information between the input and the output of an additive white Gaussian noise (AWGN) channel and the minimum mean squared error (MMSE) of non-causal estimation of the channel input based on its output. In particular, this relation, which is often called the I-MMSE relation, shows that the derivative of the mutual information with respect to (w.r.t.) the signal-to-noise (SNR) is equal to half of the MMSE, and it is intimately related to the de Bruijn identity [2, Sec. 17.7]. Later, this relation has been generalized and further developed in several directions: Guo, Shamai, and Verdú [6] and Raginsky and Coleman [12] have derived relations of the same spirit for more general additive channels. Palomar and Verdú [11] have studied relations between the covariance matrix of the MMSE estimator and arbitrary gradients of the mutual information for a general vector Gaussian channel, which allows also a linear transformation of the input signal. In [7], relations between information measures and estimation measures have been derived for Poisson channels. More recently, Verdú [13] extended the I-MMSE relation of Gaussian noise to the paradigm of mismatched conditional mean estimation, that is, to deal with an estimator that is optimally matched to a wrong probability distribution assumed on the input signal. The excess mean squared error (MSE) due to this mismatch was shown to be related to the Kullback–Leibler divergence between the channel output distributions corresponding to the true and the assumed input distributions (see also [4] for a further study in this direction). In [9], the I-MMSE relation was further investigated from a statistical physics perspective, where among other results, it was demonstrated how statistical–mechanical tools can be harnessed in order to assess the MMSE via the I-MMSE relation of [5], using the fact that in many cases, the mutual information can be viewed as the partition function of a certain physical system.

This paper is a further development in the above described direction of [9]. The main idea is that, for the purpose of evaluating the covariance matrix of the MMSE estimator, one may use a conceptually simple and more direct relationship between the MMSE covariance matrix and other information measures, that can also be presented in the form of a certain partition function and hence be analyzed using methods of statistical physics. The main advantage of the proposed approach, over those of the I-MMSE relations and its variants, is its full generality: It applies, in principle, to any joint probability function $P(\mathbf{x}, \mathbf{y})$ of the channel input signal $\mathbf{x} = (x_1, \dots, x_n)$, to be estimated, and the channel output $\mathbf{y} = (y_1, \dots, y_m)$ (where m and n are positive integers),

provided that certain technical regularity conditions hold. The channel $P(\mathbf{y}|\mathbf{x})$ does not even have to be additive, as opposed to the assumptions made in [6] and [12]. Moreover, the dimension m of the channel output vector \mathbf{y} does not have to be the same as the dimension n of the input vector \mathbf{x} .

In a nutshell, the idea is to define, for a given n -vector of real-valued parameters $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$, the ‘partition function’

$$Z(\mathbf{y}, \boldsymbol{\lambda}) = \sum_{\mathbf{x}} \exp \left\{ \sum_{i=1}^n \lambda_i x_i \right\} P(\mathbf{x}, \mathbf{y}),$$

where we have implicitly assumed that \mathbf{x} takes on discrete values, otherwise, the sum should simply be replaced by an integral. Now, it is straightforward to show that the gradient of $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$ w.r.t. $\boldsymbol{\lambda}$, computed at $\boldsymbol{\lambda} = 0$, gives the conditional mean estimator $\hat{\mathbf{x}} = \mathbf{E}\{\mathbf{X}|\mathbf{y}\}$, whereas the expectation of the Hessian of the same function, again, at $\boldsymbol{\lambda} = 0$, gives the error covariance matrix of the MMSE estimator. As we shall see in the sequel, $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$ lends itself to closed form analytic evaluation (in the spirit of a single-letter formula) in a fairly wide spectrum of situations, using methods of statistical mechanics. Thus, the MMSE estimator and its performance can quite easily be derived too in these situations. Moreover, as was demonstrated extensively in [9], the statistical-mechanical perspective on estimation-theoretic problems, may offer, not only analysis techniques, but also some important insights with regard to threshold effects (whenever existent) via the inspection of possible *phase transitions* in the parallel statistical-mechanical model.

Besides the general applicability of this approach, it has several additional advantages:

1. As mentioned in the previous paragraph, it provides, not only the MMSE error covariance matrix, but also the conditional mean estimator itself.
2. As will be seen, several variants of these relations between estimation measures and information measures can be offered. In some cases, one of the relations may be more convenient to work with than the others.
3. The approach is easy to extend to the mismatched case. Furthermore, it allows mismatch in both the source and the channel (as opposed to [13], which allows mismatch in the source only).

The remaining part of this paper is organized as follows. In Section 2, we establish notation conventions. In Section 3, we first derive the basic relations between the conditional mean estimator, as well as its error covariance matrix, and the above-mentioned partition function. In the same section, we also discuss this relation and derive a few variants that involve also information measures, like the information density, the Fisher information, etc. We also outline the extension to mismatched estimation. In Section 4, we provide three examples. In Section 5, we show how two of them set the stage to the analysis of a more general class of joint distributions, $P(\mathbf{x}, \mathbf{y})$. Finally, in Section 6, we summarize and conclude the paper.

2 Notation Conventions

Throughout this paper, scalar random variables (RV's) will be denoted by capital letters, their sample values will be denoted by the respective lower case letters, and their alphabets will be denoted by the respective calligraphic letters. A similar convention will apply to random vectors and their sample values, which will be denoted with same symbols in the bold face font. Thus, for example, \mathbf{X} will denote a random vector (X_1, \dots, X_n) , and $\mathbf{x} = (x_1, \dots, x_n)$ is a specific vector value in \mathcal{X}^n , the n -th Cartesian power of \mathcal{X} . The notations y_i^j and Y_i^j , where i and j are integers and $i \leq j$, will designate segments (y_i, \dots, y_j) and (Y_i, \dots, Y_j) , respectively.

Probability functions will be denoted generically by the letter P or Q . In particular, $P(\mathbf{x}, \mathbf{y})$ is the joint probability mass function (in the discrete case) or the joint density (in the continuous case) of the desired channel input vector $\mathbf{x} = (x_1, \dots, x_n)$ and the observed channel output vector $\mathbf{y} = (y_1, \dots, y_m)$. Accordingly, $P(\mathbf{x})$ will denote the marginal of \mathbf{x} , $P(\mathbf{y}|\mathbf{x})$ will denote the conditional probability mass (or density) of \mathbf{y} given \mathbf{x} , induced by the channel, and so on. Whenever there is room for ambiguity, these probability functions will be subscripted by the names of the random variables and the conditionings, according to standard notation conventions in probability theory and information theory. Throughout the sequel, we will assume discrete valued alphabets, mostly for the sake of simplicity and convenience. Extensions to continuous valued situations will be straightforward with summations being replaced by integrations, etc. Indeed, some of our examples will involve continuous valued random variables.

The expectation operator of a generic function $f(\mathbf{x}, \mathbf{y})$ w.r.t. the joint distribution P of (\mathbf{X}, \mathbf{Y})

will be denoted by $\mathbf{E}\{f(\mathbf{X}, \mathbf{Y})\}$. The conditional expectation of the same function given that $\mathbf{Y} = \mathbf{y}$, denoted $\mathbf{E}\{f(\mathbf{X}, \mathbf{Y})|\mathbf{Y} = \mathbf{y}\}$, and which is obviously identical to $\mathbf{E}\{f(\mathbf{X}, \mathbf{y})|\mathbf{Y} = \mathbf{y}\}$, is, of course, a function of \mathbf{y} . On substituting \mathbf{Y} in this function, this becomes then a random variable which will be denoted by $\mathbf{E}\{f(\mathbf{X}, \mathbf{Y})|\mathbf{Y}\}$. When using vectors and matrices in a linear-algebraic format, n -dimensional vectors, like \mathbf{x} (and \mathbf{X}), will be understood as column vectors, the operator $(\cdot)^T$ will denote vector or matrix transposition, and so, \mathbf{x}^T would be a row vector. For two positive sequences $\{a_n\}$ and $\{b_n\}$, the notation $a_n \doteq b_n$ means equivalence in the exponential order, i.e., $\lim_{n \rightarrow \infty} \frac{1}{n} \log(a_n/b_n) = 0$. Finally, the indicator function of an event \mathcal{A} will be denoted by $1\{\mathcal{A}\}$. I.e., $1\{\mathcal{A}\} = 1$ if \mathcal{A} occurs, and $1\{\mathcal{Z}\} = 0$ if not.

3 MMSE Estimation Relations

This section consists of two subsections. In the first, we derive the main basic relations and in the second, we show how to extend the scope to the case of mismatched estimation.

3.1 Basic Relations

Let $\mathbf{X} = (X_1, \dots, X_n)$, and $\mathbf{Y} = (Y_1, \dots, Y_m)$ (n and m being positive integers), be two random vectors, jointly distributed according to a given probability function $P(\mathbf{x}, \mathbf{y})$. It is further assumed that the alphabet \mathcal{X} , of each component of \mathbf{X} , consists of a set of real valued numbers, i.e., $\mathcal{X} \subseteq \mathbb{R}$. This assumption is obviously necessary in order to make the problem of estimating \mathbf{X} , in the MSE sense, a meaningful problem. The conditional mean estimator of \mathbf{X} based on \mathbf{Y} , i.e., $\hat{\mathbf{X}} = \mathbf{E}\{\mathbf{X}|\mathbf{Y}\}$ is well-known to be the optimum estimator in the MSE sense, i.e., it minimizes the MSE $\mathbf{E}\{(X_i - \hat{X}_i)^2\}$ for all $i = 1, 2, \dots, n$. The MMSE in estimating X_i is then $\mathbf{E}\{(X_i - \mathbf{E}\{X_i|\mathbf{Y}\})^2\}$, i.e., the expected conditional variance of X_i given \mathbf{Y} . More generally, the MMSE error covariance matrix E is an $n \times n$ matrix whose (i, j) -th element is given by $\mathbf{E}\{(X_i - \mathbf{E}\{X_i|\mathbf{Y}\})(X_j - \mathbf{E}\{X_j|\mathbf{Y}\})\}$. This matrix can be represented as the expectation (w.r.t. \mathbf{Y}) of the conditional covariance matrix of \mathbf{X} given \mathbf{Y} , henceforth denoted $\text{Cov}\{\mathbf{X}|\mathbf{Y}\}$. I.e.,

$$E = \mathbf{E}\{\text{Cov}\{\mathbf{X}|\mathbf{Y}\}\} = \mathbf{E}\{\mathbf{X}\mathbf{X}^T\} - \mathbf{E}\{\mathbf{E}\{\mathbf{X}|\mathbf{Y}\} \cdot \mathbf{E}\{\mathbf{X}^T|\mathbf{Y}\}\}.$$

Defining a column vector of n real valued parameters, $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)^T$, consider the following

function:

$$Z(\mathbf{y}, \boldsymbol{\lambda}) \triangleq \sum_{\mathbf{x} \in \mathcal{X}^n} \exp\{\boldsymbol{\lambda}^T \mathbf{x}\} P(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{x} \in \mathcal{X}^n} \exp\{\boldsymbol{\lambda}^T \mathbf{x}\} P(\mathbf{x}) P(\mathbf{y}|\mathbf{x}),$$

where it is assumed that the sum (or integral, in the continuous case) converges uniformly at least in some neighborhood of $\boldsymbol{\lambda} = 0$.¹ It is straightforward to see now that:

$$\left. \frac{\partial \ln Z(\mathbf{y}|\boldsymbol{\lambda})}{\partial \lambda_i} \right|_{\boldsymbol{\lambda}=0} = \frac{\sum_{\mathbf{x} \in \mathcal{X}^n} x_i P(\mathbf{x}, \mathbf{y})}{P(\mathbf{y})} = \sum_{x_i \in \mathcal{X}} x_i P(x_i|\mathbf{y}) = \mathbf{E}\{X_i|\mathbf{y}\}, \quad (1)$$

i.e.,

$$\mathbf{E}\{\mathbf{X}|\mathbf{y}\} = \nabla_{\boldsymbol{\lambda}} \ln Z(\mathbf{y}, \boldsymbol{\lambda}), \quad (2)$$

where $\nabla_{\boldsymbol{\lambda}}$ denotes the gradient w.r.t. $\boldsymbol{\lambda}$. Similarly, upon taking second order derivatives, one obtains

$$\left. \frac{\partial^2 \ln Z(\mathbf{y}|\boldsymbol{\lambda})}{\partial \lambda_i \partial \lambda_j} \right|_{\boldsymbol{\lambda}=0} = \mathbf{E}\{X_i X_j|\mathbf{y}\} - \mathbf{E}\{X_i|\mathbf{y}\} \cdot \mathbf{E}\{X_j|\mathbf{y}\} = \text{Cov}\{X_i, X_j|\mathbf{y}\},$$

and so,

$$\mathbf{E} = \mathbf{E} \left\{ \nabla_{\boldsymbol{\lambda}}^2 \ln Z(\mathbf{Y}, \boldsymbol{\lambda}) \Big|_{\boldsymbol{\lambda}=0} \right\}, \quad (3)$$

where $\nabla_{\boldsymbol{\lambda}}^2$ is the Hessian w.r.t. $\boldsymbol{\lambda}$, namely, the matrix of second order derivatives w.r.t. pairs of components of $\boldsymbol{\lambda}$. Note that here and throughout the sequel, we will always refer to gradients and Hessians of functions w.r.t. $\boldsymbol{\lambda}$, computed at the point $\boldsymbol{\lambda} = 0$. It will therefore be convenient to use, for a generic function g , the shorthand notations $\nabla_0 g(\boldsymbol{\lambda})$ and $\nabla_0^2 g(\boldsymbol{\lambda})$ to designate $\nabla_{\boldsymbol{\lambda}} g(\boldsymbol{\lambda}) \Big|_{\boldsymbol{\lambda}=0}$ and $\nabla_{\boldsymbol{\lambda}}^2 g(\boldsymbol{\lambda}) \Big|_{\boldsymbol{\lambda}=0}$, respectively.

Another, perhaps simpler, way to look at the relations (2) and (3) is the following: Obviously, for a given \mathbf{y} , $M(\mathbf{y}, \boldsymbol{\lambda}) = \sum_{\mathbf{x}} e^{\boldsymbol{\lambda}^T \mathbf{x}} P(\mathbf{x}|\mathbf{y})$ is the moment generating function pertaining to the conditional distribution of \mathbf{x} given \mathbf{y} and so, its derivatives relative to $\{\lambda_i\}$, computed at $\boldsymbol{\lambda} = 0$, yield the conditional moments $\mathbf{E}\{X_i|\mathbf{y}\}$, $\mathbf{E}\{X_i^2|\mathbf{y}\}$, $\mathbf{E}\{X_i X_j|\mathbf{y}\}$, etc. Therefore, $\ln M(\mathbf{y}, \boldsymbol{\lambda})$ is a generator of the corresponding conditional cumulants, $\mathbf{E}\{X_i|\mathbf{y}\}$, $\text{Var}\{X_i|\mathbf{y}\}$, $\text{Cov}\{X_i, X_j|\mathbf{y}\}$, etc. Now, observe that $\ln M(\mathbf{y}, \boldsymbol{\lambda})$ differs from $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$ merely by the additive term $\ln P(\mathbf{y})$, which does not depend on $\boldsymbol{\lambda}$ anyway and hence does not affect the gradient and Hessian w.r.t. $\boldsymbol{\lambda}$. Therefore, $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$ is a generator of conditional cumulants, exactly like $\ln M(\mathbf{y}, \boldsymbol{\lambda})$. An important point, however, is that we prefer $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$ over $\ln M(\mathbf{y}, \boldsymbol{\lambda})$ because normally, it is more convenient to

¹If this assumption is not met, one can instead, parametrize each component λ_i of $\boldsymbol{\lambda}$ as a purely imaginary number $\lambda_i = j\omega_i$ ($j = \sqrt{-1}$), as is done in the definition of the characteristic function.

work with the joint distribution $P(\mathbf{x}, \mathbf{y})$ (or equivalently, with the source $P(\mathbf{x})$ and forward channel $P(\mathbf{y}|\mathbf{x})$) rather than with the backward channel (or the posterior) $P(\mathbf{x}|\mathbf{y})$.²

We next derive several alternative versions of this relation between the error covariance matrix of the MMSE estimator and derivatives of $\ln Z$. First, observe that $Z(\mathbf{y}, \boldsymbol{\lambda})$ is proportional to $P_{\boldsymbol{\lambda}}(\mathbf{y}) \cdot \Theta(\boldsymbol{\lambda})$, where

$$\Theta(\boldsymbol{\lambda}) = \sum_{\mathbf{x} \in \mathcal{X}^n} P(\mathbf{x}) \exp\{\boldsymbol{\lambda}^T \mathbf{x}\}$$

and $P_{\boldsymbol{\lambda}}(\mathbf{y})$ is the output marginal of \mathbf{y} induced by the channel $P(\mathbf{y}|\mathbf{x})$ and the modified source distribution $P_{\boldsymbol{\lambda}}(\mathbf{x}) \triangleq e^{\boldsymbol{\lambda}^T \mathbf{x}} P(\mathbf{x}) / \Theta(\boldsymbol{\lambda})$. We therefore obtain

$$\begin{aligned} E &= \mathbf{E} \left\{ \nabla_0^2 \ln Z(\mathbf{Y}, \boldsymbol{\lambda}) \right\} \\ &= \mathbf{E} \left\{ \nabla_0^2 \ln [P_{\boldsymbol{\lambda}}(\mathbf{Y}) \cdot \Theta(\boldsymbol{\lambda})] \right\} \\ &= \nabla_0^2 \ln \Theta(\boldsymbol{\lambda}) + \mathbf{E} \left\{ \nabla_0^2 \ln P_{\boldsymbol{\lambda}}(\mathbf{Y}) \right\} \\ &= \text{Cov}\{\mathbf{X}\} - J, \end{aligned} \tag{4}$$

where $\text{Cov}\{\mathbf{X}\} = \mathbf{E}\{\mathbf{X}\mathbf{X}^T\} - \mathbf{E}\{\mathbf{X}\} \cdot \mathbf{E}\{\mathbf{X}^T\}$ is the covariance matrix of \mathbf{X} and J is the Fisher information matrix of estimating $\boldsymbol{\lambda}$ based on \mathbf{Y} , computed at the point $\boldsymbol{\lambda} = 0$. The Fisher information matrix J can also be expressed as

$$J = \mathbf{E} \left\{ \nabla_0 \ln P_{\boldsymbol{\lambda}}(\mathbf{Y}) \cdot \nabla_0^T \ln P(\mathbf{Y}|\boldsymbol{\lambda}) \right\}.$$

Equivalently, we obtained

$$J = \text{Cov}\{\mathbf{X}\} - E = \mathbf{E}\{\mathbf{E}\{\mathbf{X}|\mathbf{Y}\} \cdot \mathbf{E}\{\mathbf{X}^T|\mathbf{Y}\}\}.$$

Note that J can also be obtained as the negative expectation of the Hessian (or, equivalently, as the covariance matrix of the gradient) of the *information density* [14],

$$i_{\boldsymbol{\lambda}}(\mathbf{x}; \mathbf{y}) = \ln[P(\mathbf{y}|\mathbf{x})/P_{\boldsymbol{\lambda}}(\mathbf{y})],$$

²As a side remark, we shall mention also the physical perspective: if $Z(\mathbf{y}, \boldsymbol{\lambda})$ is thought of as the partition function of a certain statistical–mechanical model (as discussed in the Introduction), where the components of $\boldsymbol{\lambda}$ are thought of as certain generalized forces or fields that are acting on the individual particles, then the above relation between the second order derivative of $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$ w.r.t. λ_i and λ_j and the (conditional) covariances between the corresponding state variables, X_i and X_j , is known as one of the versions of the fluctuation–dissipation theorem in statistical mechanics [10, p. 32, eq. (2.44)], which relates between the linear response of the system (to an infinitesimally small perturbation in its parameters) and its fluctuations in equilibrium.

which is again, computed at $\boldsymbol{\lambda} = 0$.

Sometimes it is more convenient to square the first derivative of $\ln Z$ than to take the second derivative. In these cases, the following relationship may be useful:

$$\begin{aligned}
\Xi &\triangleq \mathbf{E} \left\{ [\nabla_0 \ln Z(\mathbf{Y}, \boldsymbol{\lambda})] \cdot [\nabla_0 \ln Z(\mathbf{Y}, \boldsymbol{\lambda})]^T \right\} \\
&= \mathbf{E} \left\{ [\nabla_0 \ln \{P(\mathbf{Y}|\boldsymbol{\lambda}) \cdot \Theta(\boldsymbol{\lambda})\}] \cdot [\nabla_0 \ln \{P(\mathbf{Y}|\boldsymbol{\lambda}) \cdot \Theta(\boldsymbol{\lambda})\}]^T \right\} \\
&= \mathbf{E} \left\{ [\nabla_0 \ln P(\mathbf{Y}|\boldsymbol{\lambda})] \cdot [\nabla_0 \ln P(\mathbf{Y}|\boldsymbol{\lambda})]^T \right\} + [\nabla_0 \ln \Theta(\boldsymbol{\lambda})] \cdot [\nabla_0 \ln \Theta(\boldsymbol{\lambda})]^T \\
&= J + \mathbf{E}\{\mathbf{X}\} \cdot \mathbf{E}\{\mathbf{X}^T\} \\
&= \text{Cov}\{\mathbf{X}\} + \mathbf{E}\{\mathbf{X}\} \cdot \mathbf{E}\{\mathbf{X}^T\} - E \\
&= \mathbf{E}\{\mathbf{X}\mathbf{X}^T\} - E
\end{aligned} \tag{5}$$

and so,

$$E = \mathbf{E}\{\mathbf{X}\mathbf{X}^T\} - \Xi.$$

Particularizing these results to the MMSE,

$$\text{mmse}(\mathbf{X}|\mathbf{Y}) \triangleq \sum_{i=1}^n \mathbf{E}\{(X_i - \mathbf{E}\{X_i|\mathbf{Y}\})^2\},$$

which is the trace of E , we have the following relations, which we formulate as a proposition.

Proposition 1. The following formulas for the MMSE hold:

$$\text{mmse}(\mathbf{X}|\mathbf{Y}) = \sum_{i=1}^n \mathbf{E} \left\{ \left. \frac{\partial^2 \ln Z(\mathbf{Y}, \boldsymbol{\lambda})}{\partial \lambda_i^2} \right|_{\boldsymbol{\lambda}=0} \right\} \tag{6}$$

$$= \sum_{i=1}^n \left[\text{Var}\{X_i\} + \mathbf{E} \left\{ \left. \frac{\partial^2 \ln P(\mathbf{Y}|\boldsymbol{\lambda})}{\partial \lambda_i^2} \right|_{\boldsymbol{\lambda}=0} \right\} \right] \tag{7}$$

$$= \sum_{i=1}^n \left[\text{Var}\{X_i\} - \mathbf{E} \left\{ \left. \left[\frac{\partial \ln P(\mathbf{Y}|\boldsymbol{\lambda})}{\partial \lambda_i} \right]^2 \right|_{\boldsymbol{\lambda}=0} \right\} \right] \tag{8}$$

$$= \sum_{i=1}^n \left[\mathbf{E}\{X_i^2\} - \mathbf{E} \left\{ \left. \left[\frac{\partial \ln Z(\mathbf{Y}, \boldsymbol{\lambda})}{\partial \lambda_i} \right]^2 \right|_{\boldsymbol{\lambda}=0} \right\} \right] \tag{9}$$

In the second and the third formulas, $\ln P(\mathbf{Y}|\boldsymbol{\lambda})$ can be replaced by $\ln i(\mathbf{X}; \mathbf{Y})$, thus relating the MMSE to the information density.

3.2 Extension to the Mismatched Case

In this short subsection, we are outlining how our approach can easily be extended to handle situations of mismatched estimation. Consider a mismatched estimator which is the conditional mean of \mathbf{X} given \mathbf{Y} , based on an incorrect joint distribution $Q(\mathbf{x}, \mathbf{y})$, whereas the true joint distribution continues to be $P(\mathbf{x}, \mathbf{y})$. Denoting by $Z_P(\mathbf{y}, \boldsymbol{\lambda})$ and $Z_Q(\mathbf{y}, \boldsymbol{\lambda})$ the corresponding partition functions, and by \mathbf{E}_P and \mathbf{E}_Q , the corresponding expectations, our approach can easily be generalized to handle this case as follows:

$$\begin{aligned}
E &= \mathbf{E}_P \left\{ (\mathbf{X} - \mathbf{E}_Q\{\mathbf{X}|\mathbf{Y}\})(\mathbf{X}^T - \mathbf{E}_Q\{\mathbf{X}^T|\mathbf{Y}\}) \right\} \\
&= \mathbf{E}_P\{\mathbf{X}\mathbf{X}^T\} - \mathbf{E}_P\{\mathbf{E}_P\{\mathbf{X}|\mathbf{Y}\}\mathbf{E}_Q\{\mathbf{X}^T|\mathbf{Y}\}\} - \\
&\quad \mathbf{E}_P\{\mathbf{E}_Q\{\mathbf{X}|\mathbf{Y}\}\mathbf{E}_P\{\mathbf{X}^T|\mathbf{Y}\}\} + \mathbf{E}_P\{\mathbf{E}_Q\{\mathbf{X}|\mathbf{Y}\}\mathbf{E}_Q\{\mathbf{X}^T|\mathbf{Y}\}\} \\
&= \mathbf{E}_P\{\mathbf{X}\mathbf{X}^T\} - \mathbf{E}_P\{[\nabla_0 \ln Z_P(\mathbf{Y}, \boldsymbol{\lambda})] \cdot [\nabla_0 \ln Z_Q(\mathbf{Y}, \boldsymbol{\lambda})]^T\} - \\
&\quad \mathbf{E}_P\{[\nabla_0 \ln Z_Q(\mathbf{Y}, \boldsymbol{\lambda})] \cdot [\nabla_0 \ln Z_P(\mathbf{Y}, \boldsymbol{\lambda})]^T\} + \mathbf{E}_P\{[\nabla_0 \ln Z_Q(\mathbf{Y}, \boldsymbol{\lambda})] \cdot [\nabla_0 \ln Z_Q(\mathbf{Y}, \boldsymbol{\lambda})]^T\}.
\end{aligned}$$

Thus, in particular, the MSE associated with the mismatched estimator is given by

$$\begin{aligned}
\text{mse}_Q(\mathbf{X}|\mathbf{Y}) &= \sum_{i=1}^n \left[\mathbf{E}_P\{X_i^2\} - 2\mathbf{E}_P \left\{ \left. \frac{\partial \ln Z_P(\mathbf{Y}, \boldsymbol{\lambda})}{\partial \lambda_i} \right|_{\boldsymbol{\lambda}=0} \cdot \left. \frac{\partial \ln Z_Q(\mathbf{Y}, \boldsymbol{\lambda})}{\partial \lambda_i} \right|_{\boldsymbol{\lambda}=0} \right\} \right. \\
&\quad \left. + \mathbf{E}_P \left\{ \left[\left. \frac{\partial \ln Z_Q(\mathbf{Y}, \boldsymbol{\lambda})}{\partial \lambda_i} \right|_{\boldsymbol{\lambda}=0} \right]^2 \right\} \right]. \tag{10}
\end{aligned}$$

4 Examples

In this section, we provide three examples, where we show how the log-partition function, $\ln Z(\mathbf{y}, \boldsymbol{\lambda})$, can be evaluated for large n , using methods of statistical mechanics. Using the relations derived in Subsection 3.1, we then show how the conditional mean estimator and the MMSE can be approximated for large n .

4.1 Example 1 – A Codeword Transmitted Over an AWGN

Our first example is taken from [9, Subsection 5.2], but here we demonstrate how to derive the conditional mean estimator and the MMSE using Proposition 1, rather than the I-MMSE relation. For the sake of completeness and convenience, we provide here the full necessary details (with

the appropriate modifications to accommodate the method proposed herein), including those that already appear in [9]. As noted in [9], the analysis of this model is intimately related to one of the statistical mechanical techniques used in the analysis of the so called random energy model (REM) of disordered magnetic materials, a.k.a. spin glasses in the statistical physics literature (see references in [9]).

Let \mathbf{X} be chosen uniformly at random from a codebook $\mathcal{C} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{M-1}\}$ of size $M = e^{nR}$. The codebook itself is also selected at random (and then revealed to the estimator) in the following manner: Each \mathbf{x}_i is selected independently and uniformly at random from the surface of a sphere of radius $\sqrt{nP_x}$ centered at the origin. The channel $P(\mathbf{y}|\mathbf{x})$ is an AWGN channel (hence $m = n$) whose noise variance is $1/\beta$ (keeping the same notation as in [9]). I.e.,

$$P(\mathbf{y}|\mathbf{x}) = \left(\frac{\beta}{2\pi}\right)^{n/2} \exp\left\{-\frac{\beta}{2}\|\mathbf{y} - \mathbf{x}\|^2\right\}.$$

Thus, for a given \mathbf{y} , we have:

$$\begin{aligned} Z(\mathbf{y}, \lambda) &= \sum_{\mathbf{x} \in \mathcal{C}} e^{-nR} \exp\{-\beta\|\mathbf{y} - \mathbf{x}\|^2/2 + \lambda^T \mathbf{x}\} \\ &= e^{-nR} \exp[-\beta\|\mathbf{y} - \mathbf{x}_0\|^2/2 + \lambda^T \mathbf{x}_0] + \sum_{\mathbf{x} \in \mathcal{C} \setminus \{\mathbf{x}_0\}} e^{-nR} \exp[-\beta\|\mathbf{y} - \mathbf{x}\|^2/2 + \lambda^T \mathbf{x}] \\ &\triangleq Z_c(\mathbf{y}, \lambda) + Z_e(\mathbf{y}, \lambda), \end{aligned} \quad (11)$$

where, without loss of generality, \mathbf{x}_0 designates the transmitted codeword. Now, since $\|\mathbf{y} - \mathbf{x}_0\|^2$ is typically around n/β , $Z_c(\mathbf{y}, \lambda)$ would typically be about $e^{-nR} e^{-\beta \cdot n/(2\beta)} e^{\lambda^T \mathbf{x}_0} = e^{-n(R+1/2) + \lambda^T \mathbf{x}_0}$. As for $Z_e(\mathbf{y}, \lambda)$, we have:

$$Z_e(\mathbf{y}, \lambda) \doteq e^{-nR} \int_{\mathbb{R}} d\epsilon N(\epsilon) e^{-\beta n \epsilon},$$

where $N(\epsilon)$ is the number of codewords $\{\mathbf{x}\}$ in $\mathcal{C} - \{\mathbf{x}_0\}$ for which $\|\mathbf{y} - \mathbf{x}\|^2/2 - \lambda^T \mathbf{x}/\beta \approx n\epsilon$, namely, between $n\epsilon$ and $n(\epsilon + d\epsilon)$. Now, given \mathbf{y} , $N(\epsilon) = \sum_{i=1}^M 1\{\mathbf{x}_i : \|\mathbf{y} - \mathbf{x}_i\|^2/2 - \lambda^T \mathbf{x}_i/\beta \approx n\epsilon\}$ is the sum of M i.i.d. Bernoulli random variables and so, its expectation is

$$\overline{N(\epsilon)} = \sum_{i=1}^M \Pr\{\|\mathbf{y} - \mathbf{X}_i\|^2/2 - \lambda^T \mathbf{X}_i/\beta \approx n\epsilon\} = e^{nR} \Pr\{\|\mathbf{y} - \mathbf{X}_1\|^2/2 - \lambda^T \mathbf{X}_1/\beta \approx n\epsilon\}. \quad (12)$$

Denoting $P_y = \frac{1}{n} \sum_{i=1}^n y_i^2$ (typically, P_y is about $P_x + 1/\beta$), the event $\|\mathbf{y} - \mathbf{x}\|^2/2 - \lambda^T \mathbf{x}/\beta \approx n\epsilon$ is equivalent to the event $\mathbf{x}^T(\mathbf{y} + \lambda/\beta) \approx [(P_x + P_y)/2 - \epsilon]n$ or equivalently,

$$\rho(\mathbf{x}, \mathbf{y}) \triangleq \frac{\mathbf{x}^T(\mathbf{y} + \lambda/\beta)}{n\sqrt{P_x P_y'}} \approx \frac{\frac{1}{2}(P_x + P_y) - \epsilon}{\sqrt{P_x P_y'}} \triangleq \frac{P_a - \epsilon}{P_g'},$$

where we have defined $P_a = (P_x + P_y)/2$ and $P'_g = \sqrt{P_x P_y}$, where $P'_y = \frac{1}{n} \sum_i (y_i + \lambda_i/\beta)^2$. The probability that a randomly chosen vector \mathbf{X} on the sphere would have an empirical correlation coefficient ρ with a given vector $\mathbf{y}' = \mathbf{y} + \boldsymbol{\lambda}/\beta$ (that is, \mathbf{X} falls within a cone of half angle $\arccos(\rho)$ around \mathbf{y}') is exponentially $\exp[\frac{n}{2} \ln(1 - \rho^2)]$. For convenience, let us define

$$\Gamma(\rho) = \frac{1}{2} \ln(1 - \rho^2)$$

so that we can write

$$\Pr\{\|\mathbf{y} - \mathbf{X}_1\|^2/2 - \boldsymbol{\lambda}^T \mathbf{X}_1/\beta \approx n\epsilon\} \doteq \exp\left\{n \Gamma\left(\frac{P_a - \epsilon}{P'_g}\right)\right\}.$$

If ϵ is such that

$$\Gamma\left(\frac{P_a - \epsilon}{P'_g}\right) > -R,$$

then the energy level ϵ will be typically populated with an exponential number of codewords, concentrated very strongly around its mean

$$\overline{N(\epsilon)} \doteq \exp\left\{n \left[R + \Gamma\left(\frac{P_a - \epsilon}{P'_g}\right) \right]\right\},$$

otherwise (which means that $\overline{N(\epsilon)}$ is exponentially small), the energy level ϵ will not be populated by any codewords typically. This means that the populated energy levels range between

$$\epsilon_1 \triangleq P_a - P'_g \sqrt{1 - e^{-2R}}$$

and

$$\epsilon_2 \triangleq P_a + P'_g \sqrt{1 - e^{-2R}},$$

or equivalently, the populated values of ρ range between $-\rho_*$ and $+\rho_*$ where $\rho_* = \sqrt{1 - e^{-2R}}$. By large deviations and saddle-point methods, it follows that for a typical realization of the randomly chosen code, we have

$$\begin{aligned} Z_e(\mathbf{y}, \boldsymbol{\lambda}) &\doteq e^{-nR} \max_{\epsilon \in [\epsilon_1, \epsilon_2]} \exp\left\{n \left[R + \Gamma\left(\frac{P_a - \epsilon}{P'_g}\right) - \beta\epsilon \right]\right\} \\ &= \max_{\epsilon \in [\epsilon_1, \epsilon_2]} \exp\left\{n \left[\Gamma\left(\frac{P_a - \epsilon}{P'_g}\right) - \beta\epsilon \right]\right\} \\ &= \exp\left\{n \left[\max_{|\rho| \leq \rho_*} \left\{ \frac{1}{2} \ln(1 - \rho^2) - \beta(P_a - \rho P'_g) \right\} \right]\right\}. \end{aligned} \quad (13)$$

The derivative of $\frac{1}{2} \ln(1 - \rho^2) + \rho\beta P'_g$ w.r.t. ρ vanishes within $[-1, 1]$ at:

$$\rho = \rho_\beta \triangleq \sqrt{1 + \theta^2} - \theta$$

where

$$\theta \triangleq \frac{1}{2\beta P'_g}.$$

This is the maximizer as long as $\sqrt{1 + \theta^2} - \theta \leq \rho_*$, namely, $\theta > e^{-2R}/2\rho_*$, or equivalently, $\beta < \rho_* e^{2R}/P'_g$, which for $P'_g = \sqrt{P_x(P_x + 1/\beta)}$ ($\|\boldsymbol{\lambda}\|$ is small), is equivalent to $\beta < \beta_R \triangleq (e^{2R} - 1)/P_x$.

Thus, for the typical code we have

$$Z_e(\beta|\mathbf{y}) \doteq \begin{cases} \exp\left\{n\left[\frac{1}{2}\ln(1 - \rho_\beta^2) - \beta(P_a - \rho_\beta P'_g)\right]\right\}, & \beta < \beta_R \\ \exp\{-n[R + \beta(P_a - \rho_* P'_g)]\}, & \beta \geq \beta_R. \end{cases}$$

Taking now into account $Z_c(\mathbf{y}, \boldsymbol{\lambda})$, it is easy to see that for $\beta \geq \beta_R$ (which means $R < C$), $Z_c(\mathbf{y}, \boldsymbol{\lambda})$ dominates $Z_e(\mathbf{y}, \boldsymbol{\lambda})$, whereas for $\beta < \beta_R$ it is the other way around. It follows then that

$$Z(\mathbf{y}, \boldsymbol{\lambda}) \doteq \begin{cases} \exp\left\{n\left[\frac{1}{2}\ln(1 - \rho_\beta^2) - \beta(P_a - \rho_\beta P'_g)\right]\right\}, & \beta < \beta_R \\ \exp\left\{-n\left(R + \frac{1}{2}\right) + \boldsymbol{\lambda}^T \mathbf{x}_0\right\}, & \beta \geq \beta_R. \end{cases}$$

A very similar analysis applies also to the derivative $\frac{\partial}{\partial \lambda_i} \ln Z(\mathbf{y}, \boldsymbol{\lambda})$, which is essentially a weighted average of x_i with weights proportional to $\overline{N(\epsilon)} e^{-\beta \epsilon}$ for all $\epsilon \in [\epsilon_1, \epsilon_2]$. Thus, the exponentially dominant weight is due to the term that maximizes the exponent. Assuming that the correct codeword \mathbf{x}_0 is dominant ($Z_c \gg Z_e$, which is the case when $R < C$), this weighted average is obviously dominated by the i -th component of \mathbf{x}_0 , in which case the MMSE essentially vanishes. Otherwise, for $R > C$, Z_e dominates the partition function and the weighted average is overwhelmingly dominated by the term corresponding to the maximizing ϵ , or equivalently, the maximizing ρ , which is ρ_β . This means that the conditional mean estimator of X_i is approximately given by:

$$\begin{aligned} \mathbf{E}\{X_i|\mathbf{y}\} &\approx \frac{\partial}{\partial \lambda_i} \left[\frac{n}{2} \ln(1 - \rho_\beta^2) + \beta \rho_\beta n P'_g \right] \\ &= -\frac{n \rho_\beta}{1 - \rho_\beta^2} \frac{\partial \rho_\beta}{\partial \lambda_i} + \beta n P'_g \frac{\partial \rho_\beta}{\partial \lambda_i} + \beta \rho_\beta n \frac{\partial P'_g}{\partial \lambda_i} \\ &= n \frac{\partial \rho_\beta}{\partial \lambda_i} \left(-\frac{\rho_\beta}{1 - \rho_\beta^2} + \beta P'_g \right) + \beta n \rho_\beta \frac{\partial P'_g}{\partial \lambda_i} \\ &= n \beta \rho_\beta \frac{\partial P'_g}{\partial \lambda_i} \\ &= \beta \rho_\beta n \cdot \frac{\sqrt{P_x}}{2\sqrt{P_y}} \cdot \frac{2y_i}{\beta n} \end{aligned}$$

$$\begin{aligned}
&= \rho\beta \sqrt{\frac{P_x}{P_x + 1/\beta}} \cdot y_i \\
&= \frac{P_x}{P_x + 1/\beta} \cdot y_i,
\end{aligned} \tag{14}$$

where in the last step we have used the identity $\rho\beta = \sqrt{P_x/(P_x + 1/\beta)}$, which can easily be verified. This is simply the linear Wiener estimator that would have been applied had the input been zero-mean, i.i.d. Gaussian, with variance $1/\beta$ (see also [9]). According to Proposition 1, the MMSE associated with X_i is given by

$$\mathbf{E}\{(X_i - \mathbf{E}\{X_i|\mathbf{Y}\})^2\} \approx P_x - \mathbf{E}\{\mathbf{E}^2(X_i|\mathbf{Y})\} = P_x - \left(\frac{P_x}{P_x + 1/\beta}\right)^2 \cdot (P_x + 1/\beta) = \frac{P_x}{1 + \beta P_x},$$

as expected.

4.2 Example 2 – The Curie–Weiss Model

Consider a binary source

$$P(\mathbf{x}) = C_n \exp \left\{ \frac{a}{2n} \left(\sum_{i=1}^n x_i \right)^2 + b \sum_{i=1}^n x_i \right\} \quad \mathbf{x} \in \{-1, +1\}^n$$

where a and b are parameters and C_n is a normalization constant, which is immaterial for our purposes (as it is going to disappear upon taking derivatives w.r.t. $\{\lambda_i\}$, and the same comment applies to the constants C'_n and C''_n below). Let the channel be binary and symmetric, i.e.,

$$P(y|x) = \frac{e^{\beta xy}}{2 \cosh(\beta)}, \quad y \in \{-1, +1\}.$$

Then, the partition function $Z(\mathbf{y}, \boldsymbol{\lambda})$ can be represented as a one-dimensional integral using the Hubbard–Stratonovich transform, which in turn can be assessed using saddle point methods, as is frequently done in the statistical physics literature. Specifically, we have the following:

$$Z(\mathbf{y}, \boldsymbol{\lambda}) = C'_n \sum_{\mathbf{x}} \exp \left\{ \frac{a}{2n} \left(\sum_{i=1}^n x_i \right)^2 + b \sum_{i=1}^n x_i + \beta \sum_{i=1}^n x_i y_i + \sum_{i=1}^n \lambda_i x_i \right\} \tag{15}$$

$$= C'_n \sum_{\mathbf{x}} \exp \left\{ \sum_{i=1}^n x_i (\beta y_i + \lambda_i + b) + \frac{a}{2n} \left(\sum_{i=1}^n x_i \right)^2 \right\} \tag{16}$$

$$= C''_n \sum_{\mathbf{x}} \exp \left\{ \sum_{i=1}^n x_i (\beta y_i + \lambda_i + b) \right\} \cdot \int_{-\infty}^{+\infty} d\theta \exp \left\{ -\frac{n\theta^2}{2a} + \theta \sum_{i=1}^n x_i \right\} \tag{17}$$

$$= C_n'' \int_{-\infty}^{+\infty} d\theta e^{-n\theta^2/(2a)} \sum_{\mathbf{x}} \exp \left\{ \sum_{i=1}^n x_i (\beta y_i + \lambda_i + b + \theta) \right\} \quad (18)$$

$$= C_n'' \int_{-\infty}^{+\infty} d\theta e^{-n\theta^2/(2a)} \prod_{i=1}^n [2 \cosh(\beta y_i + \lambda_i + b + \theta)] \quad (19)$$

$$= 2^n C_n'' \int_{-\infty}^{+\infty} d\theta \exp \left\{ -\frac{n\theta^2}{2a} + \sum_{i=1}^n \ln \cosh(\beta y_i + \lambda_i + b + \theta) \right\}. \quad (20)$$

Thus,

$$\begin{aligned} \frac{\partial \ln Z(\mathbf{y}, \boldsymbol{\lambda})}{\partial \lambda_i} &= \frac{\int_{-\infty}^{+\infty} d\theta \tanh(\beta y_i + \lambda_i + b + \theta) \exp \left\{ -\frac{n\theta^2}{2a} + \sum_{i=1}^n \ln \cosh(\beta y_i + \lambda_i + b + \theta) \right\}}{\int_{-\infty}^{+\infty} d\theta \exp \left\{ -\frac{n\theta^2}{2a} + \sum_{i=1}^n \ln \cosh(\beta y_i + \lambda_i + b + \theta) \right\}} \\ &\approx \tanh(\beta y_i + \lambda_i + b + \theta_*), \end{aligned} \quad (21)$$

where θ_* is the maximizer of the expression at the exponent, i.e., it is the solution to the zero-derivative equation:

$$\theta = \frac{a}{n} \sum_{i=1}^n \tanh(\beta y_i + \lambda_i + b + \theta).$$

Thus, the MMSE estimator is:

$$\mathbf{E}\{X_i|\mathbf{y}\} = \frac{\int_{-\infty}^{+\infty} d\theta \tanh(\beta y_i + b + \theta) \exp \left\{ -\frac{n\theta^2}{2a} + \sum_{i=1}^n \ln \cosh(\beta y_i + b + \theta) \right\}}{\int_{-\infty}^{+\infty} d\theta \exp \left\{ -\frac{n\theta^2}{2a} + \sum_{i=1}^n \ln \cosh(\beta y_i + b + \theta) \right\}} \quad (22)$$

$$\approx \tanh(\beta y_i + b + \theta_*), \quad (23)$$

where now θ_* is understood to be taken with $\boldsymbol{\lambda} = 0$. For $b \neq 0$, the asymptotic MMSE is then given by

$$\lim_{n \rightarrow \infty} \frac{\text{mmse}(\mathbf{X}|\mathbf{Y})}{n} = 1 - \mathbf{E}\{\tanh^2(\beta Y + b + \theta_0)\},$$

where θ_0 is the solution to the equation

$$\theta = a \mathbf{E}\{\tanh(\beta Y + b + \theta)\},$$

and where Y is a binary $\{\pm 1\}$ RV, with mean $m^* \tanh(\beta)$, m^* being the dominant solution to the equation $m = \tanh(am + b)$, i.e., the maximizer of $h_2((1+m)/2) + am^2/2 + bm$, where $h_2(\cdot)$ is the binary entropy function. When $b = 0$, θ_0 becomes a random variable which takes on, with equal probabilities, one of two values, each one being the solution to the above displayed equation, except that in one of them Y has mean $m^* \tanh(\beta)$ and in the other, its mean is $-m^* \tanh(\beta)$.

This calculation is intimately related to the Curie–Weiss model of magnetic spins [10, Subsection 2.5.2, pp. 40–44], where the parameter m plays the role of magnetization.

4.3 Example 3 – The Generalized Multivariate Cauchy Noise Model

Let $X_i \sim \mathcal{N}(0, \sigma^2)$ be i.i.d. RV's, and let the additive noise have a generalized multivariate Cauchy distribution, i.e.,

$$P(\mathbf{y}|\mathbf{x}) = \frac{C_{n,k}}{[1 + (\mathbf{y} - \mathbf{x})^T S (\mathbf{y} - \mathbf{x})]^k}$$

where $C_{n,k}$ is a normalization constant, S is a positive definite matrix, and $k > 0$ is chosen large enough (as a function of n) such $\int_{\mathbb{R}^n} d\mathbf{z}/[1 + \mathbf{z}^T S \mathbf{z}]^k < \infty$, i.e., $k > n/2$. The choice $k = (n + 1)/2$ corresponds to the ordinary multivariate Cauchy distribution. Here, however, we will require moreover that k is even large enough such that the second moments exist, i.e., $\int_{\mathbb{R}^n} d\mathbf{z} \cdot \mathbf{z}^T \mathbf{z}/[1 + \mathbf{z}^T S \mathbf{z}]^k < \infty$, which means $k > n/2 + 1$. For simplicity, we will take S to be the identity matrix. However, our analysis easily extends to a general positive matrix S , as well as to a general Gaussian vector \mathbf{X} , not necessarily with i.i.d. components. Using the Laplace transform identity $\int_0^\infty dt \cdot t^{k-1} e^{-st} = \Gamma(k)/s^k$, we have:

$$Z(\mathbf{y}, \boldsymbol{\lambda}) = \int_{\mathbb{R}^n} d\mathbf{x} P(\mathbf{x}) e^{\boldsymbol{\lambda}^T \mathbf{x}} \cdot \frac{C_{n,k}}{[1 + \sum_{i=1}^n (y_i - x_i)^2]^k} \quad (24)$$

$$= C_{n,k} \int_{\mathbb{R}^n} d\mathbf{x} P(\mathbf{x}) e^{\boldsymbol{\lambda}^T \mathbf{x}} \int_0^\infty dt \cdot \frac{t^{k-1}}{\Gamma(k)} \cdot e^{-t[1 + \sum_i (y_i - x_i)^2]} \quad (25)$$

$$= C'_{n,k} \int_0^\infty dt \cdot t^{k-1} e^{-t} \int_{\mathbb{R}^n} d\mathbf{x} P(\mathbf{x}) e^{\boldsymbol{\lambda}^T \mathbf{x}} \cdot e^{-t \sum_i (y_i - x_i)^2} \quad (26)$$

$$= C''_{n,k} \int_0^\infty dt \cdot t^{k-1} e^{-t} \prod_{i=1}^n \int_{\mathbb{R}} dx_i e^{-x_i^2/2\sigma^2} e^{\lambda_i x_i} \cdot e^{-t(y_i - x_i)^2} \quad (27)$$

$$= C''_{n,k} \int_0^\infty dt \cdot t^{k-1} e^{-t} \left(\frac{2\pi\sigma^2}{1 + 2t\sigma^2} \right)^{n/2} \cdot \exp \left\{ -t \sum_{i=1}^n y_i^2 + \sum_{i=1}^n \frac{(ty_i + \lambda_i/2)^2}{t + 1/2\sigma^2} \right\}. \quad (28)$$

and so,

$$\left. \frac{\partial \ln Z(\mathbf{y}, \boldsymbol{\lambda})}{\partial \lambda_i} \right|_{\boldsymbol{\lambda}=0} = \frac{\int_0^\infty dt \cdot \frac{ty_i}{t+1/2\sigma^2} e^{-t} t^{k-1} \exp \left\{ -\frac{n}{2} \ln(1 + 2t\sigma^2) - \frac{t}{1+2t\sigma^2} \sum_i y_i^2 \right\}}{\int_0^\infty dt e^{-t} t^{k-1} \exp \left\{ -\frac{n}{2} \ln(1 + 2t\sigma^2) - \frac{t}{1+2t\sigma^2} \sum_i y_i^2 \right\}}$$

which can be approximated by $\hat{t}y_i/(\hat{t} + 1/2\sigma^2)$, where \hat{t} is the value of t that dominates the integral, i.e.,

$$\hat{t} = \operatorname{argmax}_t \left[(k-1) \ln t - \frac{n}{2} \ln(1 + 2t\sigma^2) - \frac{t}{1 + 2t\sigma^2} \sum_i y_i^2 \right].$$

The derivation of the MMSE can be done in a similar manner as in the previous examples.

5 Joint Distributions with Generalized Spherical Symmetry

Examples 2 and 3 of the previous section have one idea in common. In both of them we expressed either the source or the channel as a one-dimensional integral over a variable (t or θ , in those examples), where for each value of this variable, we have a product form measure, which enables, after applying saddle point analysis on this integral, to pass to a closed-form formula, which has the flavor of a single-letter characterization. In this section, we generalize this idea to establish a somewhat more general framework.

Suppose that $m = n$ and the joint distribution of \mathbf{X} and \mathbf{Y} is of the form

$$P(\mathbf{x}, \mathbf{y}) = F_n\left(\sum_i \phi(x_i, y_i)\right).$$

Let $f_n(t)$ be the inverse Laplace transform of $F_n(s)$. Then, we have

$$\begin{aligned} Z(\mathbf{y}, \boldsymbol{\lambda}) &= \int_{\mathbb{R}^n} d\mathbf{x} e^{\boldsymbol{\lambda}^T \mathbf{x}} P(\mathbf{x}, \mathbf{y}) \\ &= \int_{\mathbb{R}^n} d\mathbf{x} e^{\boldsymbol{\lambda}^T \mathbf{x}} \int_0^\infty dt f_n(t) \exp\left\{-t \sum_i \phi(x_i, y_i)\right\} \\ &= \int_0^\infty dt f_n(t) \int_{\mathbb{R}^n} d\mathbf{x} e^{\boldsymbol{\lambda}^T \mathbf{x}} \exp\left\{-t \sum_i \phi(x_i, y_i)\right\} \\ &= \int_0^\infty dt f_n(t) \prod_i \int_{\mathbb{R}} dx_i e^{\lambda_i x_i} \exp\{-t \phi(x_i, y_i)\}. \end{aligned} \quad (29)$$

Before proceeding, we should note that by using the Laplace transform, we have essentially represented the joint distribution of \mathbf{X} and \mathbf{Y} as a mixture of product form measures, indexed by t , each being proportional to $\exp\{-t \sum_i \phi(x_i, y_i)\}$. If we normalize these measures by $Z_t^n = [\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \exp\{-t \phi(x, y)\}]^n$, and define the i.i.d. probability distribution

$$P(\mathbf{x}, \mathbf{y}|t) = \frac{\exp\{-t \sum_i \phi(x_i, y_i)\}}{Z_t^n}$$

then $P(\mathbf{x}, \mathbf{y})$ is essentially expressed here as a mixture of i.i.d. probability functions $\{P(\mathbf{x}, \mathbf{y}|t)\}$, where t can be thought of as a random parameter whose prior is given by $w_n(t) = f_n(t)Z_t^n$. However, it should be kept in mind that this integral representation goes somewhat further than being a mixture of i.i.d. distributions because $f_n(t)$, and hence also $w_n(t)$, may be negative for some ranges of t even when $F_n(s)$ is strictly positive for all s . For example, recall that the inverse Laplace transform of $s^2/(s^2 + \alpha^2)$ is $\sin(\alpha t)$ ($t \geq 0$), and so, for $F_n(s) = \alpha^n/(s^2 + \alpha^2)^n$, $f_n(t)$ is

given by the n -fold convolution of $\sin(\alpha t)$ with itself. In such cases, $P(\mathbf{x}, \mathbf{y})$ cannot be considered a mixture of i.i.d. distributions.

Let us now denote

$$\begin{aligned}\rho(\lambda, \mathbf{y}, t) &= \ln \left[\int_{-\infty}^{\infty} dx e^{\lambda x - t\phi(x, \mathbf{y})} \right], \\ \rho_0(\mathbf{y}, t) &= \rho(0, \mathbf{y}, t) = \ln \left[\int_{-\infty}^{\infty} dx e^{-t\phi(x, \mathbf{y})} \right],\end{aligned}$$

and

$$\zeta(\mathbf{y}, t) = \left. \frac{\partial \rho(\lambda, \mathbf{y}, t)}{\partial \lambda} \right|_{\lambda=0} = \frac{\int_{\mathbb{R}} dx \cdot x e^{-t\phi(x, \mathbf{y})}}{\int_{\mathbb{R}} dx \cdot e^{-t\phi(x, \mathbf{y})}}.$$

Then,

$$Z(\mathbf{y}, \boldsymbol{\lambda}) = \int_0^{\infty} dt f_n(t) e^{\sum_i \rho(\lambda_i, \mathbf{y}_i, t)},$$

and so,

$$\mathbf{E}\{X_i | \mathbf{y}\} = \left. \frac{\partial \ln Z(\mathbf{y}, \boldsymbol{\lambda})}{\partial \lambda_i} \right|_{\boldsymbol{\lambda}=0} = \frac{\int_0^{\infty} dt f_n(t) \zeta(\mathbf{y}_i, t) e^{\sum_i \rho_0(\mathbf{y}_i, t)}}{\int_0^{\infty} dt f_n(t) e^{\sum_i \rho_0(\mathbf{y}_i, t)}}$$

which is approximated by $\zeta(\mathbf{y}_i, \hat{t})$, where \hat{t} is the maximizer of the expression

$$\ln |f_n(t)| + \sum_i \rho_0(\mathbf{y}_i, t).$$

The MMSE of estimating X_i is given by

$$\text{mmse}(X_i | \mathbf{Y}) \approx \mathbf{E}\{X_i^2\} - \mathbf{E}\{\zeta^2(Y_i, t_0(t))\}$$

where the second term is computed as follows:

$$\mathbf{E}\{\zeta^2(Y_i, t_0(t))\} = \int_0^{\infty} dt w_n(t) \mathbf{E}\{\zeta^2(Y_i, t_0(t)) | t\}$$

with the inner expectation being

$$\mathbf{E}\{\zeta^2(Y_i, t_0(t)) | t\} = \frac{\int_{\mathbb{R}} dy e^{\rho_0(\mathbf{y}, t)} \zeta^2(\mathbf{y}, t_0(t))}{\int_{\mathbb{R}} dy e^{\rho_0(\mathbf{y}, t)}}$$

and with $t_0(t)$ being the value of t' that maximizes

$$\left[\ln |f_n(t')| + n \cdot \frac{\int_{\mathbb{R}} dy \cdot e^{\rho_0(\mathbf{y}, t')} \rho(\mathbf{y}, t')}{\int_{\mathbb{R}} dy \cdot e^{\rho_0(\mathbf{y}, t')}} \right].$$

Thus, we have characterized both the conditional mean estimator and the MMSE in the spirit of a single-letter formula for this class of joint distributions.

The following further extensions of this formalism are conceptually straightforward:

1. The range of the variable t may not necessarily be $[0, \infty)$. Our above analysis applies to whatever range as long as the integrals exist.
2. The joint distribution $P(\mathbf{x}, \mathbf{y})$ may be a function of more than one statistic $\sum_i \phi(x_i, y_i)$, i.e.,

$$P(\mathbf{x}, \mathbf{y}) = F_n \left(\sum_{i=1}^n \phi_1(x_i, y_i), \dots, \sum_{i=1}^n \phi_k(x_i, y_i) \right).$$

In this case, one may apply a Laplace transform of a higher dimension

$$F(s_1, \dots, s_k) = \int_0^\infty \dots \int_0^\infty dt_1 \dots dt_k f(t_1, \dots, t_k) e^{-s_1 t_1 - \dots - s_k t_k},$$

where $s_i = \sum_i \phi_i(x_i)$, $i = 1, 2, \dots, k$.

3. The assumption that the i -th term of $\sum_i \phi(x_i, y_i)$ depends only on the i -th coordinate of \mathbf{y} is not really necessary. The derivation continues to hold, for example, if we allow more generally the form $\sum_i \phi(x_i, y_i, y_{i-1}, \dots, y_{i-k})$.
4. The case where ϕ is a quadratic form can be extended to allow a quadratic form that involves all coordinates of \mathbf{x} and \mathbf{y} collectively, using a positive definite matrix S for weighting. In other words, joint distributions with elliptic symmetry are allowed, with the form $P(\mathbf{x}, \mathbf{y}) = F_n[(\mathbf{x}, \mathbf{y})^T S(\mathbf{x}, \mathbf{y})]$, where (\mathbf{x}, \mathbf{y}) denotes the concatenated column vector of dimension $(n+m)$ formed by \mathbf{x} and \mathbf{y} , and the matrix S is of dimension $(n+m) \times (n+m)$. In this case, the kernel is Gaussian and hence the estimator is linear for a given t .

6 Conclusion

In this paper, we have proposed a simple relation between MMSE estimation measures and a certain expression, which can be viewed as a partition function, and hence be analyzed using methods of statistical mechanics. This partition function is also related to several information measures, like the information density and the Fisher information. The proposed approach has several advantages over the I-MMSE relation and its variants:

1. It is conceptually simple and direct.
2. It applies in full generality, for every joint distribution of the desired random vector \mathbf{X} and its noisy observation vector \mathbf{Y} .

3. It provides, not only the MMSE error covariance matrix, but also the conditional mean estimator itself $\hat{\mathbf{x}} = \mathbf{E}\{\mathbf{X}|\mathbf{y}\}$.
4. It offers several alternative expressions of the MMSE (see Proposition 1).
5. The approach is easy to extend to the mismatched case and it allows mismatch, not only in the marginal of \mathbf{X} , but in the entire joint density $P(\mathbf{x}, \mathbf{y})$.

Finally, considering earlier work on the I-MMSE relation and its various variants that were discussed in the Introduction, it would be natural to seek relations between MMSE estimation to the Hessian of the mutual information. One can show, using the same techniques as in Subsection 3.1, that the following relation holds:

$$E = \nabla_0^2 I_\lambda(\mathbf{X}; \mathbf{Y}) + \text{Cov}\{\mathbf{X}\} - \text{Cov}\left\{(\mathbf{X} - \mathbf{E}\{\mathbf{X}\})(\mathbf{X} - \mathbf{E}\{\mathbf{X}\})^T, \ln \frac{P(\mathbf{Y}|\mathbf{X})}{Z(\mathbf{Y}, \lambda)}\right\},$$

where $I_\lambda(\mathbf{X}; \mathbf{Y})$ is the mutual information induced by the joint distribution

$$P_\lambda(\mathbf{x}, \mathbf{y}) = \frac{e^{\boldsymbol{\lambda}^T \mathbf{x}} P(\mathbf{x}, \mathbf{y})}{\Theta(\boldsymbol{\lambda})}.$$

Unfortunately, this relation seems somewhat more complicated and not as useful as the I-MMSE relation of [5] or the relations proposed in Subsection 3.1 herein.

References

- [1] R. S. Bucy, “Information and filtering,” *Information Sciences*, vol. 18, pp. 179–187, 1979.
- [2] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, John Wiley & Sons, Hoboken, NJ, U.S.A., 2006.
- [3] T. E. Duncan, “On the calculation of mutual information,” *SIAM Journal on Applied Mathematics*, vol. 19, no. 1, pp. 215–220, 1970.
- [4] D. Guo, “Relative entropy and score function: new information–estimation relationships through arbitrary additive perturbations,” *Proc. ISIT 2009*, Seoul, South Korea, June–July 2009.
- [5] D. Guo, S. Shamai, and S. Verdú, “Mutual information and minimum mean–square error in Gaussian channels,” *IEEE Trans. Inform. Theory*, vol. 51, no. 4, pp. 1261–1282, April 2005.
- [6] D. Guo, S. Shamai, and S. Verdú, “Additive non–Gaussian noise channels: mutual information and conditional mean estimation,” *Proc. 2005 IEEE Symp. on Inform. Theory (SIT 2005)*, pp. 719–723, Adelaide, Australia, September 2005.
- [7] D. Guo, S. Shamai, and S. Verdú, “Mutual information and conditional mean estimation in Poisson channels,” *IEEE Trans. Inform. Theory*, vol. 54, no. 5, pp. 1187–1849, May 2008.
- [8] T. Kailath, “The innovations approach to detection and estimation theory,” *Proc. of the IEEE*, vol. 58, no. 5, pp. 680–695, May 1970.
- [9] N. Merhav, D. Guo, and S. Shamai (Shitz), “Statistical physics of signal estimation in Gaussian noise: theory and examples of phase transitions,” to appear in *IEEE Trans. Inform. Theory*, March 2010.
- [10] M. Mézard and A. Montanari, *Information, Physics, and Computation*, Oxford University Press, 2009.
- [11] D. P. Palomar and S. Verdú, “Gradient of mutual information in linear vector Gaussian channels,” *IEEE Trans. Inform. Theory*, vol. 52, no. 1, pp. 141–154, January 2006.

- [12] M. Raginsky and T. P. Coleman, “Mutual information and posterior estimates in channels of exponential family type,” *Proc. 2009 IEEE Workshop on Inform. Theory*, pp. 399–403, Taormina, Sicily, October 2009.
- [13] S. Verdú, “Mismatched estimation and relative entropy,” *Proc. ISIT 2009*, Seoul, South Korea, June–July 2009.
- [14] S. Verdú and T. S. Han, “A general formula for channel capacity,” *IEEE Trans. Inform. Theory*, vol. IT-40, no. 4, pp. 1147–1157, July 1994.