

A General Derivative Identity for the Conditional Mean Estimator in Gaussian Noise and Some Applications

Alex Dytso*, H. Vincent Poor†, and Shlomo Shamai (Shitz)‡

* †Department of Electrical Engineering, Princeton University

‡Department of Electrical Engineering, Technion – Israel Institute of Technology

Email: adytso@princeton.edu*, poor@princeton.edu† and sshlomo@ee.technion.ac.il‡

Abstract—This paper provides a general derivative identity for the conditional mean estimator of an arbitrary vector signal in Gaussian noise with an arbitrary covariance matrix. This new identity is used to recover and generalize many known identities in the literature and derive some new identities. For example, a new identity is discovered, which shows that an arbitrary higher-order conditional moment is completely determined by the first conditional moment.

Several applications of the identities are shown. For instance, by using one of the identities, a simple proof of the uniqueness of the conditional mean estimator as a function of the distribution of the signal is shown. Moreover, one of the identities is used to extend the notion of empirical Bayes to higher-order conditional moments. Specifically, based on a random sample of noisy observations, a consistent estimator for a conditional expectation of any order is derived.

Index Terms—Conditional mean estimator, empirical Bayes, Gaussian Noise.

A full version of this paper is accessible at [1].

I. INTRODUCTION

There are several derivative identities in the literature that relate conditional mean estimators to other quantities such as the score function and the conditional variance. Such identities are often used in information theory to give way to estimation theoretic arguments (e.g., I-MMSE relationship [2]). In estimation theory such identities are often used to design new estimation procedures (e.g., empirical Bayes [3]). In this work, it is shown that many of the known identities in the literature can be derived systematically from a single unifying derivative identity. Moreover, we use this new identity to derive several generalizations of the previously known identities and discover some new identities. Furthermore, the derived identities are used to propose a generalization of an empirical Bayes procedure.

Contribution: The contribution and the outline of the paper are as follows:

- Section II presents the system model;
- In Section III, Theorem 1 presents a new identity for the Jacobian of the conditional mean;

This work was supported in part by the U. S. National Science Foundation under Grant CCF-1908308 and by the United States-Israel Binational Science Foundation, under Grant BSF-2018710.

- Section IV is dedicated to the Tweedie-Robbins-Esposito (TRE) formula. The TRE formula relates the conditional expectation with the score function of the output random variable;
- In Section V, Proposition 1 presents a simple proof of a vector version of Hatsel-Nolte identity, which relates the Jacobian of the conditional expectation to the conditional variance;
- In Section VI, Proposition 2 and Proposition 3 show two vector generalizations of the recursive Jaffer's identity, which relates higher-order conditional moments to the derivatives of the lower-order conditional moments. Moreover, in Section VI-B, Proposition 4 provides an equivalent integral version of Jaffer's identity which shows that every higher-order conditional moment is completely determined by the first-order conditional moment;
- Section VII is dedicated to applications of the derived identities. In Section VII-A, Proposition VII-A shows a small application of Hatsel-Nolte identity concerning the maximum and minimum ‘slope’ of the conditional mean estimator. In Section VII-B, Theorem 2 uses the new integral generalization of Jaffer's identity to show the uniqueness of the conditional mean estimator as a function of the distribution of the signal. Finally, Section VII-C uses a new integral generalization of Jaffer's identity to extend the notion of empirical Bayes to higher-order conditional moments. Specifically, Theorem 3, based on a random sample of noisy observations, proposes a consistent estimator for a conditional expectation of any order; and
- Section VIII concludes the paper.

Notation: The set of all positive integers is denoted by \mathbb{N} , $[n]$ is the set of integers $\{1, \dots, n\}$, and \mathbb{R}^n is the set of all n -dimensional real-valued vectors.

Deterministic scalar quantities are denoted by lowercase letters, scalar random variables are denoted by uppercase letters, vectors are denoted by bold lowercase letters, random vectors by bold uppercase letters, and matrices by bold uppercase sans serif letters (e.g., x , X , \mathbf{x} , \mathbf{X} , \mathbf{X}). All vectors in the paper are

column vectors.

The standard basis vectors for \mathbb{R}^n are denoted by \mathbf{e}_i , $i \in [n]$. For a matrix \mathbf{A} , we use $[\mathbf{A}]_{ij}$ to denote the entry of the row i and column j . The Euclidian norm of a vector $\mathbf{x} \in \mathbb{R}^n$ in this paper is denoted by $\|\mathbf{x}\|$.

The gradient of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is denoted by $\nabla_{\mathbf{x}} f(\mathbf{x}) \in \mathbb{R}^n$, and the Jacobian matrix of a function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is denoted by $\mathbf{J}_{\mathbf{x}} \mathbf{f}(\mathbf{x}) \in \mathbb{R}^{m \times n}$, that is

$$[\mathbf{J}_{\mathbf{x}} \mathbf{f}(\mathbf{x})]_{ij} = \frac{\partial f_i(\mathbf{x})}{\partial x_j}, i \in [m], j \in [n]. \quad (1)$$

The Hadamard product between $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{B} \in \mathbb{R}^{n \times m}$ will be denoted by $\mathbf{A} \odot \mathbf{B}$. Moreover, $\mathbf{A}^{\odot n}$, $n \in \mathbb{N}$ denotes the Hadamard product repeated n times on the matrix \mathbf{A} . For a function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, we use the following version of the higher-order gradient $\nabla^{\odot v} \mathbf{f}$, $v \in \mathbb{N}$, which is defined as

$$[\nabla^{\odot v} \mathbf{f}(\mathbf{x})]_i = \frac{\partial^v}{\partial x_i^v} [\mathbf{f}(\mathbf{x})]_i, i \in [n]. \quad (2)$$

II. MODEL

The underlying model considered through this paper is described by the following input-output relationship:

$$\mathbf{Y} = \mathbf{X} + \mathbf{N}, \quad (3)$$

where $\mathbf{N} \in \mathbb{R}^n$ is a zero mean, normally distributed with the covariance matrix $\mathbf{K}_{\mathbf{N}}$, and independent of $\mathbf{X} \in \mathbb{R}^n$. Throughout the paper $\mathbf{K}_{\mathbf{N}}$ is assumed to be a positive definite matrix, and we make no assumptions about the probability distribution of \mathbf{X} .

III. A NEW IDENTITY FOR THE CONDITIONAL EXPECTATION

The main result of this paper is the following theorem.

Theorem 1. Let $\mathbf{U} \in \mathbb{R}^m$ be an arbitrary function of $\mathbf{X} \in \mathbb{R}^n$. Then,

$$\mathbf{J}_{\mathbf{y}} \mathbb{E}[\mathbf{U} | \mathbf{Y} = \mathbf{y}] = \mathbf{Cov}(\mathbf{U}, \mathbf{X} | \mathbf{Y} = \mathbf{y}) \mathbf{K}_{\mathbf{N}}^{-1}, \mathbf{y} \in \mathbb{R}^n. \quad (4)$$

Proof: See Appendix. ■

Conversely, Theorem 1 can be re-written as

$$\begin{aligned} \mathbb{E}[U_i | \mathbf{Y} = \mathbf{y}] &= \mathbb{E}[U_i | \mathbf{Y} = \mathbf{0}] + \oint_0^{\mathbf{y}} \mathbf{Cov}(U_i, \mathbf{X} | \mathbf{Y} = \mathbf{t}) \mathbf{K}_{\mathbf{N}}^{-1} \cdot d\mathbf{t}, \end{aligned} \quad (5)$$

where $i \in [m]$, and \oint is a line integral over an arbitrary path between $\mathbf{0}$ and \mathbf{y} .

We now discuss some consequences of the identity in Theorem 1. Specifically, this would be done by evaluating Theorem 1 with different choices of \mathbf{U} such as $\mathbf{U} = \mathbf{X}$, and $\mathbf{U} = (\mathbf{X} \mathbf{X}^T)^{k-1} \mathbf{X}$, $k \in \mathbb{N}$.

IV. TWEEDIE-ROBBINS-ESPOSITO IDENTITY

The proof of Theorem 1 will rely on the following identity:

$$\mathbb{E}[\mathbf{X} | \mathbf{Y} = \mathbf{y}] = \mathbf{y} + \mathbf{K}_{\mathbf{N}} \frac{\nabla_{\mathbf{y}} f_{\mathbf{Y}}(\mathbf{y})}{f_{\mathbf{Y}}(\mathbf{y})}, \quad (6)$$

where $f_{\mathbf{Y}}(\mathbf{y})$ is the probability density function (pdf) of \mathbf{Y} . We note that the quantity $\frac{\nabla_{\mathbf{y}} f_{\mathbf{Y}}(\mathbf{y})}{f_{\mathbf{Y}}(\mathbf{y})}$ is commonly known as the score function.

Literature Review: The scalar version of the identity in (6) has been derived by Robbins in [4] where he credits Maurice Tweedie for the derivation. The vector version of the identity in (6) was derived by Esposito in [5]. Therefore, throughout this paper, we refer to the identity in (6) as Tweedie-Robbins-Esposito identity or TRE for short.

The application of the TRE identity in (6) can considerably simplify the computation of $\mathbb{E}[\mathbf{X} | \mathbf{Y}]$ as we do not need to derive the conditional distribution $P_{\mathbf{X} | \mathbf{Y}}$ and only require to compute $f_{\mathbf{Y}}(\mathbf{y})$ and the gradient of $f_{\mathbf{Y}}(\mathbf{y})$. For an example of such an application, the interested reader is referred to [6] where the TRE identity was used to compute $\mathbb{E}[\mathbf{X} | \mathbf{Y}]$ for the case where \mathbf{X} is uniform on a sphere.

The observation that, via the TRE identity, the conditional expectation can be represented only in terms of the marginal distribution of the output \mathbf{Y} has led to the development of the empirical Bayes procedure [4]; the interested reader is referred to [3] for an overview of the empirical Bayes procedure.

The TRE identity has also been used in the proofs of the scalar and vector versions of the I-MMSE relationship in [7] and [8], respectively.

V. THE IDENTITY OF HATSELL AND NOLTE

By setting $\mathbf{U} = \mathbf{X}$ in (4) we arrive at the following identity.

Proposition 1. For $\mathbf{y} \in \mathbb{R}^n$.

$$\mathbf{J}_{\mathbf{y}} \mathbb{E}[\mathbf{X} | \mathbf{Y} = \mathbf{y}] = \mathbf{Var}(\mathbf{X} | \mathbf{Y} = \mathbf{y}) \mathbf{K}_{\mathbf{N}}^{-1}. \quad (7)$$

Literature Review: The identity in (7) has been first derived by Hatsell and Nolte in [9] for the case of $\mathbf{K}_{\mathbf{N}} = \mathbf{I}$. The general version in (7) was first derived in [8] where it was used, together with the TRE identity in (6), to give a proof of the vector version of the I-MMSE relationship.

In [10], the scalar version of the identity in (7), was used to show that the minimum mean squared error is Lipschitz continuous with respect to the Wasserstein distance.

VI. JAFFER'S IDENTITY

In [11], Jaffer has shown the following identity: for $k \in \mathbb{N} \cup \{0\}$

$$\begin{aligned} \mathbb{E}[X^{k+1} | Y = y] &= \sigma^2 \frac{d}{dy} \mathbb{E}[X^k | Y = y] + \mathbb{E}[X^k | Y = y] \mathbb{E}[X | Y = y], \end{aligned} \quad (8)$$

where σ^2 is the variance of the noise. To the best of our knowledge, Jaffer's identity in (8) has had little applications and is not well known. In what follows, we develop several vector generalizations of Jaffer's identity. Moreover, we derive an alternative but equivalent integral version of the Jaffer's identity and show how this new identity can be used to prove uniqueness of the conditional mean estimator. We also use this integral identity to extend the notion of empirical Bayes to higher order conditional moments.

A. Vector Generalization of Jaffer's Identity

Given the fact that there is no unique generalization of higher moments to the vector case, several vector generalizations of the identity in (8) are possible. Next, we present two such generalizations.

The first generalization of (8) allows to have different exponents across elements of \mathbf{X} .

Proposition 2. For every $m \in [n]$, $v_i \in \mathbb{N} \cup \{0\}$, $i \in [n]$ and $\mathbf{y} \in \mathbb{R}^n$

$$\begin{aligned} & \frac{d}{dy_m} \mathbb{E} \left[\prod_{i=1}^n (\mathbf{e}_i^\top \mathbf{K}_N^{-1} \mathbf{X})^{v_i} \mid \mathbf{Y} = \mathbf{y} \right] \\ &= \mathbb{E} \left[\prod_{i=1: i \neq m}^n (\mathbf{e}_i^\top \mathbf{K}_N^{-1} \mathbf{X})^{v_i} (\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X})^{v_{m+1}} \mid \mathbf{Y} = \mathbf{y} \right] \\ & \quad - \mathbb{E} \left[\prod_{i=1}^n (\mathbf{e}_i^\top \mathbf{K}_N^{-1} \mathbf{X})^{v_i} \mid \mathbf{Y} = \mathbf{y} \right] \mathbb{E} [\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{y}]. \quad (9) \end{aligned}$$

Proof: The proof follows by evaluating (4) with $\mathbf{U} = \prod_{i=1}^n (\mathbf{e}_i^\top \mathbf{K}_N^{-1} \mathbf{X})^{v_i}$. \blacksquare

In the case when \mathbf{K}_N is a diagonal matrix with $[\mathbf{K}_N]_{ii} = \sigma_i^2$ the identity (9) reduces to

$$\begin{aligned} & \mathbb{E} \left[X_m^{v_{m+1}} \prod_{i=1: i \neq m}^n X_i^{v_i} \mid \mathbf{Y} = \mathbf{y} \right] = \sigma_i^2 \frac{d}{dy_i} \mathbb{E} \left[\prod_{i=1}^n X_i^{v_i} \mid \mathbf{Y} = \mathbf{y} \right] \\ & \quad + \mathbb{E} \left[\prod_{i=1}^n X_i^{v_i} \mid \mathbf{Y} = \mathbf{y} \right] \mathbb{E} [X_m \mid \mathbf{Y} = \mathbf{y}]. \quad (10) \end{aligned}$$

Furthermore, by setting $v_i = 0$ for $i \neq m$ the identity in (9) can be re-written in terms of the Hadamard product as follows: for $\mathbf{y} \in \mathbb{R}^n$, $m \in [n]$ and $v \in \mathbb{N} \cup \{0\}$

$$\begin{aligned} & \frac{d}{dy_m} \mathbb{E} [(\mathbf{K}_N^{-1} \mathbf{X})^{\odot v} \mid \mathbf{Y} = \mathbf{y}] = \mathbb{E} [(\mathbf{K}_N^{-1} \mathbf{X})^{\odot(v+1)} \mid \mathbf{Y} = \mathbf{y}] \\ & \quad - \mathbb{E} [(\mathbf{K}_N^{-1} \mathbf{X})^{\odot v} \mid \mathbf{Y} = \mathbf{y}] \odot \mathbb{E} [\mathbf{K}_N^{-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{y}]. \quad (11) \end{aligned}$$

The second generalization of (8) is in terms of powers of a matrix.

Proposition 3. For $k \in \mathbb{N}$ and $\mathbf{y} \in \mathbb{R}^n$

$$\begin{aligned} & \mathbb{E} [(\mathbf{X} \mathbf{X}^\top)^k \mid \mathbf{Y} = \mathbf{y}] = \mathbf{K}_N \mathbf{J}_y \mathbb{E} [(\mathbf{X} \mathbf{X}^\top)^{k-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{y}] \\ & \quad + \mathbb{E} [(\mathbf{X} \mathbf{X}^\top)^{k-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{y}] \mathbb{E} [\mathbf{X}^\top \mid \mathbf{Y} = \mathbf{y}]. \quad (12) \end{aligned}$$

Proof: The proof follows by evaluating (4) with $\mathbf{U} = (\mathbf{X} \mathbf{X}^\top)^{k-1} \mathbf{X}$. \blacksquare

B. A New Perspective on Jaffer's Identity

Next, we show that Jaffer's identity has an alternative integral version. This new integral representation leads to several interesting consequences. The following auxiliary lemma, the proof of which can be found in [1], would be useful.

Lemma 1. Let $f_k : \mathbb{R} \rightarrow \mathbb{R}$ be a sequence of functions $k = 0, 1, 2, \dots$ with $f_0 \equiv 1$. Then, the recurrence relationship

$$f_k(x) = \frac{d}{dx} f_{k-1}(x) + f_{k-1}(x) f_1(x), k = 1, 2, \dots \quad (13)$$

is equivalent to

$$f_k(x) = e^{- \int_0^x f_1(t) dt} \frac{d^k}{dx^k} e^{\int_0^x f_1(t) dt}. \quad (14)$$

Next, we present the integral alternative of Jaffer's identity.

Proposition 4. Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\mathbf{f}^{-\odot} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be defined as follows: for $i \in [n]$ and $\mathbf{y} \in \mathbb{R}^n$

$$[\mathbf{f}(\mathbf{y})]_i = e^{\int_0^{y_i} \mathbb{E} [\mathbf{e}_i^\top \mathbf{K}_N^{-1} \mathbf{X} \mid \mathbf{Y} = (y_1, \dots, y_{i-1}, t, y_{i+1}, \dots, y_n)] dt}, \quad (15)$$

$$[\mathbf{f}^{-\odot}(\mathbf{y})]_i = [\mathbf{f}(\mathbf{y})]_i^{-1}. \quad (16)$$

Then, for $k \in \mathbb{N}$ and $\mathbf{y} \in \mathbb{R}^n$

$$\mathbb{E} [(\mathbf{K}_N^{-1} \mathbf{X})^{\odot k} \mid \mathbf{Y} = \mathbf{y}] = \mathbf{f}^{-\odot}(\mathbf{y}) \odot \nabla^{\odot k} \mathbf{f}(\mathbf{y}). \quad (17)$$

Proof: By setting $v_i = k$ for $i = m$ and $v_i = 0$ for $i \neq m$ the identity in (9) reduces to

$$\begin{aligned} & \frac{d}{dy_m} \mathbb{E} [(\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X})^k \mid \mathbf{Y} = \mathbf{y}] \\ &= \mathbb{E} [(\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X})^{k+1} \mid \mathbf{Y} = \mathbf{y}] \\ & \quad - \mathbb{E} [(\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X})^k \mid \mathbf{Y} = \mathbf{y}] \mathbb{E} [\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{y}]. \quad (18) \end{aligned}$$

Now, by defining

$$f_k(y_m) = \mathbb{E} [(\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X})^k \mid \mathbf{Y} = \mathbf{y}], \quad (19)$$

and applying Lemma 1 to (18) we arrive at

$$\begin{aligned} & \mathbb{E} [(\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X})^k \mid \mathbf{Y} = \mathbf{y}] \\ &= e^{- \int_0^{y_m} \mathbb{E} [\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{t}] dt_m} \frac{d^k}{dy_m^k} e^{\int_0^{y_m} \mathbb{E} [\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} \mid \mathbf{Y} = \mathbf{t}] dt_m} \quad (20) \end{aligned}$$

where $[\mathbf{t}]_i = y_i$ for $i \neq m$. The proof is concluded by using the definition of the Hadamard product and $\nabla^{\odot v}$ in (2). \blacksquare

In the case when \mathbf{K}_N is a diagonal matrix with $[\mathbf{K}_N]_{ii} = \sigma_i^2$ the identity (17) reduces to: for $y \in \mathbb{R}$ and $k \in \mathbb{N}$

$$\begin{aligned} & \mathbb{E} [X_m^k \mid Y_m = y_m] \\ &= \sigma_i^{2k} e^{-\frac{1}{\sigma_m^2} \int_0^{y_m} \mathbb{E} [X_m \mid Y_m = t] dt} \frac{d^k}{dy_m^k} e^{\frac{1}{\sigma_m^2} \int_0^{y_m} \mathbb{E} [X_m \mid Y_m = t] dt}. \quad (21) \end{aligned}$$

An important feature of the integral version of Jaffer's identity is that any higher-order conditional moment is determined by the first conditional moment. This observation would be used in Section VII-B to show that the conditional expectation is uniquely determined by the distribution of the input X .

Another identity equivalent to that in (21) is shown in the following corollary.

Corollary 1. For any $k \in \mathbb{N}$ and $y \in \mathbb{R}$

$$\mathbb{E} [X^k \mid Y = y] = \sigma^{2k} \frac{\frac{d^k}{dy^k} \left(f_Y(y) e^{\frac{y^2}{2\sigma^2}} \right)}{f_Y(y) e^{\frac{y^2}{2\sigma^2}}}. \quad (22)$$

Alternatively, let $t \mapsto H_{e_m}(t)$, $m \in [0, 1, \dots]$ be a probabilistic Hermite polynomial, then

$$\mathbb{E}[X^k | Y = y] = \sigma^{2k} \frac{\sum_{m=0}^k \binom{k}{m} f_Y^{(k-m)}(y) \frac{(-i)^m}{\sigma^m} H_{e_m}(i \frac{y}{\sigma})}{f_Y(y)}. \quad (23)$$

Proof: First, observe that using the scalar TRE identity in (6) we have that

$$\int_0^y \frac{\mathbb{E}[X|Y=t]}{\sigma^2} dt = \int_0^y \left(\frac{t}{\sigma^2} + \frac{d}{dt} \log(f_Y(t)) \right) dt \quad (24)$$

$$= \frac{y^2}{2\sigma^2} + \log(f_Y(y)) - \log(f_Y(0)). \quad (25)$$

Inserting (25) into (21) concludes the proof of (22). The proof of (23) follows from the generalized product rule. ■

The identity in (22) can be thought of as a generalization of the TRE identity in (6) to the higher-order moments. Similarly, to the TRE identity the important feature of the identity in (22) is that $\mathbb{E}[X^k | Y]$ depends on the joint distribution $P_{X,Y}$ only through the marginal pdf of Y . In Section VII-C, the identity in (22) will be used to extended the empirical Bayes procedure to higher-order conditional moments.

VII. APPLICATIONS

In this section, we demonstrate three applications of the above identities. For ease of exposition, we mostly focus on the scalar case.

A. A Small Application of the Identity in (7)

The identity of Hatsell and Nolte in (7) can be used to make various statements about the minimum and maximum ‘slope’ of $\mathbb{E}[\mathbf{X}|\mathbf{Y}]$.

Corollary 2. For every $\mathbf{y} \in \mathbb{R}^n$

$$0 \leq \text{Tr}(\mathbf{J}_y \mathbb{E}[\mathbf{X}|\mathbf{Y} = \mathbf{y}]). \quad (26)$$

In addition, if $\|\mathbf{X}\| \leq R$, then

$$\text{Tr}(\mathbf{J}_y \mathbb{E}[\mathbf{X}|\mathbf{Y} = \mathbf{y}]) \leq R^2 \text{Tr}(\mathbf{K}_N^{-1}). \quad (27)$$

Proof: The proof of the lower bound follows by using (7) together with the properties that both variances are positive definite matrices and that the trace of product of two positive definite matrices is positive. To show the upper bound in (27) we use (7) together with the Cauchy-Schwarz inequality

$$\text{Tr}(\mathbf{J}_y \mathbb{E}[\mathbf{X}|\mathbf{Y} = \mathbf{y}]) = \text{Tr}(\mathbf{Var}(\mathbf{X}|\mathbf{Y} = \mathbf{y}) \mathbf{K}_N^{-1}) \quad (28)$$

$$\leq R^2 \text{Tr}(\mathbf{K}_N^{-1}). \quad (29)$$

■

B. Uniqueness of the Conditional Expectation via Integral Version of Jaffer’s Identity

The identity in (21) leads to an interesting observation that all of the higher conditional moments are completely determined by the first conditional moment. In the next result, we use this observation to establish the uniqueness of the conditional expectation as a function of the distribution of X .

Theorem 2. The conditional expectation $\mathbb{E}[X|Y]$ is a bijective operator of the distribution of X . In other words, let P_X be the distribution of X , then

$$\mathbb{E}[X_1|Y_1 = y] = \mathbb{E}[X_2|Y_2 = y], \forall y \in \mathbb{R} \iff P_{X_1} = P_{X_2}. \quad (30)$$

Proof: Let $P_{X_1} = P_{X_2}$, then it is immediate that

$$\mathbb{E}[X_1|Y_1 = y] = \mathbb{E}[X_2|Y_2 = y], \forall y \in \mathbb{R}. \quad (31)$$

Now suppose that $\mathbb{E}[X_1|Y_1 = y] = \mathbb{E}[X_2|Y_2 = y], \forall y \in \mathbb{R}$, and the goal is to show that $P_{X_1} = P_{X_2}$. First, observe that by the identity in (21) we have that for all $k \in \mathbb{N}$ and all $y \in \mathbb{R}$

$$\mathbb{E}[X_1^k|Y_1 = y] = \mathbb{E}[X_2^k|Y_2 = y]. \quad (32)$$

We now use (32) to establish that $P_{X_1|Y_1=y} = P_{X_2|Y_2=y}$ for all $y \in \mathbb{R}$. First, without loss of generality assume that $\sigma^2 = 1$. Second, fix some $y \in \mathbb{R}$ in (32) and let

$$m_{1,k} = \mathbb{E}[X_1^k|Y_1 = y], m_{2,k} = \mathbb{E}[X_2^k|Y_1 = y]. \quad (33)$$

The identity (32) implies that $m_{1,k} = m_{2,k}, \forall k \in \mathbb{N}$. The question of whether a distribution of a real valued random variables is determined by its moment is known as Hamburger moment problem [12]. We now use Carleman’s sufficient condition to check whether the moments uniquely determine the distribution, this requires verifying that the following sum of moments is divergent:

$$\sum_{k=1}^{\infty} m_{i,2k}^{-\frac{1}{2k}} = \infty. \quad (34)$$

To show that the left side of (34) diverges we will need the following upper bound on the conditional moments shown in [13]: for $i \in \{1, 2\}$

$$m_{i,2k} \leq c_i k 2^{k+1} \sqrt{(2k-1)!}, k \in [n] \quad (35)$$

with $c_i = \frac{e^{\frac{y^2}{2}}}{f_{Y_i}(y)}$. Applying (35) to (34) we have that

$$\sum_{k=1}^{\infty} m_{i,2k}^{-\frac{1}{2k}} \geq \sum_{k=1}^{\infty} \frac{1}{c_i^{\frac{1}{2k}} k^{\frac{1}{2k}} 2^{\frac{k+1}{2k}} ((2k-1)!)^{\frac{1}{4k}}} \quad (36)$$

$$\geq \frac{1}{2c_i^{\frac{1}{2}}} \sum_{k=1}^{\infty} \frac{1}{k^{\frac{1}{2k}} ((2k)!)^{\frac{1}{4k}}}, \quad (37)$$

which diverges by the comparison test. Therefore, the Carleman condition in (34) is satisfied, and the moments determine the distribution. In other words, we have demonstrated that (32) implies that for all $y \in \mathbb{R}$

$$P_{X_1|Y_1=y} = P_{X_2|Y_2=y}. \quad (38)$$

Now the equality in (38) implies that $P_{X_1} = P_{X_2}$. To see this choose some measurable set $A \subseteq \mathbb{R}$ and observe that

$$P_{X_1}(A) = \mathbb{E}[1_A(X_1)] \quad (39)$$

$$= \mathbb{E}[\mathbb{E}[1_A(X_1)|Y_1]] \quad (40)$$

$$= \mathbb{E}[\mathbb{E}[1_A(X_2)|Y_2]] = P_{X_2}(A). \quad (41)$$

This concludes the proof. ■

Theorem 2 has been previously shown in [10]. However, our proof here is different than that in [10], and our method is more akin to the proof of the uniqueness of the conditional expectation for the Poisson noise used in [14].

C. Empirical Bayes for Higher-Order Conditional Moments

An interesting application of the original TRE identity is the idea of empirical Bayes proposed by Robbins in [4]. Consider an independent and identically distributed sequence Y_1, \dots, Y_n according to f_Y , and assume that we have a perfect knowledge of σ . Because the conditional estimator in the TRE formula depends only on the marginal distribution of the output Y , from the Y_i observations, we can build an empirical estimate of $\mathbb{E}[X|Y = y]$ by mimicking the TRE formula

$$\hat{m}(y) = y + \sigma^2 \frac{\hat{f}'_Y(y)}{\hat{f}_Y(y)}, \quad y \in \mathbb{R}, \quad (42)$$

where $\hat{f}_Y(y)$ and $\hat{f}'_Y(y)$ are estimates of $f_Y(y)$ and $f'_Y(y)$, respectively. In other words, we are able to estimate the $\mathbb{E}[X|Y]$ without the knowledge of the prior distribution on X . The interested reader is referred to [3, Chapter 6.1] for a historical account and impact of the empirical Bayes formula.

Let $\hat{f}_Y^{(k)}(y)$ denote the estimate of $f_Y^{(k)}(y)$ based on a random sample Y_1, \dots, Y_n . Now, inspired by (22), define the estimator of $\mathbb{E}[X^k|Y = y]$ as follows:

$$\hat{m}_k(y) = \sigma^{2k} \frac{\sum_{m=0}^k \binom{k}{m} \hat{f}_Y^{(k-m)}(y) \frac{(-i)^m}{\sigma^m} H_{e_m} \left(\frac{i y}{\sigma} \right)}{\hat{f}_Y(y)}, \quad y \in \mathbb{R}. \quad (43)$$

To estimate $f_Y^{(k)}(y)$ we use the following steps. First, estimate f_Y by using a kernel-density estimator

$$\hat{f}_Y(t) = \frac{1}{n} \sum_{i=1}^n \frac{1}{a} k \left(\frac{t - Y_i}{a} \right), \quad (44)$$

where $a > 0$ is the bandwidth parameter. We take the kernel to be $k(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$. Second, estimate $f_Y^{(k)}(t)$ by taking the derivative of (44) k times.

We conclude this section by showing that the estimator in (43) is consistent.

Theorem 3. Let $a = \frac{1}{n^u}$ and $t_n = \frac{\sigma^2 \sqrt{w \log(n)}}{3}$ for some $u \in (0, \frac{1}{2k+6})$ and $w \in (0, u)$. Moreover, assume that $\mathbb{E}[X^2] < \infty$. Then, for every $k \in \mathbb{N}$ and $\sigma^2 > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\sup_{|y| \leq t_n} |\hat{m}_k(y) - \mathbb{E}[X^k|Y = y]| \geq \frac{C_{k,\sigma}}{n^{u-w}} \right] = 0,$$

where $C_{k,\sigma}$ is a constant that depends only on k and σ .

The proof is omitted and can be found in [1].

VIII. CONCLUSION AND OUTLOOK

This work has derived a general derivative identity for a conditional mean estimator. This identity has been used to recover several known derivative identities. Moreover, some generalizations of known identities and some new identities have been derived. For example, a new integral version of Jaffer's identity has been shown and used to prove the uniqueness of the conditional expectation. Moreover, the identities are used to extend the notion of empirical Bayes to higher-order moments.

An interesting future direction would be to see if the main identity in Theorem 1 can shed light on the vector generalization of the single crossing property in [15]. It would also be interesting to see if a vector version of integral Jaffer's identity in Proposition 2 can be used to show the uniqueness of the vector conditional expectation as has been done for the scalar case in Theorem 2.

APPENDIX

First, denote the pdf of \mathbf{N} by $\phi_{\mathbf{K}_N}(\mathbf{y})$ and observe that

$$\begin{aligned} \frac{d}{dy_m} \phi_{\mathbf{K}_N}(\mathbf{y} - \mathbf{X}) \\ = -\frac{1}{2} \phi_{\mathbf{K}_N}(\mathbf{y} - \mathbf{X}) \frac{d}{dy_m} (\mathbf{y} - \mathbf{X})^\top \mathbf{K}_N^{-1} (\mathbf{y} - \mathbf{X}) \end{aligned} \quad (45)$$

$$= \phi_{\mathbf{K}_N}(\mathbf{y} - \mathbf{X}) \mathbf{e}_m^\top \mathbf{K}_N^{-1} (\mathbf{X} - \mathbf{y}). \quad (46)$$

Second, observe the following sequence of steps:

$$\begin{aligned} \frac{d}{dy_m} \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \\ = \frac{d}{dy_m} \mathbb{E} \left[\mathbf{U} \frac{\phi_{\mathbf{K}_N}(\mathbf{y} - \mathbf{X})}{f_{\mathbf{Y}}(\mathbf{y})} \right] \end{aligned} \quad (47)$$

$$= \mathbb{E} \left[\mathbf{U} \frac{d}{dy_m} \frac{\phi_{\mathbf{K}_N}(\mathbf{y} - \mathbf{X})}{f_{\mathbf{Y}}(\mathbf{y})} \right] \quad (48)$$

$$= \mathbb{E} \left[\mathbf{U} \frac{\frac{d}{dy_m} \phi_{\mathbf{K}_N}(\mathbf{y} - \mathbf{X})}{f_{\mathbf{Y}}(\mathbf{y})} \right] - \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \frac{\frac{d}{dy_m} f_{\mathbf{Y}}(\mathbf{y})}{f_{\mathbf{Y}}(\mathbf{y})} \quad (49)$$

$$\begin{aligned} &= \mathbb{E} [\mathbf{U} \mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} | \mathbf{Y} = \mathbf{y}] \\ &- \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{y} - \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \frac{\frac{d}{dy_m} f_{\mathbf{Y}}(\mathbf{y})}{f_{\mathbf{Y}}(\mathbf{y})} \end{aligned} \quad (50)$$

$$\begin{aligned} &= \mathbb{E} [\mathbf{U} \mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} | \mathbf{Y} = \mathbf{y}] \\ &- \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \mathbf{e}_m^\top \left(\mathbf{K}_N^{-1} \mathbf{y} + \frac{\nabla_{\mathbf{y}} f_{\mathbf{Y}}(\mathbf{y})}{f_{\mathbf{Y}}(\mathbf{y})} \right) \end{aligned} \quad (51)$$

$$\begin{aligned} &= \mathbb{E} [\mathbf{U} \mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} | \mathbf{Y} = \mathbf{y}] \\ &- \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \mathbb{E} [\mathbf{e}_m^\top \mathbf{K}_N^{-1} \mathbf{X} | \mathbf{Y} = \mathbf{y}] \end{aligned} \quad (52)$$

$$\begin{aligned} &= \mathbb{E} [\mathbf{U} \mathbf{X}^\top | \mathbf{Y} = \mathbf{y}] \mathbf{K}_N^{-1} \mathbf{e}_m \\ &- \mathbb{E} [\mathbf{U} | \mathbf{Y} = \mathbf{y}] \mathbb{E} [\mathbf{X}^\top | \mathbf{Y} = \mathbf{y}] \mathbf{K}_N^{-1} \mathbf{e}_m \end{aligned} \quad (53)$$

$$= \mathbf{Cov}(\mathbf{U}, \mathbf{X} | \mathbf{Y} = \mathbf{y}) \mathbf{K}_N^{-1} \mathbf{e}_m, \quad (54)$$

where the equalities follow from: (47) using Bayes' theorem; (50) using the expression in (46); and (52) using the TRE identity in (6). This concludes the proof.

REFERENCES

- [1] A. Dytso, H. V. Poor, and S. Shamai (Shitz), “A general derivative identity for the conditional mean estimator in Gaussian noise and some applications,” 2020. [Online]. Available: <http://www.princeton.edu/~7Eadytso/papers/EB.pdf>
- [2] D. Guo, S. Shamai, and S. Verdú, *The Interplay Between Information and Estimation Measures*. now Publishers Incorporated, 2013.
- [3] B. Efron and T. Hastie, *Computer Age Statistical Inference*. Cambridge University Press, 2016, vol. 5.
- [4] H. Robbins, “An empirical Bayes approach to statistics,” in *Proc. Third Berkeley Symp. Math Statist. Probab.*. CiteSeer, 1956.
- [5] R. Esposito, “On a relation between detection and estimation in decision theory,” *Inf. Control*, vol. 12, no. 2, pp. 116–120, February 1968.
- [6] A. Dytso, M. Al, H. V. Poor, and S. Shamai (Shitz), “On the capacity of the peak power constrained vector Gaussian channel: An estimation theoretic perspective,” *IEEE Trans. Inf. Theory*, vol. 65, no. 6, pp. 3907–3921, 2019.
- [7] D. Guo, S. Shamai, and S. Verdú, “Mutual information and minimum mean-square error in Gaussian channels,” *IEEE Trans. Inf. Theory*, vol. 51, no. 4, pp. 1261–1282, 2005.
- [8] D. P. Palomar and S. Verdú, “Gradient of mutual information in linear vector Gaussian channels,” *IEEE Trans. Inf. Theory*, vol. 52, no. 1, p. 141, 2006.
- [9] C. Hatsell and L. Nolte, “Some geometric properties of the likelihood ratio (corresp.),” *IEEE Trans. Inf. Theory*, vol. 17, no. 5, pp. 616–618, 1971.
- [10] Y. Wu and S. Verdú, “Functional properties of minimum mean-square error and mutual information,” *IEEE Trans. Inf. Theory*, vol. 58, no. 3, pp. 1289–1301, 2012.
- [11] A. G. Jaffer, “A note on conditional moments of random signals in Gaussian noise (corresp.),” *IEEE Trans. Inf. Theory*, vol. 18, no. 4, pp. 513–514, 1972.
- [12] J. A. Shohat and J. D. Tamarkin, *The Problem of Moments*. American Mathematical Soc., 1943, no. 1.
- [13] D. Guo, Y. Wu, S. Shamai, and S. Verdú, “Estimation in Gaussian noise: Properties of the minimum mean-square error,” *IEEE Trans. Inf. Theory*, vol. 57, no. 4, pp. 2371–2385, 2011.
- [14] A. Dytso and H. V. Poor, “Estimation in Poisson noise: Properties of the conditional mean estimator,” *IEEE Trans. Inf. Theory*, 2020, to appear.
- [15] R. Bustin, M. Payaró, D. P. Palomar, and S. Shamai, “On MMSE crossing properties and implications in parallel vector Gaussian channels,” *IEEE Trans. Inf. Theory*, vol. 59, no. 2, pp. 818–844, 2013.