



# Gaussian approximations in high dimensional estimation



Vivek S. Borkar<sup>\*</sup>, Raaz Dwivedi<sup>1</sup>, Neeraja Sahasrabudhe<sup>\*</sup>

Department of Electrical Engineering, Indian Institute of Technology, Powai, Mumbai 400076, India

## ARTICLE INFO

### Article history:

Received 21 April 2015

Accepted 3 March 2016

Available online 29 March 2016

### Keywords:

Gaussian approximations

High dimensional estimation

Log-concave distributions

Stochastic sparsity

Random projections

## ABSTRACT

Several estimation techniques assume validity of Gaussian approximations for estimation purposes. Interestingly, these ensemble methods have proven to work very well for high-dimensional data even when the distributions involved are not necessarily Gaussian. We attempt to bridge the gap between this oft-used computational assumption and the theoretical understanding of why this works, by employing some recent results on random projections on low dimensional subspaces and concentration inequalities.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

In many geophysical or meteorological signal processing applications, Ensemble Kalman Filter (EnKF) or data assimilation has become a popular methodology. Unlike the extended Kalman filter, it does not linearize the dynamics around a nominal trajectory. Instead, it propagates state-observation dynamics as per the original nonlinear rule, but estimates the next state as though it were conditionally Gaussian, using empirical estimates of covariances based on simulated transitions [1]. The Gaussianity hypothesis remains ad hoc, nevertheless the methodology has been found to be very useful by practitioners. Trying to make sense of this ‘unreasonable effectiveness of Gaussianity’ (to borrow a phrase from Wigner) is the motivation behind this work. We do not, however, address the dynamic situation handled by EnKF, but consider the simpler problem of estimating a random variable, given another, in a high dimensional set-up and justify the Gaussian approximation thereof.

Traditionally, Gaussian assumption has been justified either by invoking the classical central limit theorem, postulating that the observed randomness is the cumulative effect of a large number of independent small events (e.g., shot noise), or by the maximum entropy principle, which is a ‘worst case’ analysis. (The two philosophies are not unrelated, as we now know from [2].) What we propose here is a third alternative, also a central limit

theorem, but in large dimension asymptotics rather than large sample asymptotics as in the classical case. The key tool is a result regarding approximate Gaussianity of low dimensional marginals of a class of high dimensional distributions due to Klartag and others. The details follow in subsequent sections.

While EnKF remains our original motivation, the applicability and relevance of this work to other domains is not ruled out. More generally, this work is the first step towards providing a rigorous basis for using Gaussian approximations in high dimensional inference, wherever it occurs, subject to the log-concavity and sparsity hypotheses. We use ideas from compressive sensing to claim that given an  $n$ -dimensional stochastically sparse random vector, one can recover it from samples or measurements that are fewer than  $n$  in number. Compressive sensing essentially deals with the problem of reconstructing a sparse vector from underdetermined measurements. One aims to minimize the  $l_1$ -error between the coefficients of the original vector and the reconstructed one. See [3] for details.

The paper is organized as follows. We outline the problem and the notation in the next section. In Section 2, we present our result for the special case of stochastically sparse vectors. As mentioned earlier, this requires some results from the theory of compressive sensing. Section 3 recalls the key result of Klartag and Eldan on low dimensional projection with nearly Gaussian densities and points out its implications in the present context. Throughout,  $\|\cdot\|$  denotes the standard Euclidean norm in  $\mathbb{R}^n$ .

### 1.1. Outline of the problem

We first show that given a random vector  $(X, Y) \in \mathbb{R}^{n_1+n_2}$ , if  $X, Y$  are sparse,  $E[Y|X]$  can be approximated by projections

<sup>\*</sup> Corresponding authors.

E-mail addresses: [borkar.vs@gmail.com](mailto:borkar.vs@gmail.com) (V.S. Borkar), [dwivediraz@gmail.com](mailto:dwivediraz@gmail.com) (R. Dwivedi), [neeraja.sigma@gmail.com](mailto:neeraja.sigma@gmail.com) (N. Sahasrabudhe).

<sup>1</sup> Raaz Dwivedi is now with Department of EECS, University of California, Berkeley.

on a smaller dimensional subspace. For the final step, where we show that suitable conditional densities can be approximated by Gaussian densities, an additional assumption of log-concavity of conditional density of  $Y$  given  $X$  is required.

We assume throughout that

- (A1)  $E[\|Y\|^2]^{1/2}$  and  $E[\|X\|^2]^{1/2}$  are bounded by some constant  $M < \infty$ , and,  
 (A2) the regular conditional law  $\Psi(\cdot|x)$  of  $Y$  given  $X = x$  has a Lipschitz version as a map  $x \in \mathbb{R}^{n_1} \mapsto \Psi(\cdot|x) \in \mathcal{P}_1(\mathbb{R}^{n_2})$  with Lipschitz constant  $L$ , where  $\mathcal{P}_1(\mathbb{R}^{n_2})$  is the space of probability measures  $\mu$  on  $\mathbb{R}^{n_2}$  with  $\int |x| \mu(dx) < \infty$  under the Wasserstein-1 metric  $\rho(\mu', \mu'') := \inf E[\|Y' - Y''\|]$ . Here the infimum is over all pairs of random variables  $(Y', Y'')$  with law of  $Y'$ , resp.  $Y''$ , being  $\mu'$ , resp.  $\mu''$ . We work with this version throughout.

Let  $\tilde{X}, \tilde{Y}$  denote the orthogonal projection of  $X, Y$  on random  $k_1$  and  $k_2$  dimensional subspaces respectively. By suitable choice of basis, we take this to be the first  $k_1$  co-ordinates of  $X$  and first  $k_2$  co-ordinates of  $Y$ . Let  $\hat{X} = \tilde{X}$  and  $\hat{Y} = \sqrt{\frac{n_2}{k_2}} \tilde{Y}$  be the scaled projection.

We denote by  $\tilde{X}$  and  $\tilde{Y}$  the vectors obtained by padding  $\hat{X}$  and  $\hat{Y}$  by  $n_1 - k_1$  and  $n_2 - k_2$  zeros respectively.

In Section 2 we show that using results in [4] and under suitable conditions of sparsity of  $X$  and  $Y$ , one can approximate  $E[Y|X]$  by  $E[Y^*|X^*]$ , where  $Y^*$  and  $X^*$  are ‘‘good’’ reconstructions of  $X, Y$  from only  $k_1$ , resp.  $k_2$  observations, where a ‘‘good’’ reconstruction means that  $Y$  and  $Y^*$  (resp.  $X$  and  $X^*$ ) are close in standard Euclidean norm with high probability. Furthermore, in Section 3, we show that  $E[Y^*|X^*]$  and therefore  $E[Y|X]$  can be computed approximately using a Gaussian density under a log-concavity assumption on  $\Psi$ .

Let  $G_{n,l}$  denote the Grassmannian of all  $l$ -dimensional subspaces of  $\mathbb{R}^n$ , and let  $\sigma_{n,l}$  stand for the unique rotationally invariant probability measure on  $G_{n,l}$  [5].

## 2. Sparse vectors

Our aim is to estimate  $E[Y|X]$  by  $E[\hat{Y}|\hat{X}]$ , thereby reducing the cost of computation. Using the results in [4], we now show that this aim can be achieved for a ‘‘stochastically sparse’’ vector. In [4], the authors show that it is possible to reconstruct a sparse vector to high accuracy from a small number of random measurements. Let  $|v|_n$  denote the  $n$ th largest entry of the vector  $v$ , or the  $n$ th largest coefficient in a fixed basis. Consider a vector  $v \in \mathbb{R}^N$  such that either  $|v|_n \leq R \cdot n^{-1/p}$  for some  $R > 0$  and some  $0 < p < 1$ , or  $\|v\|_1 \leq R$  for some  $R > 0$  and  $p = 1$ . Consider a random orthonormal basis  $\{\phi_m\}_{m=1}^\infty$  of  $\mathbb{R}^N$  and let  $\theta(g) = [\langle g, \phi_1 \rangle, \dots, \langle g, \phi_N \rangle]^T$  for  $g \in \mathbb{R}^N$ . Suppose we observe only the first  $K$  coefficients in this basis. Let  $F_\Omega$  be the submatrix enumerating those sampled vectors, i.e., the projection operator. One can solve the following optimization problem

$$(P) \min_{g \in \mathbb{R}^N} \|\theta(g)\|_1 \quad \text{subject to } F_\Omega g = F_\Omega v. \quad (1)$$

The solution  $v^*$  is such that for  $\beta > 0$  sufficiently small

$$\|v - v^*\| \leq C_{p,\beta} \cdot R \cdot (K/\log N)^{-r}$$

with probability at least  $1 - O(N^{-\rho/\beta})$ , where  $r = 1/p - 1/2$  and  $\rho > 0$  is a universal constant. Also, as noted in [4], the choice of basis is in fact irrelevant. All that is needed is that the vector  $v$  be sparse in some fixed basis.

The above optimization problem can be reduced to a linear program by the standard technique of replacing each variable (say)  $x$  by  $x^+ - x^-$  and defines a map  $h : v \in \mathbb{R}^N \mapsto v^* \in \mathbb{R}^N$  whenever the solution  $v^*$  is unique. Since the latter holds for a.e.  $v$ ,  $h$  is well

defined as a measurable function. From the 1-homogeneity of the objective function and the constraints, it is easy to see that  $h$  has linear growth.

To use the above results for random vectors, we define the idea of ‘‘stochastically sparse’’ random vectors.

**Definition 2.1.** Let  $Z \in \mathbb{R}^m$  be a random vector. We say that  $Z$  is stochastically sparse if for a prescribed  $\eta_1 > 0$

$$P\left(\sup_n \frac{|Z|_n}{n^{-1/p}} > R\right) < \eta_1 \quad (2)$$

for some  $R > 0$  and  $0 < p < 1$ .

Let  $X^*$  and  $Y^*$  denote the solution to the optimization problem (P) corresponding to stochastically sparse random vectors  $X$  and  $Y$  respectively. Then from the above discussion we have that,  $Y^* = h(\tilde{Y})$  and  $X^* = h(\tilde{X})$ . Define  $H(x) = \int y \Psi(y|x) dy$  and let  $\tilde{\Psi}(\cdot|x)$  denote the image of  $\Psi(\cdot|x)$  under the projection  $\mathbb{R}^{n_1} \mapsto \mathbb{R}^{k_2}$ . Now we can prove the following approximation result.

**Theorem 2.2.** Let  $X \in \mathbb{R}^{n_1}$  and  $Y \in \mathbb{R}^{n_2}$  be stochastically sparse. Let  $\eta_1, \rho$  and  $\beta$  be as defined above. Then, given  $\epsilon > 0$ ,

$$P\left(\left|E[Y|X] - \int h(\tilde{y}) \tilde{\Psi}(\tilde{y}|X^*) d\tilde{y}\right| > \epsilon\right) \leq \frac{4}{\epsilon} \left(\delta_1 + \frac{L\delta_2}{2} + M(2+L)\sqrt{1-q}\right)$$

where,  $q = 1 - 2\eta_1 - O(n_2^{-\rho/\beta}) - O(n_1^{-\rho/\beta})$ .

**Proof.** Using the result in [4], we have that on a set  $B$  with  $P(B) \geq q = 1 - 2\eta_1 - O(n_2^{-\rho/\beta}) - O(n_1^{-\rho/\beta})$ ,

$$\|Y - Y^*\| \leq \delta_1 \quad \text{and} \quad \|X - X^*\| \leq \delta_2 \quad (3)$$

where,

$$\delta_1 = C_{p,\beta} \cdot R \cdot (k_2/\log n_2)^{-r} \quad \text{and} \quad \delta_2 = C_{p,\beta} \cdot R \cdot (k_1/\log n_1)^{-r}$$

for  $r = 1/p - 1/2$ . We have,

$$P\left(\left|E[Y|X] - \int h(\tilde{y}) \tilde{\Psi}(\tilde{y}|X^*) d\tilde{y}\right| > \epsilon\right) \leq P(|E[Y|X] - H(X^*)| > \epsilon/2) + P\left(\left|H(X^*) - \int h(\tilde{y}) \tilde{\Psi}(\tilde{y}|X^*) d\tilde{y}\right| > \epsilon/2\right).$$

Note that

$$P\left(\left|H(X^*) - \int h(\tilde{y}) \tilde{\Psi}(\tilde{y}|X^*) d\tilde{y}\right| > \epsilon/2\right) \leq P(E[\|Y - Y^*\| I\{B\} | X = x]_{x=X^*} > \epsilon/4) + P(E[\|Y - Y^*\| I\{B^c\} | X = x]_{x=X^*} > \epsilon/4).$$

From stochastic sparsity of  $Y$ , we get

$$P(E[\|Y - Y^*\| I\{B\} | X = x]_{x=X^*} > \epsilon/4) \leq \frac{4}{\epsilon} \delta_1 \quad (4)$$

and,

$$P(E[\|Y - Y^*\| I\{B^c\} | X = x] > \epsilon/4) \leq \frac{4}{\epsilon} E[\|Y - Y^*\|^2]^{1/2} \sqrt{P(B^c)} \leq \frac{8}{\epsilon} M \sqrt{(1-q)}. \quad (5)$$

Download English Version:

<https://daneshyari.com/en/article/751932>

Download Persian Version:

<https://daneshyari.com/article/751932>

[Daneshyari.com](https://daneshyari.com)