

Representation formulae for score functions

Ivan Nourdin, Giovanni Peccati and Yvik Swan*

Département de Mathématique, Université de Liège

July 2, 2014

- 1 Score
- 2 Stein and Fisher
- 3 Controlling the relative entropy
- 4 Key identity
- 5 Cattywampus Stein's method
- 6 Extension
- 7 Coda

Scoooores



- 1 Score
- 2 Stein and Fisher**
- 3 Controlling the relative entropy
- 4 Key identity
- 5 Cattywampus Stein's method
- 6 Extension
- 7 Coda

Let X be a centered d -random vector with covariance $B > 0$.

Definition

The *Stein kernel* of X is a $d \times d$ matrix $\tau_X(X)$ such that

$$E [\tau_X(X) \nabla \varphi(X)] = E [X \varphi(X)]$$

for all $\varphi \in C_c^\infty(\mathbb{R}^d)$.

Definition

The *score* of X is the $d \times 1$ vector $\rho_X(X)$ such that

$$E [\rho_X(X) \varphi(X)] = -E [\nabla \varphi(X)]$$

for all $\varphi \in C_c^\infty(\mathbb{R}^d)$.

In the Gaussian case $Z \sim \mathcal{N}_d(0, C)$ the **Stein identity**

$$E[Z\varphi(Z)] = E[C\nabla\varphi(Z)]$$

gives

$$\rho_Z(Z) = -C^{-1}Z \text{ and } \tau_Z(Z) = C.$$

Intuitively, a measure of proximity

$$\rho_X(X) \approx -B^{-1}X$$

and

$$\tau_X(X) \approx B$$

should provide an assessment of “Gaussianity”.

Definition

The *standardised Fisher information* of X is

$$J_{st}(X) = BE \left[(\rho_X(X) + B^{-1}X) (\rho_X(X) + B^{-1}X)^T \right].$$

A simple computation gives

$$J_{st}(X) = BJ(X) - Id$$

with $J(X) = E [\rho_X(X)\rho_X(X)^T]$ the **Fisher information matrix**.

Definition

The *Stein discrepancy* is

$$S(X) = E [\|\tau_X(X) - B\|_{H.S.}^2].$$

Control on $J_{st}(X)$ and $S(X)$ provides control on several distances (Kullback-Leibler, Kolmogorov, Wasserstein, Hellinger, Total Variation, ...) between the law of X and the Gaussian.

Controlling $J_{st}(X)$:

- Johnson and Barron through careful analysis of the score function (PTRF, 2004)
- Artstein, Ball, Barthe, Naor through “variational tour de force” (PTRF, 2004)

Controlling $S(X)$:

- Cacoullos Papathanassiou and Utev (AoP 1994) in a number of settings
- Nourdin and Peccati through their infamous Malliavin/Stein fourth moment theorem (PTRF, 2009)
- Extension to abstract settings (Ledoux, AoP 2012)

- 1 Score
- 2 Stein and Fisher
- 3 Controlling the relative entropy**
- 4 Key identity
- 5 Cattywampus Stein's method
- 6 Extension
- 7 Coda

Let Z be centered Gaussian with density $\phi = \phi_d(\cdot; C)$.

Definition

The *relative entropy* between X and Z is

$$D(F \parallel Z) = E [\log(f(X)/\phi(X))] = \int_{\mathbb{R}^d} f(x) \log \left(\frac{f(x)}{\phi(x)} \right) dx.$$

The Pinsker-Csiszar-Kullback inequality yields

$$2TV(X, Z) \leq \sqrt{2D(X \parallel Z)}.$$

In other words

$$D(X \parallel Z) \Rightarrow TV(X, Z)^2.$$

Usefulness of $J_{st}(X)$ can be seen via the **de Bruijn identity**.

Let $X_t = \sqrt{t}X + \sqrt{1-t}Z$ and $\Gamma_t = tB + (1-t)C$. Then

$$\begin{aligned} D(X \parallel Z) &= \int_0^1 \frac{1}{2t} \operatorname{tr} (C\Gamma_t^{-1} J_{st}(X_t)) dt \\ &\quad + \frac{1}{2} (\operatorname{tr}(C^{-1}B) - d) + \int_0^1 \frac{1}{2t} \operatorname{tr} (C\Gamma_t^{-1} - I_d) dt \end{aligned}$$

In other words

$$J_{st}(X_t) \Rightarrow D(X \parallel Z) \Rightarrow (TV(X, Z))^2.$$

Usefulness of $S(X)$ can be seen via **Stein's method**.

Fix $d = 1$. Then, given $h : \mathbb{R} \rightarrow \mathbb{R}$ such that $\|h\|_\infty \leq 1$ seek g_h solution of the Stein equation to get

$$\begin{aligned} E[h(X)] - E[h(Z)] &= E[g'_h(X) - Xg_h(X)] \\ &= E[(1 - \tau_X(X))g'_h(X)] \end{aligned}$$

so that

$$\begin{aligned} TV(X, Z) &= \frac{1}{2} \sup_{\|h\|_\infty \leq 1} |E[h(X)] - E[h(Z)]| \\ &\leq \left(\frac{1}{2} \sup_{\|h\|_\infty \leq 1} \|g'_h\| \right) \sqrt{S(X)}. \end{aligned}$$

In other words

$$S(X) \Rightarrow TV(X, Z)^2.$$

If h is not smooth there is no way of obtaining sufficiently precise estimates on the quantity “ ∇g_h ” in dimension greater than 1.

For the moment Stein’s method only works in dimension 1 for total variation distance.

The IT approach via de Bruijn’s identity does not suffer from this “dimensionality issue” .

We aim to mix the Stein method approach and the IT approach.

To this end we need one final ingredient : **a representation formulae for the score in terms of the Stein kernel.**

- 1 Score
- 2 Stein and Fisher
- 3 Controlling the relative entropy
- 4 Key identity**
- 5 Cattywampus Stein's method
- 6 Extension
- 7 Coda

Theorem

Let $X_t = \sqrt{t}X + \sqrt{1-t}Z$ with X and Z independent. Then

$$\rho_t(X_t) + C^{-1}X_t = -\frac{t}{\sqrt{1-t}}E[(I_d - C^{-1}\tau_X(X))Z | X_t] \quad (1)$$

for all $0 < t < 1$.

Proof when $d = 1$ and $C = 1$.

$$\begin{aligned} E[E[(1 - \tau_X(X))Z | X_t] \phi(X_t)] &= E[(1 - \tau_X(X))Z \phi(X_t)] \\ &= \sqrt{1-t}E[\phi'(X_t)] - \sqrt{1-t}E[\tau_X(X)\phi'(X_t)] \\ &= \sqrt{1-t}E[\phi'(X_t)] - \sqrt{\frac{1-t}{t}}E[X\phi(X_t)] \\ &= \sqrt{1-t}E[\phi'(X_t)] - \frac{\sqrt{1-t}}{t}E[X_t\phi(X_t)] + \frac{1-t}{t}E[Z\phi(X_t)] \\ &= \sqrt{1-t}E[\phi'(X_t)] - \frac{\sqrt{1-t}}{t}E[X_t\phi(X_t)] + \frac{1-t}{t}\sqrt{1-t}E[\phi'(X_t)] \\ &= -\frac{\sqrt{1-t}}{t}(E[\phi'(X_t)] - E[X_t\phi(X_t)]) \end{aligned}$$

This formula provides a **nearly one-line argument**.

Define

$$\Delta(X, t) = E [(I_d - C^{-1}\tau_X(X))Z | X_t].$$

Take $d = 1$ and all variances set to 1. Then

$$J_{st}(X_t) = E [(\rho_t(X_t) + X_t)^2] = \frac{t^2}{1-t} E [\Delta(X, t)^2]$$

so that

$$D(X || Z) = \frac{1}{2} \int_0^1 \frac{t}{1-t} E [\Delta(X, t)^2] dt.$$

Also,

$$E [\Delta(X, t)^2] \leq E [(1 - \tau_X(X))^2] = S(X).$$

This yields

$$D(X||Z) \leq \frac{1}{2}S(X) \int_0^1 \frac{t}{1-t} dt$$

which is useless.

There is hope, nevertheless :

$$\int_0^1 \frac{t}{1-t} dt$$

is barely infinity.

Recall $X_t = \sqrt{t}X + \sqrt{1-t}Z$. Then

$$\Delta(X, t) = E[(1 - \tau_X(X))Z | X_t]$$

is such that

$$\Delta(X, 0) = \Delta(X, 1) = 0 \text{ a.s.}$$

Hence we need to identify conditions under which

$$\frac{t}{1-t} E[\Delta(X, t)^2]$$

is integrable at $t = 1$.

The behaviour of $\Delta(X, t)$ around $t \approx 1$ is central to the understanding of the law of X .

The behaviour of

$$E [\Delta(X, t)^2] \text{ at } t \approx 1$$

is closely connected to the so-called **MMSE dimension** studied by the IT community.

This quantity revolves around the remarkable “MMSE formula”

$$\frac{d}{dr} I(X; \sqrt{r}X + Z) = E [(X - E[X | \sqrt{r}X + Z])^2]$$

due to Guo, Shamai and Verdu (IEEE, 2005)

The connexion is explicitly stated in NPSb (IEEE, 2014).

- 1 Score
- 2 Stein and Fisher
- 3 Controlling the relative entropy
- 4 Key identity
- 5 Cattywampus Stein's method**
- 6 Extension
- 7 Coda

In NPSa (JFA, 2014) we suggest the following IT alternative to Stein's method.

First cut the integral :

$$\begin{aligned} & 2D(X||Z) \\ & \leq E [(1 - \tau_X(X))^2] \int_0^{1-\epsilon} \frac{t}{1-t} dt + \int_{1-\epsilon}^1 \frac{t}{1-t} E [\Delta(X, t)^2] dt \\ & \leq E [(1 - \tau_X(X))^2] |\log \epsilon| + \int_{1-\epsilon}^1 \frac{t}{1-t} E [\Delta(X, t)^2] dt. \end{aligned}$$

Next suppose that when t is close to 1 we have

$$E [\Delta(X, t)^2] \leq C_\kappa t^{-1} (1-t)^\kappa \quad (2)$$

for some $\kappa > 0$.

We deduce

$$\begin{aligned} 2D(X \parallel Z) &\leq S(X)|\log \epsilon| + C_\eta \int_{1-\epsilon}^1 (1-t)^{-1+\kappa} dt \\ &= S(X)|\log \epsilon| + \frac{C_\kappa}{\kappa} \epsilon^\kappa. \end{aligned}$$

The optimal choice is $\epsilon = E[(1 - \tau_X(X))^2]^{1/\kappa}$ which leads to

$$D(X \parallel Z) \leq \frac{1}{2\kappa} S(X) \log S(X) + \frac{C_\kappa}{2\kappa} S(X)$$

which provides a bound on the total variation distance in terms of $S(X)$ which is of the correct order up to a logarithmic factor.

Under what conditions do we have (2)?

It is relatively easy to show (via Hölder's inequality) that

$$E \left[|\tau_X(X)|^{2+\eta} \right] < \infty \text{ and } E [|\Delta(X, t)|] \leq ct^{-1}(1-t)^\delta \quad (3)$$

implies (2).

It now remains to identify under which conditions we have (3).

Lemma (Poly's first lemma)

Let X be an integrable random variable and let Y be a \mathbb{R}^d -valued random vector having an absolutely continuous distribution. Then

$$E |E[X | Y]| = \sup E [Xg(Y)],$$

where the supremum is taken over all $g \in C_c^1$ such that $\|g\|_\infty \leq 1$

Thus

$$E |E [Z(1 - \tau_X(X)) | X_t]| = \sup E [Z(1 - \tau_X(X))g(X_t)].$$

Now choose $g \in C_c^1$ such that $\|g\|_\infty \leq 1$. Then

$$\begin{aligned} & E [Z(1 - \tau_X(X))g(X_t)] \\ &= E [Zg(X_t)] - E [Zg(X_t)\tau_X(X)] \\ &= E [Zg(X_t)] - \sqrt{1-t}E [\tau_X(X)g'(X_t)] \\ &= E [Z(g(X_t) - g(X))] - \sqrt{\frac{1-t}{t}}E [g(X_t)X] \end{aligned}$$

and thus

$$|E [Z(1 - \tau_X(X))g(X_t)]| \leq |E [Z(g(X_t) - g(X))]| + t^{-1}\sqrt{1-t}.$$

Also

$$\begin{aligned} & \sup |E [Z (g(X_t) - g(X))]| \\ &= \sup \left| \int_{\mathbb{R}} x E \left[g(\sqrt{t}X + \sqrt{1-t}x) - g(X) \right] \phi_1(x) dx \right| \\ &\leq 2 \int_{\mathbb{R}} |x| TV(\sqrt{t}X + \sqrt{1-t}x, X) \phi_1(x) dx. \end{aligned}$$

Wrapping up we get

$$\begin{aligned} & E |E [Z(1 - \tau_X(X)) | X_t]| \\ &\leq 2E \left[|Z| TV(\sqrt{t}X + \sqrt{1-t}Z, X) \right] + t^{-1} \sqrt{1-t}. \end{aligned}$$

It therefore all boils down to a condition on

$$TV(\sqrt{t}X + \sqrt{1-t}x, X).$$

Recall that we want

$$E |E [Z(1 - \tau_X(X)) | X_t]| \leq ct^{-1}(1 - t)^\delta. \quad (3)$$

As it turns out, in view of previous results, a sufficient condition for (3) is

$$TV(\sqrt{t}X + \sqrt{1 - t}x, X) \leq \kappa(1 + |x|)t^{-1}(1 - t)^\alpha.$$

This condition – and its multivariate extension – is satisfied by a wide family of random vectors including those for which they can apply their fourth moment bound

$$S(X) \leq c(E [X^4] - 3).$$

Theorem (Entropic CLTs on Wiener chaos)

Let $d \geq 1$ and $q_1, \dots, q_d \geq 1$ be fixed integers. Consider vectors

$$F_n = (F_{1,n}, \dots, F_{d,n}) = (I_{q_1}(h_{1,n}), \dots, I_{q_d}(h_{d,n})), \quad n \geq 1,$$

with $h_{i,n} \in \mathfrak{H}^{\odot q_i}$. Let C_n denote the covariance matrix of F_n and let $Z_n \sim \mathcal{N}_d(0, C_n)$ be a centered Gaussian random vector in \mathbb{R}^d with the same covariance matrix as F_n .

Let

$$\Delta_n := E[\|F_n\|^4] - E[\|Z_n\|^4],$$

Assume that $C_n \rightarrow C > 0$ and $\Delta_n \rightarrow 0$, as $n \rightarrow \infty$.

Then, the random vector F_n admits a density for n large enough, and

$$D(F_n \| Z_n) = O(1) \Delta_n |\log \Delta_n| \quad \text{as } n \rightarrow \infty, \quad (4)$$

where $O(1)$ indicates a bounded numerical sequence depending on d, q_1, \dots, q_d , as well as on the sequence $\{F_n\}$.

- 1 Score
- 2 Stein and Fisher
- 3 Controlling the relative entropy
- 4 Key identity
- 5 Cattywampus Stein's method
- 6 Extension**
- 7 Coda

Let X_i , $i = 1, \dots, n$ be independent random vectors with Stein kernels $\tau_i(X_i)$ and score functions $\rho_i(X_i)$, $i = 1, \dots, n$.

For all $t = (t_1, \dots, t_n) \in [0, 1]^d$ such that $\sum_{i=1}^n t_i = 1$ we define

$$W_t = \sum_{i=1}^n \sqrt{t_i} X_i$$

and denote Γ_t the corresponding covariance matrix. Then

$$\rho_t(W_t) + \Gamma_t^{-1} W_t = \sum_{i=1}^n \frac{t_i}{\sqrt{t_{i+1}}} E \left[(Id - \Gamma_t^{-1} \tau_i(X_i)) \rho_{i+1}(X_{i+1}) | W_t \right]$$

where we identify $X_{n+1} = X_1$ and $t_{n+1} = t_1$.

Lemma (Poly's second lemma)

Let X and Y be square-integrable random variables with mean $E[X] = 0$. Then

$$E \left[(E[X | Y])^2 \right] = \sup_{\varphi \in \mathcal{H}(Y)} (E[X\varphi(Y)])^2,$$

where the supremum is taken over the collection $\mathcal{H}(Y)$ of functions φ such that $E[\varphi(Y)] = 0$ and $E[\varphi(Y)^2] \leq 1$.

Theorem

Let $W_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$ where the X_i are independent random variables with Stein factor $\tau_i(X_i)$ and score function $\rho_i(X_i)$. Then

$$J_{st}(W_n) = \sup_{\varphi \in \mathcal{H}(W_n)} (E[\varphi'(W_n) - W_n\varphi(W_n)])^2.$$

There seem to be many applications of the last formula.

For instance the difference

$$J_{st}(W_{n+1}) - J_{st}(W_n)$$

can be studied in quite some detail.

We had hoped to obtain the “entropy jump inequality” as well as the “increasingness of entropy”.

There is, however, some work left before we can hooray.

- 1 Score
- 2 Stein and Fisher
- 3 Controlling the relative entropy
- 4 Key identity
- 5 Cattywampus Stein's method
- 6 Extension
- 7 Coda**

Just a final word to say

thank you

to Janna, Jay and Larry for the great conference.

The key is a generalisation Carbery and Wright inequality : there is a universal constant $c > 0$ such that, for any polynomial $Q : \mathbb{R}^n \rightarrow \mathbb{R}$ of degree at most d and any $\alpha > 0$ we have

$$E[Q(X_1, \dots, X_n)^2]^{\frac{1}{2d}} P(|Q(X_1, \dots, X_n)| \leq \alpha) \leq cd\alpha^{\frac{1}{d}},$$

where X_1, \dots, X_n are independent random variables with common distribution $\mathcal{N}(0, 1)$.

Explicit conditions : fix $d, q_1, \dots, q_d \geq 1$,

- 1 let $F = (F_1, \dots, F_d)$ be a random vector such that $F_i = I_{q_i}(h_i)$ with $h_i \in \mathfrak{H}^{\odot q_i}$
- 2 set $N = 2d(q - 1)$ with $q = \max_{1 \leq i \leq d} q_i$
- 3 Let C be the covariance matrix of F

Let $\Gamma = \Gamma(F)$ denote the Malliavin matrix of F , and assume that $E[\det \Gamma] > 0$ (which is equivalent to assuming that F has a density).

There exists a constant $c_{q,d,\|C\|_{H.S.}} > 0$ (depending only on q, d and $\|C\|_{H.S.}$ — with a continuous dependence in the last parameter) such that, for any $\mathbf{x} \in \mathbb{R}^d$ and $t \in [\frac{1}{2}, 1]$,

$$\begin{aligned} & \mathbf{TV}(\sqrt{t}F + \sqrt{1-t}\mathbf{x}, F) \\ & \leq c_{q,d,\|C\|_{H.S.}} \left(\beta^{-\frac{1}{N+1}} \wedge 1 \right) (1 + \|\mathbf{x}\|_1) (1-t)^{\frac{1}{2(2N+4)(d+1)+2}}. \end{aligned}$$

Theorem (Entropic fourth moment bound)

Let $F_n = (F_{1,n}, \dots, F_{d,n})$ be a sequence of d -dimensional random vectors such that: (i) $F_{i,n}$ belongs to the q_i th Wiener chaos of \mathbf{G} , with $1 \leq q_1 \leq q_2 \leq \dots \leq q_d$; (ii) each $F_{i,n}$ has variance 1, (iii) $E[F_{i,n}F_{j,n}] = 0$ for $i \neq j$, and (iv) the law of F_n admits a density f_n on \mathbb{R}^d . Write

$$\Delta_n := \int_{\mathbb{R}^d} \|\mathbf{x}\|^4 (f_n(\mathbf{x}) - \phi_d(\mathbf{x})) d\mathbf{x},$$

where $\|\cdot\|$ stands for the Euclidean norm, and assume that $\Delta_n \rightarrow 0$, as $n \rightarrow \infty$. Then,

$$\int_{\mathbb{R}^d} f_n(\mathbf{x}) \log \frac{f_n(\mathbf{x})}{\phi_d(\mathbf{x})} d\mathbf{x} = O(1) \Delta_n |\log \Delta_n|, \quad (5)$$

where $O(1)$ stands for a bounded numerical sequence, depending on d, q_1, \dots, q_d and on the sequence $\{F_n\}$.

Corollary

Let $d \geq 1$ and $q_1, \dots, q_d \geq 1$ be fixed integers. Consider vectors

$$F_n = (F_{1,n}, \dots, F_{d,n}) = (I_{q_1}(h_{1,n}), \dots, I_{q_d}(h_{d,n})), \quad n \geq 1,$$

with $h_{i,n} \in \mathfrak{H}^{\odot q_i}$. Let C_n denote the covariance matrix of F_n and let $Z_n \sim \mathcal{N}_d(0, C_n)$ be a centered Gaussian random vector in \mathbb{R}^d with the same covariance matrix as F_n . Assume that $C_n \rightarrow C > 0$. Then, the following three assertions are equivalent, as $n \rightarrow \infty$:

- (i) $\Delta_n \rightarrow 0$;
- (ii) F_n converges in distribution to $Z \sim \mathcal{N}_d(0, C)$;
- (iii) $D(F_n \| Z_n) \rightarrow 0$.