

Stein's (magic) method

A. D. Barbour and Louis H. Y. Chen

Universität Zürich and National University of Singapore

Written in celebration of Charles Stein's 90th birthday in 2010

1. Introduction.

One of the greatest achievements of probability has been its success in approximating the distributions of arbitrarily complicated random processes in terms of a rather small number of 'universal' processes — Brownian motion, the Poisson process, the Ewens sampling formula, Airy processes, stochastic Loewner evolution and so on. The standard approach is to consider sequences of processes, indexed by a parameter n , and to establish that suitably normalized versions of the processes converge in distribution to one of the standard processes as n tends to infinity. However, for practical purposes, it is much more important to know how accurate such an approximation is for a particular process with a fixed value of n , and this is a more difficult question to answer. For instance, central limit theorems were known already around 1715, and in full generality by 1900, whereas the corresponding approximation theorem of Berry and Esseen was only proved in 1941. Stein's method, introduced in 1970, offers a general means of solving such problems. By constructing and exploiting a novel characteristic operator associated with a random system — most often, the one used as the approximation — it turns out to be possible to make precise assessments of the approximation error in a wide variety of circumstances.

Stein's original application was in the context of central limit approximation to partial sums of random variables having a stationary dependence structure, a problem involving the normal distribution and the real line. However, his method has a big advantage over most other techniques, in that it can in principle be used for approximation in terms of any distribution on any space, including random variables on the real line, processes on a space of sequences, functions or measures, and combinatorial structures on discrete spaces. A further big advantage over its competitors is that strong independence requirements are not needed to make the method work (though they may of course simplify many arguments and the form of the bounds that can be attained). As a result of this considerable freedom, its uses have proliferated, with approximations not only to the normal distribution, but also

to the Poisson distribution, to multivariate normal distributions, to diffusions, to Poisson processes, to the Ewens sampling formula, to the Wigner semi-circle law, and more.

The method continues to produce new developments in a wide variety of settings. The last five years alone have seen great progress in problems concerning large deviations and concentration of measure inequalities; in the application of Stein’s method to problems having an essentially algebraic component; in proving bounds for normal and gamma approximations to the distributions of functionals of infinite–dimensional Gaussian fields, using a combination of Stein’s method and Malliavin calculus; and in a range of problems involving random geometrical graphs. It is a tribute to the importance of Stein’s original and amazing idea that it continues to inspire vigorous research 40 years after its inception. In this article, we can do no more than scratch the surface of a wide-ranging and still mysterious topic.

2. Normal approximation

Stein originally introduced his method in a course of lectures at Stanford. Dissatisfied with the then available proofs of the combinatorial central limit theorem, he devised a new one for his class, which entirely dispensed with the use of Fourier analysis. His paper (Stein, 1972) in the VI’t h Berkeley Symposium contained the first published version of his method, giving Berry–Esseen bounds for the accuracy of the normal approximation to the distribution of a sum of a stationary sequence of random variables; already, the flexibility of the method is in evidence, since the dependence setting is quite different from that of the combinatorial central limit theorem, with which he began. One way of arriving at his new approach to normal approximation is as follows.

Let Z be a random variable on \mathbf{R} with differentiable probability density p . Our main emphasis here is for p the standard normal density, but this argument works more generally. Now let \mathcal{F} be the set of differentiable real functions f such that $f(x)p(x) \rightarrow 0$ as $|x| \rightarrow \infty$. Then, clearly,

$$\begin{aligned} 0 &= \int_{-\infty}^{\infty} \frac{d}{dx} \{f(x)p(x)\} dx \\ &= \int_{-\infty}^{\infty} \{f'(x) + \psi(x)f(x)\}p(x) dx, \end{aligned} \tag{2.1}$$

with $\psi(x) := p'(x)/p(x)$: for p the standard normal density, $\psi(x) = -x$. Hence any function h of the form $f' + \psi f$ with $f \in \mathcal{F}$ automatically has $\mathbf{E}h(Z) = 0$. Conversely, given any continuous h with $\mathbf{E}|h(Z)| < \infty$, the function $h - \mathbf{E}h(Z)$ can be written in the

form $f' + \psi f$ with $f \in \mathcal{F}$. To do so is simple. First observe that, if $g = f' + \psi f$, then

$$\int_{-\infty}^X g(x)p(x) dx = \int_{-\infty}^X \frac{d}{dx} \{f(x)p(x)\} dx = f(X)p(X), \quad (2.2)$$

and that this argument can be reversed: f so defined, for an arbitrary g for which $\mathbf{E}|g(Z)| < \infty$, is such that $g = f' + \psi f$. Hence we can take $f = f_h$ to be defined by

$$\begin{aligned} f_h(X)p(X) &= \int_{-\infty}^X \{h(x) - \mathbf{E}h(Z)\}p(x) dx \\ &= - \int_X^{\infty} \{h(x) - \mathbf{E}h(Z)\}p(x) dx, \end{aligned} \quad (2.3)$$

noting that it is then directly checked that $f_h \in \mathcal{F}$. This allows us to write

$$\mathbf{E}h(W) - \mathbf{E}h(Z) = \mathbf{E}\{f'_h(W) + \psi(W)f_h(W)\}, \quad (2.4)$$

for any random variable W for which the expectations exist; in particular, for standard normal approximation,

$$\mathbf{E}h(W) - \mathbf{E}h(Z) = \mathbf{E}\{f'_h(W) - Wf_h(W)\}. \quad (2.5)$$

Taking the supremum of the left hand side in (2.4) over test functions h in some suitable class \mathcal{H} gives a (very concrete) measure of the distance between the distributions of W and Z . This distance can in turn be computed by taking the supremum of the right hand side of (2.4) over $h \in \mathcal{H}$. The intuition is then that, since the right hand side is exactly zero for *all* $f \in \mathcal{F}$ if $\mathcal{L}(W) = \mathcal{L}(Z)$, it should ‘automatically’ be close to zero if $\mathcal{L}(W) \approx \mathcal{L}(Z)$, again more or less irrespective of f , and hence the supremum of the right hand side of (2.4) for $h \in \mathcal{H}$ may indeed be shown to be small, as it were by right. The remarkable power of the method derives from the fact that, in many many circumstances, expressing the difference $\mathbf{E}h(W) - \mathbf{E}h(Z)$ in the entirely equivalent form $\mathbf{E}\{f'_h(W) + \psi(W)f_h(W)\}$ makes it much easier to bound.

A typical setting, reasonable for standard normal approximation, is one in which W is a sum $\sum_{i=1}^n X_i$ of many individually small and only weakly dependent random variables X_i , with zero means and with $\text{Var } W = 1$. Since, for standard normal approximation, $\psi(W) = -\sum_{i=1}^n X_i$, the expression $\mathbf{E}\{\psi(W)f_h(W)\}$ can be broken up into a sum of terms $\mathbf{E}\{-X_i f_h(W)\}$, in which X_i has only limited influence on the whole sum W . This can then be exploited, together with a Taylor expansion, to show that $\mathbf{E}\{f'_h(W) + \psi(W)f_h(W)\}$ is

small. For instance, for *independent* summands X_i , one can write

$$\begin{aligned} \mathbf{E}\{Wf_h(W)\} &= \sum_{i=1}^n (\mathbf{E}\{X_i f_h(W - X_i)\} + \mathbf{E}\{X_i [f_h(W) - f_h(W - X_i)]\}) \\ &= \sum_{i=1}^n (\mathbf{E}\{X_i^2 f_h'(W - X_i)\} + \mathbf{E}\{X_i [f_h(W) - f_h(W - X_i) - X_i f_h'(W - X_i)]\}), \end{aligned}$$

using the independence of $W - X_i$ and X_i with $\mathbf{E}X_i = 0$, and compare the result with

$$\mathbf{E}f_h'(W) = \sum_{i=1}^n \mathbf{E}\{X_i^2\} \mathbf{E}f_h'(W),$$

because $\mathbf{E}W^2 = 1$. This immediately gives

$$|\mathbf{E}h(W) - \mathbf{E}h(Z)| = |\mathbf{E}\{f_h'(W) - Wf(W)\}| \leq \frac{3}{2} \sum_{i=1}^n \mathbf{E}|X_i|^3 \|f_h''\|_\infty, \quad (2.6)$$

yielding an explicit Lyapounov bound for classes of test functions \mathcal{H} for which the supremum $\sup_{h \in \mathcal{H}} \|f_h''\|_\infty < \infty$. This elementary argument can of course be improved in many ways, to deal with less restrictive classes of test functions: see Section 6. What is more important to notice is that independence can immediately be replaced by some form of local dependence in the argument, without any great change in spirit, by using the dissection

$$X_i f_h(W) = X_i f_h(W - Y_i) + X_i [f_h(W) - f_h(W - Y_i)],$$

where now X_i and $W - Y_i$ are (almost) independent and Y_i is still small enough not to have too big an influence on W . This illustrates that the method can be much more robust with respect to dependence than the classical transform methods.

The discussion above is intentionally kept as simple as possible, and this to some extent disguises the power of the method. Even for normal approximation, Stein's method has had many startling successes. One of the earliest was Bolthausen's (1984) Berry–Esseen Lyapounov bound for the error in the combinatorial central theorem, which had been the subject of much research in the previous decades. Götze's (1991) multivariate normal approximation theorem is another: using a proof based on Stein's method, he was able to establish an error bound of order $O(n^{-1/2})$ for the probabilities of an extremely large collection of sets, including all convex sets, for sums of independent random vectors. Rinott & Rotar (1996) developed ideas originating in Stein's proof of (5.3) and in Götze (1991), in proving an explicit and very effective error bound, typically of order $n^{-1/2} \log n$, for the normal approximation of sums of bounded, dependent random vectors

X_j in \mathbf{R}^d . The dependence structure envisaged there is one in which only relatively few other “neighbouring” random vectors are significantly dependent on any given X_j , with the meaning of neighbouring to be chosen to suit: see also Chen & Shao (2004).

3. Stein’s method in a nutshell.

Stein’s method can be simply described as follows. Let W and Z be random elements taking values in a space \mathcal{S} , and let \mathcal{X} and \mathcal{Y} be classes of bounded real-valued functions defined on \mathcal{S} . In approximating the distribution $\mathcal{L}(W)$ of W by the distribution $\mathcal{L}(Z)$ of Z , we write $\mathbf{E}h(W) - \mathbf{E}h(Z) = \mathbf{E}\{Lf_h(W)\}$ for $h \in \mathcal{Y}$, where L is a linear operator from \mathcal{X} into \mathcal{Y} , and f_h a bounded solution of the equation $Lf = h - \mathbf{E}h(Z)$. The error $\mathbf{E}\{Lf_h(W)\}$ can then be bounded by studying the solution f_h and by exploiting the probabilistic properties of W .

Usually, we take as test functions a class $\mathcal{H} \subset \mathcal{Y}$, large enough to be separating; that is, large enough so that, if $\mathbf{E}h(W) = \mathbf{E}h(Z)$ for all $h \in \mathcal{H}$, then $\mathcal{L}(W) = \mathcal{L}(Z)$. \mathcal{X} is then assumed to contain all f_h for $h \in \mathcal{H}$. In this case, $\mathbf{E}\{Lf(W)\} = 0$ for all $f \in \mathcal{X}$ if and only if $\mathcal{L}(W) = \mathcal{L}(Z)$. Such an L characterizes $\mathcal{L}(Z)$; the equation $Lf = h - \mathbf{E}h(Z)$ is called a Stein equation for $\mathcal{L}(Z)$, and L a Stein operator for $\mathcal{L}(Z)$. In the particular case of normal approximation, where W and Z are real-valued random variables and $\mathcal{L}(Z) = \mathcal{N}(0, 1)$, the operator L used by Stein (1972) is given by $Lf(w) = f'(w) - wf(w)$, as noted above.

4. Stein identities.

The discussion so far has been rather vague as to why it might be easier to show that $\mathbf{E}\{Lf_h(W)\}$ is small than it is to show the same for $\mathbf{E}h(W) - \mathbf{E}h(Z)$ directly. For the normal distribution, the motivation came by way of an analytic plausibility argument, and it was demonstrated by example that it was true for the chosen Stein operator L . However, it is implied in Stein (1992) that one can directly exploit the probabilistic properties of W , using auxiliary randomization, in order to find a linear operator L for which the error $\mathbf{E}\{Lf_h(W)\}$ is manageable, and indeed thereby deducing the distribution $\mathcal{L}(Z)$ that should be used to approximate $\mathcal{L}(W)$. The idea is to begin by finding an operator \tilde{L} such that $\mathbf{E}\{\tilde{L}f(W)\} = 0$ for all $f \in \mathcal{X}$; that is, to begin by finding a Stein identity for $\mathcal{L}(W)$. Once this has been done, one can look for an operator L that characterizes a better known distribution $\mathcal{L}(Z)$, with the property that Lf is close to $\tilde{L}f$ for all $f \in \mathcal{X}$. The error $\mathbf{E}\{Lf_h(W)\}$ in the approximation is then equivalently expressed as $\mathbf{E}\{Lf_h(W) - \tilde{L}f_h(W)\}$, because $\mathbf{E}\{\tilde{L}f(W)\} = 0$ for all f , and this difference is small because Lf_h is close to $\tilde{L}f_h$.

This approach is most simply illustrated when the random variable W is known to have the equilibrium distribution of some Markov process \tilde{Z} with generator \tilde{L} . In that case, for suitable functions f , $\mathbf{E}\{\tilde{L}f(W)\} = 0$ by Dynkin's formula, and if \tilde{L} is close to the generator L of another Markov process whose limit distribution $\mathcal{L}(Z)$ is well known, approximations can be deduced. More frequently, W has the equilibrium distribution of some function $g(\tilde{Z})$, so that $\tilde{L}f(W)$ is replaced by $\tilde{L}(f \circ g)(\tilde{Z})$, which need not be a function of $g(\tilde{Z})$ alone; however, it may still be possible to find an operator L on \mathcal{X} such that $Lf(g(\tilde{Z}))$ is close to $\tilde{L}(f \circ g)(\tilde{Z})$, which is all that is needed to make the method work. Indeed, having guessed such an L , one can then consider the expressions $\mathbf{E}\{Lf_h(W)\}$ directly, with a well founded hope that they can be shown to be small. This idea has become known as the generator method (Barbour 1988). One of its greatest successes has been the multivariate normal approximation theorem of Götze (1991), referred to above. The use of the generator method in Poisson process approximation is discussed in Section 7, and its application to diffusion approximation can be found in Barbour (1990).

As a simpler example, suppose that $W = \sum_{i=1}^n X_i$, with the X_1, \dots, X_n independent and having zero means, and with $\mathbf{E}W^2 = 1$. Define a Markov chain \tilde{Z} as follows. Let I_j , $j \geq 1$, be independent (also of the X_i 's) and uniformly distributed on $\{1, 2, \dots, n\}$, and let $\tilde{Z}(0) = (X_1, \dots, X_n)$. Then, if $\tilde{Z}(j-1) = z$, define $\tilde{Z}(j) = z + \mathbf{e}_{I_j}(X_{I_j}^{(j)} - z_{I_j})$, where the $X_{i(j)}$, $j \geq 1$, are independent copies of X_i , $1 \leq i \leq n$, and all are independent of everything else: \mathbf{e}_i denotes the i -th unit vector. Clearly, \tilde{Z} is a stationary Markov chain, and its distribution at any time is that of independent random variables with the distributions of the X_i . The generator is given by

$$\tilde{L}f(z) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}\{f(z + \mathbf{e}_i(X_i - z_i)) - f(z)\}. \quad (4.1)$$

We are actually interested in $W = g(\tilde{Z}(0))$, where $g(z) = \sum_{i=1}^n z_i$. For functions $f \circ g$ for this g , we have

$$\begin{aligned} \tilde{L}(f \circ g)(z) &= \frac{1}{n} \sum_{i=1}^n \mathbf{E}\{f(g(z + \mathbf{e}_i(X_i - z_i))) - f(g(z))\} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{E} \left\{ (X_i - z_i) f'(g(z)) + \frac{1}{2} (X_i - z_i)^2 f''(g(z)) \right\} + \frac{1}{n} E(f, g, z), \end{aligned} \quad (4.2)$$

where

$$\begin{aligned} E(f, g, z) &= \sum_{i=1}^n \mathbf{E} \left\{ f(g(z + \mathbf{e}_i(X_i - z_i))) - f(g(z)) - (X_i - z_i) f'(g(z)) - \frac{1}{2} (X_i - z_i)^2 f''(g(z)) \right\}. \end{aligned} \quad (4.3)$$

Using the fact that $\mathbf{E}X_i = 0$ for each i and that $\mathbf{E}W^2 = 1$, and absorbing the factor $1/n$ (which plays no real part) into the definition of \tilde{L} , this gives

$$\tilde{L}(f \circ g)(z) = -g(z)f'(g(z)) + \frac{1}{2} \left(1 + \sum_{i=1}^n z_i^2 \right) f''(g(z)) + E(f, g, z). \quad (4.4)$$

This shows that $\tilde{L}(f \circ g)(z)$ is close to $Lf(g(z))$, with the Stein operator given by

$$Lf(w) = -wf'(w) + f''(w),$$

provided that $(\sum_{i=1}^n z_i^2 - 1)$ and $E(f, g, z)$ are small. Note that the Stein operator is the same as that derived earlier for the standard normal distribution, except that the derivative f' now plays the part of f . It can thus be exploited exactly as before. Note also that L is itself the generator of a stationary Markov process \hat{Z} , the Ornstein–Uhlenbeck diffusion.

For the purposes of normal approximation, the Stein identity $\mathbf{E}\{\tilde{L}(f \circ g)(\tilde{Z}(0))\} = 0$, for all suitable f and for g the component sum, as above, yields

$$-\mathbf{E}Lf(W) = \frac{1}{2}\mathbf{E} \left\{ \left(\sum_{i=1}^n X_i^2 - 1 \right) f''(W) \right\} + \mathbf{E}E(f, g, X), \quad (4.5)$$

because of (4.4). If the right hand side of (4.5) can be uniformly bounded for $f = f_h$ and $h \in \mathcal{H}$, a corresponding bound for the accuracy of the normal approximation is obtained. In particular, it is necessary to show that

$$\mathbf{E} \left| \sum_{i=1}^n X_i^2 - 1 \right| \quad \text{and} \quad \mathbf{E}|E(f, g, X)| \quad (4.6)$$

are both small. The latter can be estimated by a Lyapounov term, if $\|f'''\|_\infty$ is finite, but requires more work if this cannot be assumed, and the former can also be expected to be small if the Lyapounov ratio is small.

The concrete arguments above have always involved properties of the solution f_h to the Stein equation $Lf = h$. For the standard normal distribution, Stein (1972) gave bounds for f_h and f'_h in the supremum norm, when h is an indicator function. Indeed, the explicit form of the solution f_h given in (2.3) greatly facilitates such calculations. For other Stein equations, controlling the solutions can be much more difficult. However, for L the generator of a stationary Markov process \hat{Z} with stationary distribution $\mathcal{L}(Z)$, the equation $Lf = h$ is the so-called Poisson equation, and the solution, under appropriate conditions, is the recurrent potential

$$f_h(w) = - \int_0^\infty \mathbf{E}\{h(\hat{Z}(t)) \mid \hat{Z}(0) = w\} dt, \quad (4.7)$$

assuming as always that $\mathbf{E}h(Z) = 0$.

5. Exchangeable pairs.

In his 1986 monograph ‘Approximate Computation of Expectations’, Stein begins by introducing the notion of an exchangeable pair of random variables for the construction of Stein identities. As usual, the aim is to approximate the distribution of some random element W . Suppose that W' is another random variable, defined on the same probability space as W , such that $(W, W') =_d (W', W)$. Then (W, W') is called an exchangeable pair. Hence, if F is any antisymmetric real function on $\mathcal{S} \times \mathcal{S}$, it is immediate that $\mathbf{E}F(W, W') = 0$ (if the expectation exists), and hence also that $\mathbf{E}\{\mathbf{E}^W F(W, W')\} = 0$, where \mathbf{E}^W denotes the conditional expectation given W . If \mathcal{S} is finite, and connected in the sense that w_1 and w_2 are neighbours if $\mathbf{P}[(W, W') = (w_1, w_2)] > 0$, Stein showed that every function g such that $\mathbf{E}g(W) = 0$ can be represented in the form $g(W) = \mathbf{E}^W F(W, W')$, characterizing the distribution $\mathcal{L}(W)$ in terms of conditional expectations of the form $\mathbf{E}^W F(W, W')$.

To convert this statement into a characterizing operator \tilde{L} on \mathcal{X} , one needs to define an appropriate subcollection of antisymmetric functions F . For normal approximation, Stein proposed the functions

$$F(w, w') = (w - w')(f(w) + f(w')), \quad (5.1)$$

for $f \in \mathcal{X}$. If the linear regression condition

$$\mathbf{E}^W W' = (1 - \lambda)W \quad (5.2)$$

is satisfied for some $0 < \lambda < 1$ (already implying that $\mathbf{E}W = 0$), Stein (1986, Theorem III.1) showed that

$$\begin{aligned} |\mathbf{P}[W \leq w] - \Phi(w)| &\leq 2\sqrt{\mathbf{E}\left(1 - \frac{1}{2\lambda}\mathbf{E}^W\{(W' - W)^2\}\right)^2} \\ &\quad + \frac{1}{(2\pi)^{1/4}}\sqrt{\frac{1}{\lambda}\mathbf{E}|W - W'|^3}. \end{aligned} \quad (5.3)$$

Thus, provided that W and W' are typically close to one another, and that the conditional expectation $\mathbf{E}^W\{(W' - W)^2\}$ is concentrated near the value 2λ , the linear regression property implies that the distribution of W is close to standard normal. This theorem, and its refinements (in particular allowing the linear regression condition to be only approximately satisfied), provide a powerful and flexible tool for estimating the accuracy of normal approximation in a wide variety of settings, with complicated dependence structures.

How is this theorem proved? The basic idea is simple. Take the antisymmetric function F from (5.1), giving

$$\mathbf{E}\{(W - W')(f(W) + f(W'))\} = 0. \quad (5.4)$$

This, with liberal use of (5.2), can equivalently be written as

$$\begin{aligned} \mathbf{E}\{Wf(W) - f'(W)\} &= \mathbf{E}\left\{f'(W) \left(\frac{(W - W')^2}{2\lambda} - 1\right)\right\} \\ &\quad + \frac{1}{2\lambda}\mathbf{E}\{(W - W')\{f(W) - f(W') - (W - W')f'(W)\}\}. \end{aligned} \quad (5.5)$$

Note that the left hand side is just $-\mathbf{E}\{Lf(W)\}$, for L the Stein operator for the standard normal distribution. The first term on the right hand side is small if $\mathbf{E}^W\{(W' - W)^2\}/2\lambda$ is concentrated close to 1 (and (5.2) implies that $\mathbf{E}\{(W' - W)^2\} = 2\lambda\mathbf{E}W^2$, so that having $\mathbf{E}W^2$ close to 1 is certainly a good idea), and, for twice differentiable functions f , the second term is plausibly small if $\frac{1}{\lambda}\mathbf{E}|W - W'|^3$ is small. The detailed argument leading to (5.3) is somewhat more complicated, but Stein's method has already accomplished the hard work.

The choice of antisymmetric function $F(w, w') = (w - w')(f(w) + f(w'))$ is not the simplest: one could instead have chosen $F(w, w') = f(w') - f(w)$. This would in similar style give

$$\begin{aligned} 0 &= \mathbf{E}\{f(W') - f(W)\} \\ &= \mathbf{E}\left\{(W' - W)f'(W) + \frac{1}{2}(W' - W)^2f''(W)\right\} + \mathbf{E}E'(f, W, W'), \end{aligned} \quad (5.6)$$

where

$$E'(f, W, W') = f(W') - f(W) - (W' - W)f'(W) - \frac{1}{2}(W' - W)^2f''(W).$$

Using the linear regression condition $\mathbf{E}^W W' = (1 - \lambda)W$ of (5.2), which as above also implies that $\mathbf{E}\{(W' - W)^2\} = 2\lambda\mathbf{E}W^2$, it follows that

$$0 = -\lambda\mathbf{E}\{-Wf'(W) + f''(W)\} + \lambda\mathbf{E}\left\{\left(\frac{\mathbf{E}^W(W' - W)^2}{2\lambda} - 1\right)f''(W)\right\} + \mathbf{E}E'(f, W, W'). \quad (5.7)$$

Cancelling the factor λ , this gives a close parallel to (5.5), but with f' replacing f in the analogy, as was the case when using the generator approach in Section 4. Note that, in this argument, exchangeability is not used: $\mathbf{E}\{f(W') - f(W)\} = 0$ requires only that $W =_d W'$ (Röllin, 2008).

In many examples, the exchangeable pair can be realized as a pair of successive states in a stationary *reversible* Markov chain. In these circumstances, (5.6) can also be interpreted as the Stein identity from the generator approach. This is indeed the case for the sums of independent random variables considered in the previous section. The linear regression condition is satisfied with $\lambda = 1/n$, and expression (4.5) is an exact parallel to (5.7).

The approach using exchangeable pairs has proved extremely useful in many contexts. In the book Diaconis and Holmes (2004), for instance, there are chapters showing how Stein’s method and exchangeable pairs can be effectively exploited in a variety of quite disparate settings. They are used, amongst other things, to reduce variance in simulation experiments (Stein, Diaconis, Holmes & Reinert 2004), to estimate rates of mixing for some ergodic Markov chains (Diaconis 2004), to prove a central limit theorem for the number of descents in a random permutation (Fulman 2004), and to establish the large sample properties of the bootstrap, including a variant of the blockwise bootstrap for some dependent samples (Holmes & Reinert 2004). Other examples that have since been influential include those in Rinott & Rotar (1997). The combination of linear regression condition and the exchangeable pair has recently been effectively developed in the setting of multivariate normal approximation by Reinert & Röllin (2009).

6. The concentration inequality approach.

The Berry–Esseen theorem is framed in terms of the Kolmogorov distance, in which the class \mathcal{H} of test functions h consists of the indicators of half lines. The corresponding functions f_h defined by (2.3) do not have $\|f_h''\|_\infty < \infty$, since their derivatives have a jump. This means that the simple argument given in Section 2 (c.f. Erickson 1974) cannot be used directly to prove normal approximation with respect to this distance. However, the effect of the jump on the expectations to be bounded can be controlled, if it is known that the concentration function of the distribution of W is suitably bounded. This was an essential step in Stein’s original proof of the Berry–Esseen theorem for sums of independent and identically distributed random variables. Astonishingly, his proof of the concentration inequality also made use of his method in a most ingenious way. We now reproduce his proof in detail, in the very simplest context.

Let X_1, X_2, \dots, X_n be a sequence of independent and identically distributed random variables such that $\mathbf{E}X_1 = 0$, $\mathbf{E}X_1^2 = 1/n$ and $\mathbf{E}|X_1|^3 < \infty$; set $\beta = n^{3/2}\mathbf{E}|X_1|^3$. The Berry–Esseen theorem states that, for every real z ,

$$|F_n(z) - \Phi(z)| \leq C\beta n^{-1/2}, \tag{6.1}$$

where F_n is the distribution function of $\sum_{i=1}^n X_i$, Φ is the standard normal distribution function, and C is an absolute constant.

Stein's proof, to be found in Ho & Chen (1978), proceeds as follows. Let $W_n = \sum_{i=1}^n X_i$, $W_{n-1} = \sum_{i=1}^{n-1} X_i$. Then, for any absolutely continuous function f with bounded derivative, we have

$$\mathbf{E}\{W_n f(W_n)\} = \sum_{i=1}^n \mathbf{E} \left\{ X_i f \left(\sum_{j \neq i} X_j + X_i \right) \right\}, \quad (6.2)$$

which by independence, symmetry and the fact that $\mathbf{E}X_n = 0$, gives

$$\begin{aligned} \mathbf{E}\{W_n f(W_n)\} &= n \mathbf{E}[X_n f(W_{n-1} + X_n)] = n \mathbf{E}[X_n (f(W_{n-1} + X_n) - f(W_{n-1}))] \\ &= n \mathbf{E} \left\{ X_n \int_0^{X_n} f'(W_{n-1} + t) dt \right\} \\ &= \mathbf{E} \int_{-\infty}^{\infty} f'(W_{n-1} + t) K(t) dt, \end{aligned} \quad (6.3)$$

where

$$K(t) = n \mathbf{E}X_n [I(0 < t < X_n) - I(X_n < t < 0)]. \quad (6.4)$$

Hence we obtain the identity

$$\mathbf{E}\{W_n f(W_n)\} = \mathbf{E} \int_{-\infty}^{\infty} f'(W_{n-1} + t) K(t) dt. \quad (6.5)$$

Now the function $K(t)$ is nonnegative, satisfying $K(-\infty) = 0$ and $K(\infty) = 0$, and, by setting $f(w) = w$, (6.5) yields

$$\int K(t) dt = \mathbf{E}W_n^2 = 1, \quad (6.6)$$

showing that K is in fact a probability density function; we note also, for future reference, that

$$\int |t| K(t) dt = \beta/2\sqrt{n}. \quad (6.7)$$

Hence, from (6.5), the Stein operator for standard normal approximation applied to W_n can be expressed as

$$\mathbf{E}\{f'(W_n) - W_n f(W_n)\} = \mathbf{E} \int_{-\infty}^{\infty} \{f'(W_{n-1} + X_n) - f'(W_{n-1} + t)\} K(t) dt, \quad (6.8)$$

and it remains only to show that the expression on the right hand side is small.

Clearly, as before, there is a very simple argument to do so if f_h'' is uniformly bounded for our test functions h , but f_h' has a jump of size 1 at a when $h = \mathbf{1}_{(-\infty, a]}$, and its derivative may be big close to a if a is large. Thus it is necessary to show that the random variable W_{n-1} rarely takes a value close enough to a to make the expected difference $\mathbf{E}^{W_{n-1}}|f'(W_{n-1} + X_n) - f'(W_{n-1} + t)|$ appreciable. Note that, from (6.7), the relevant values of t here are of order $1/\sqrt{n}$, and this is also the scale appropriate to values of X_n ; hence it is necessary to be able to show that W_{n-1} has suitably small probability of being close to a on this scale. It turns out that the right hand side of (6.8) can indeed be shown to be of the required order $O(\beta/\sqrt{n})$, by using the following concentration inequality.

Lemma. *For all real a and b such that $a < b$,*

$$\mathbf{P}(a \leq W_{n-1} \leq b) \leq b - a + 2\beta/\sqrt{n}.$$

Proof. We begin with a simple inequality. By (6.7), we have

$$\int_{t > \beta/\sqrt{n}} K(t) dt \leq (n^{1/2}/\beta) \int_{t > \beta/\sqrt{n}} |t|K(t) dt \leq (n^{1/2}/\beta) \int |t|K(t) dt = \frac{1}{2}.$$

This and (6.6) yield

$$\int_{t \leq \beta/\sqrt{n}} K(t) dt = \int K(t) dt - \int_{t > \beta/\sqrt{n}} K(t) dt \geq \frac{1}{2}. \quad (6.9)$$

Now let a and b be two real numbers such that $a < b$. For any real $x > 0$, define

$$g_x(w) = \begin{cases} -\frac{1}{2}(b-a) - x & \text{if } w \leq a - x \\ w - \frac{1}{2}(a+b) & \text{if } a - x \leq w \leq b + x \\ \frac{1}{2}(b-a) + x & \text{if } b + x \leq w. \end{cases} \quad (6.10)$$

Clearly, g_x is an absolutely continuous function, and $g'_x(w) = I(a - x \leq w \leq b + x)$. Taking $f = g_{\beta/\sqrt{n}}$, we have

$$\begin{aligned} \mathbf{E} \int f'(W_{n-1} + t)K(t) dt &= \mathbf{E} \int I(a - \beta/\sqrt{n} \leq W_{n-1} + t \leq b + \beta/\sqrt{n})K(t) dt \\ &\geq \mathbf{E}I(a \leq W_{n-1} \leq b) \int I(|t| \leq \beta/\sqrt{n})K(t) dt. \end{aligned} \quad (6.11)$$

By (6.9), the right hand side is at least $\frac{1}{2}\mathbf{P}(a \leq W_{n-1} \leq b)$. But now, in view of (6.5), this gives

$$\mathbf{P}(a \leq W_{n-1} \leq b) \leq 2\mathbf{E}W_n f(W_n) \leq 2\mathbf{E}|W_n f(W_n)| \leq b - a + 2\beta/\sqrt{n}, \quad (6.12)$$

because $\|f\|_\infty \leq \frac{1}{2}(b-a) + \beta/\sqrt{n}$ from (6.10), and $\mathbf{E}|W_n| \leq (\mathbf{E}W_n^2)^{1/2} = 1$. □

Stein's idea of concentration inequality was extended beyond the case of i.i.d. random variables in Chen (1986). For locally dependent and not necessarily identically distributed random variables X_1, X_2, \dots, X_n and a random vector ζ depending on a relatively small subset of $\{X_1, X_2, \dots, X_n\}$, a randomized concentration inequality can be proved for $P^\zeta(a\zeta \leq W \leq b\zeta)$. This is used to derive good bounds on the difference with respect to Kolmogorov distance between the distribution the standardized sum W of the X_i 's and the standard normal distribution, thereby extending the classical Berry–Esseen theorem. Non-uniform concentration inequalities were also developed in Chen & Shao (2001, 2004), and were used to obtain a non-uniform Berry–Esseen theorem.

If a random variable T can be expressed as $T = W + \Delta$, where the Kolmogorov distance between $\mathcal{L}(W)$ and $\mathcal{L}(Z)$ is known to be small and where Δ is a relatively small random variable which may depend on W , a randomized concentration inequality for $\mathbf{P}(z \leq W \leq z + |\Delta|)$ again gives good bound for the Kolmogorov distance between $\mathcal{L}(W)$ and $\mathcal{L}(Z)$. This idea was used in Chen & Shao (2007) to prove normal approximation for nonlinear statistics, and in Barbour & Chen (2005) for the permutation distribution of matrix correlation statistics.

7. Poisson approximation.

One of the biggest successes of Stein's method has been in Poisson approximation. The classical limit theorem, that $\text{Bi}(n, \lambda/n) \rightarrow \text{Po}(\lambda)$ is taught in most first courses on probability theory, and the proof uses explicit computation of the point probabilities. It was only rather recently (Prohorov, 1953) that a good bound on the difference between the two distributions was given, in the form

$$d_{\text{TV}}(\text{Bi}(n, p), \text{Po}(np)) \leq cp \min\{1, np\}, \quad (7.1)$$

for an explicit constant c , where the total variation distribution $d_{\text{TV}}(P, Q)$ between two probability measures is the supremum of $|P(A) - Q(A)|$ over measurable sets A . Note that, if W and Z are random variables on some \mathcal{S} , then

$$d_{\text{TV}}(\mathcal{L}(W), \mathcal{L}(Z)) = \frac{1}{2} \sup_{h \in \mathcal{H}_0} |\mathbf{E}h(W) - \mathbf{E}h(Z)|, \quad (7.2)$$

where \mathcal{H}_0 is the set of measurable real functions h on \mathcal{S} with $\|h\|_\infty \leq 1$. This representation suggests that Stein's method may be well suited to total variation approximation: its development in the case of Poisson approximation is due to Chen (1975).

The first step is to derive a suitable Stein operator on real valued functions on the non-negative integers \mathbf{Z}_+ . By analogy with the derivation of the Stein operator for standard normal approximation in Section 2, letting $p(j)$ denote the Poisson $\text{Po}(\lambda)$ probability of the point j , we have

$$\begin{aligned} 0 &= \sum_{j=0}^{\infty} \Delta_* \{f(j)p(j)\} \\ &= \sum_{j=0}^{\infty} \{f(j) - jf(j-1)/\lambda\}p(j) \end{aligned} \tag{7.3}$$

for all $f \in \mathcal{F}$, the set of real functions f on \mathbf{Z}_+ such that $f(j)p(j) \rightarrow 0$ as $j \rightarrow \infty$, where $\Delta_* g(j) = g(j) - g(j-1)$ denotes the backward difference operator, and $\Delta_* \{f(0)g(0)\} = f(0)p(0)$. Thus any function h of the form $h(j) = f(j) - jf(j-1)/\lambda$ with $f \in \mathcal{F}$ automatically has $\mathbf{E}h(Z) = 0$, if $Z \sim \text{Po}(\lambda)$. Conversely, for any h such that $\mathbf{E}|h(Z)| < \infty$, the function $h(j) - \mathbf{E}h(Z)$ can be written in the form $f(j) - jf(j-1)/\lambda$, with $f \in \mathcal{F}$. This is because, if $g(j) = f(j) - jf(j-1)/\lambda$, then

$$\sum_{j=0}^J g(j)p(j) = \sum_{j=0}^J \Delta_* \{f(j)p(j)\} = f(J)p(J), \tag{7.4}$$

and, in reverse, defining f by $f(J)p(J) = \sum_{j=0}^J g(j)p(j)$ for any g such that $\mathbf{E}|g(Z)| < \infty$ yields $f(j) - jf(j-1)/\lambda = g(j)$. Hence we can take $f = f_h$ to be defined by

$$\begin{aligned} f_h(J)p(J) &= \sum_{j=0}^J \{h(j) - \mathbf{E}h(Z)\}p(j) \\ &= - \sum_{j=J+1}^{\infty} \{h(j) - \mathbf{E}h(Z)\}p(j) \end{aligned}, \tag{7.5}$$

noting that it is then directly checked that $f_h \in \mathcal{F}$. As before, this allows us to write

$$\mathbf{E}h(W) - \mathbf{E}h(Z) = \mathbf{E}\{f_h(W) - Wf_h(W-1)\}, \tag{7.6}$$

for any random variable W on \mathbf{Z}_+ for which the expectations exist; the value of $f_h(-1)$ can be chosen arbitrarily, since it plays no role. Thus we are led to a Stein operator for the Poisson distribution, usually expressed, by shifting the argument of f and multiplying through by λ , in the equivalent form

$$Lf(w) = \lambda f(w+1) - wf(w), \tag{7.7}$$

for functions f such that $f(j+1)p(j) \rightarrow 0$ as $j \rightarrow \infty$.

To illustrate its use, let X_1, \dots, X_n be independent Bernoulli random variables, with $X_i \sim \text{Be}(p_i)$, and let $W = \sum_{i=1}^n X_i$. If the p_i are not too large, Le Cam (1960) gave an analogue of Prohorov's (1953) approximation for the total variation distance between $\mathcal{L}(W)$ and $\mathcal{L}(Z)$, where $Z \sim \text{Po}(\lambda)$ and $\lambda = \sum_{i=1}^n p_i$: he showed that

$$d_{\text{TV}}(\mathcal{L}(W), \mathcal{L}(Z)) \leq 8 \min(1, \lambda^{-1}) \sum_{i=1}^n p_i^2, \quad (7.8)$$

if $\max_i p_i \leq 1/4$, using an operator approach that relies heavily on independence. To apply Stein's method, we need to bound the right hand side of (7.6). As for normal approximation, the term $\mathbf{E}\{W f_h(W)\}$ breaks up into elements of the form $\mathbf{E}\{X_i f_h(W)\}$, and since $X_i f(W) = X_i f(W_i + 1)$ a.s., where $W_i = W - X_i$, it follows by the independence of X_i and W_i that

$$\mathbf{E}\{X_i f_h(W)\} = p_i \mathbf{E}f_h(W_i + 1); \quad \mathbf{E}\{W f_h(W)\} = \sum_{i=1}^n p_i \mathbf{E}f_h(W_i + 1).$$

This is to be compared with

$$\lambda \mathbf{E}f_h(W + 1) = \sum_{i=1}^n p_i \mathbf{E}f_h(W_i + X_i + 1).$$

Taking the difference, it follows easily that

$$|\mathbf{E}h(W) - \mathbf{E}h(Z)| \leq \sum_{i=1}^n p_i \mathbf{E}|f_h(W_i + X_i + 1) - f_h(W_i + 1)| \leq \sum_{i=1}^n p_i^2 \|\Delta f_h\|_\infty, \quad (7.9)$$

where $\Delta f(j) = f(j+1) - f(j)$ denotes the forward difference. Since $\sup_{h \in \mathcal{H}_0} \|\Delta f_h\|_\infty \leq 2 \min(1, \lambda^{-1})$, it follows immediately that

$$d_{\text{TV}}(\mathcal{L}(W), \mathcal{L}(Z)) \leq \min(1, \lambda^{-1}) \sum_{i=1}^n p_i^2, \quad (7.10)$$

improving the constant in Le Cam's bound, and removing the restriction on the maximum value of the p_i 's. The argument was also refined to establish a *lower* bound for the error, having exactly the same order (Barbour & Hall, 1984).

Of course, the real advance of the method in this context is the ease with which dependence can be accommodated. For instance, if W can be written as $X_i + Y_i + \widetilde{W}_i$, with Y_i a random variable in \mathbf{Z}_+ and \widetilde{W}_i almost independent of X_i , then one can write

$$\begin{aligned} \mathbf{E}\{X_i f_h(W)\} &= \mathbf{E}\{X_i f_h(1 + Y_i + \widetilde{W}_i)\} \\ &= \mathbf{E}\{X_i f_h(1 + \widetilde{W}_i)\} + \mathbf{E}\{X_i (f_h(1 + Y_i + \widetilde{W}_i) - f_h(1 + \widetilde{W}_i))\}, \end{aligned} \quad (7.11)$$

with the first term close to $p_i \mathbf{E}f_h(\widetilde{W}_i + 1)$ and the second bounded by $\mathbf{E}(X_i Y_i) \|\Delta f_h\|_\infty$. Comparing this to

$$p_i \mathbf{E}f_h(W + 1) = p_i \mathbf{E}f_h(X_i + Y_i + \widetilde{W}_i),$$

it follows immediately that

$$\begin{aligned} & |\mathbf{E}h(W) - \mathbf{E}h(Z)| \\ & \leq \sum_{i=1}^n \{p_i^2 + p_i \mathbf{E}Y_i + \mathbf{E}(X_i Y_i)\} \min\{1, \lambda^{-1}\} + \sum_{i=1}^n \mathbf{E}|\mathbf{E}(X_i | \widetilde{W}_i) - p_i| \|f_h\|_\infty, \end{aligned} \quad (7.12)$$

and $\sup_{h \in \mathcal{H}_0} \|f_h\|_\infty \leq 2 \min(1, \lambda^{-1/2})$. This leads to useful Poisson approximation in total variation for a wide variety of sums of dependent Bernoulli random variables.

Poisson approximation for the sum of independent Bernoulli random variables could also have been approached by deriving a Stein identity for W , and then deducing the Stein operator appropriate to the problem. A stationary Markov chain \widetilde{Z} with $\widetilde{Z}(0) = (X_1, \dots, X_n)$ can be constructed exactly as in Section 4, and its generator is

$$\widetilde{L}f(z) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}\{f(z + \mathbf{e}_i(X_i - z_i)) - f(z)\}; \quad (7.13)$$

specializing to functions $f \circ g$ with $g(z) = \sum_{i=1}^n z_i$, this gives, after some manipulation,

$$\begin{aligned} & n\widetilde{L}(f \circ g)(z) \\ & = \lambda\{f(g(z) + 1) - f(g(z))\} - g(z)\{f(g(z)) - f(g(z) - 1)\} + \sum_{i=1}^n z_i p_i (\Delta^2 f)(g(z) - 1). \end{aligned} \quad (7.14)$$

Thus, writing $g(z) = w$, and dropping the factor n , this expression is close to $Lf(w)$, where

$$Lf(w) = \lambda \Delta f(w) - w \Delta f(w - 1), \quad (7.15)$$

the same generator as in (7.7), except that now $\Delta f(w - 1)$ plays the part of $f(w)$. Note that L is the generator of a well-known stationary, reversible Markov process, the immigration–death process with immigration rate λ and with unit *per capita* death rate. Poisson approximation for sums of independent Bernoulli random variables can now be deduced directly from (7.14), since it gives

$$|\mathbf{E}Lf_h(W)| \leq \sum_{i=1}^n p_i^2 \|\Delta^2 f_h\|_\infty, \quad (7.16)$$

which is precisely the bound given in (7.9).

In more complicated circumstances, with dependent summands, it may still sometimes be easier to start with the version of Lf given in (7.7) or in (7.15) and to make direct calculations, rather than to set up a Stein identity for W as the first step. On the other hand, one step of the chain \tilde{Z} in the independent setting would yield an exchangeable pair, and constructing an exchangeable pair also offers many possibilities in dependent settings in which direct calculation may seem unattractive.

The identification of the distribution $\text{Po}(\lambda)$ as that of a stationary immigration–death chain can be very simply extended to Poisson processes. A Poisson point process Π with intensity measure λ on a space \mathcal{U} has distribution $\text{PP}(\lambda)$, the equilibrium distribution of a spatial immigration–death process $\tilde{\Pi}$ with immigration intensity measure λ over \mathcal{U} and with unit *per capita* death rate. Letting δ_u denote the point mass at u , ξ a Radon measure on \mathcal{U} and f a bounded real function on the space of such Radon measures, the generator of the process $\tilde{\Pi}$ can be expressed as

$$Lf(\xi) = \int_{\mathcal{U}} \{f(\xi + \delta_u) - f(\xi)\} \lambda(du) + \int_{\mathcal{U}} \{f(\xi - \delta_u) - f(\xi)\} \xi(du). \quad (7.17)$$

The corresponding Stein equation $Lf = h$, for bounded functions h for which $\mathbf{E}h(\Pi) = 0$, can be solved in the form

$$f_h(\xi) = - \int_0^\infty \mathbf{E}\{h(\tilde{\Pi}(t)) \mid \tilde{\Pi}(0) = \xi\} dt, \quad (7.18)$$

as an example of (A4.9), and coupling arguments can be used to bound the first and second differences of f_h . However, at least for the class of test functions appropriate to total variation approximation, the bounds on these differences do not in general decrease with increasing $\lambda(\mathcal{U})$, in contrast to the case of the Poisson distribution $\text{Po}(\lambda)$, which corresponds to \mathcal{U} consisting of a single point. Much more on this topic can be found in Barbour & Brown (1992) and in Barbour, Holst & Janson (1992, Chapter X).

The simplest application is once again in the independent setting. With X_1, \dots, X_n as defined earlier in the section, let $\mathcal{U} = \{1, 2, \dots, n\}$, and set $\Xi = \sum_{i=1}^n X_i \delta_i$. To see how close the distribution of Ξ is to that of a Poisson process Π with intensity $\lambda = \sum_{i=1}^n p_i \delta_i$, we need to examine $\mathbf{E}Lf(\Xi)$. Now, much as for the random variable case,

$$\mathbf{E} \left\{ \int_{\mathcal{U}} \{f(\Xi - \delta_u) - f(\Xi)\} \Xi(du) \right\} = - \sum_{i=1}^n \mathbf{E}(X_i \{f(\Xi_i + \delta_i) - f(\Xi_i)\})$$

where $\Xi_i = \Xi - X_i \delta_i$ is independent of X_i , and

$$\mathbf{E} \left\{ \int_{\mathcal{U}} \{f(\Xi + \delta_u) - f(\Xi)\} \lambda(du) \right\} = \sum_{i=1}^n p_i \mathbf{E}\{f(\Xi + \delta_i) - f(\Xi)\}.$$

Adding the two yields

$$\mathbf{E}L f(\Xi) = \sum_{i=1}^n p_i \{f(\Xi_i + \delta_i + X_i \delta_i) - f(\Xi_i + X_i \delta_i) - f(\Xi_i + \delta_i) + f(\Xi_i)\}, \quad (7.19)$$

giving an immediate bound

$$|\mathbf{E}h(\Xi) - \mathbf{E}h(\Pi)| = |\mathbf{E}L f_h(\Xi)| \leq \sum_{i=1}^n p_i^2 \sup_{u, \xi} |f_h(\xi + 2\delta_u) - 2f_h(\xi + \delta_u) + f_h(\xi)|, \quad (7.20)$$

and the second difference can be bounded starting from (7.18); for total variation, the supremum is just 2. Note that the argument is in every sense equivalent to that for random variables.

For dependent random variables X_i , there is again a parallel. Suppose now that, for each i , we can write $\Xi = X_i \delta_i + H_i + \tilde{\Xi}_i$, where H_i is a Radon measure and X_i and $\tilde{\Xi}_i$ are almost independent. Then

$$\begin{aligned} \mathbf{E} \left\{ \int_{\mathcal{U}} \{f(\Xi - \delta_u) - f(\Xi)\} \Xi(du) \right\} &= - \sum_{i=1}^n \mathbf{E}(X_i \{f(\tilde{\Xi}_i + \delta_i) - f(\tilde{\Xi}_i)\}) \\ &\quad - \sum_{i=1}^n \mathbf{E}(X_i \{f(H_i + \tilde{\Xi}_i) - f(H_i + \tilde{\Xi}_i + \delta_i) - f(\tilde{\Xi}_i + \delta_i) + f(\tilde{\Xi}_i)\}), \end{aligned} \quad (7.21)$$

with the first term, much as for (7.11), close to $\sum_{i=1}^n p_i \mathbf{E} \Delta_i f_h(\tilde{\Xi}_i)$, and the second bounded by $\mathbf{E}(X_i Y_i) \sup_{i,j} \|\Delta_{ij} f_h\|_\infty$: here, Δ_i denotes the first forward difference in the direction i , and Δ_{ij} the corresponding second differences, and $Y_i = H_i(\mathcal{U})$. This is to be added to

$$\begin{aligned} &\mathbf{E} \left\{ \int_{\mathcal{U}} \{f(\Xi + \delta_u) - f(\Xi)\} \lambda(du) \right\} \\ &= \sum_{i=1}^n p_i \mathbf{E} \{f(X_i \delta_i + H_i + \tilde{\Xi}_i + \delta_i) - f(X_i \delta_i + H_i + \tilde{\Xi}_i)\} \\ &= \sum_{i=1}^n p_i \mathbf{E} \Delta_i f(\tilde{\Xi}_i) + \sum_{i=1}^n p_i \mathbf{E} \{ \Delta_i f(X_i \delta_i + H_i + \tilde{\Xi}_i) - \Delta_i f(\tilde{\Xi}_i) \}, \end{aligned} \quad (7.22)$$

whose second term is bounded by

$$\sum_{i=1}^n p_i \{p_i + \mathbf{E}Y_i\} \sup_{i,j} \|\Delta_{ij} f_h\|_\infty. \quad (7.23)$$

Thus, in parallel to (7.12), we obtain

$$|\mathbf{E}h(\Xi) - \mathbf{E}h(\Pi)| \leq \sum_{i=1}^n \{p_i^2 + p_i \mathbf{E}Y_i + \mathbf{E}(X_i Y_i)\} c_1(h) + \sum_{i=1}^n \mathbf{E} |\mathbf{E}(X_i | \tilde{\Xi}_i) - p_i| c_0(h); \quad (7.24)$$

again, in contrast with the situation for random variables, the best general bounds for $c_0(h)$ and $c_1(h)$ for functions h with $\|h\|_\infty \leq 1$ are each 2, so that there is in general no decrease when $\lambda(\mathcal{U})$ increases.

The total variation bound for the approximation of $\mathcal{L}(\Xi)$ by PP (λ) implied by (7.24) was demonstrated in a wide variety of examples in Arratia, Goldstein & Gordon (1989, 1990). They also used the bound as a basis for compound Poisson approximation, by showing that random variables that might be expected to have approximately such a distribution can frequently be represented as functions of Poisson processes — a clump of observations of size j at point x is associated with the point mass $\delta_{(x,j)}$, and \mathcal{U} consists of the set of such pairs. Their technique turns out to be widely applicable in practice. Compound Poisson approximation can also be approached directly, using a Stein operator derived in Barbour, Chen & Loh (1992) and the method for exploiting it derived in Roos (1994). For this operator, there are difficulties in bounding the solutions of the corresponding Stein equation — see Barbour & Utev (1998, 1999) and Barbour & Xia (1999). The expository paper of Barbour & Chryssaphinou (2001) gives a discussion of the various approaches to compound Poisson approximation using Stein’s method, and a more detailed overview of Poisson approximation can also be found in Barbour (2001).

8. *Ramifications.*

This article is intended more as a personal view of the extraordinary impact that Stein’s method has had, rather than as an exhaustive survey; the subject is just too large. Much more information is to be found in the books and expository articles that have already been written on Stein’s method. The primary source for the general method is in those by Stein himself (1986, 1992). For Poisson approximation and related themes, the papers of Arratia, Goldstein & Gordon (1989, 1990) and the book by Barbour, Holst & Janson (1992) are good starting points. The Diaconis & Holmes (2004) collection has lots of ideas concerned with the use of exchangeable pairs. The tutorial lectures in Barbour & Chen (2005a) also furnish a digestible introduction to many aspects of Stein’s method, and the accompanying workshop proceedings (Barbour & Chen 2005b) illustrate how wide the scope of Stein’s method has become.

This latter aspect, the enormous breadth of research that has developed out of Stein’s ideas, should nonetheless be briefly emphasized. For instance, the method has been applied in the context of approximation by distributions unrelated to either Poisson or normal. Luk (1994) studied approximation by Gamma distributions, Loh (1992) by the multinomial distribution, Peköz (1996) by the geometric distribution, Brown & Xia (2001) by the equilibrium distribution of a birth and death process, and Götze & Tikhomirov (2005) by

the Wigner semi-circle law. This latter paper is but one of a number of applications of Stein’s method in random matrix theory: for example, there is the recent work of Meckes (2008) giving sharp rates for normal approximation to the distribution of linear functions of random orthogonal and unitary matrices, Fulman’s (2009) normal approximation to the distribution of the trace of a random matrix from a compact Lie group, and Chatterjee’s (2009) results discussed below. Gaussian and Poisson process approximation by Stein’s method have already been briefly mentioned; Reinert (1995) showed how to use the method to prove weak laws of large numbers for empirical processes.

Another fascinating development coming from the original normal and Poisson applications is the relationship between the method and certain biasing constructions. The ‘coupling’ approach to Poisson approximation can be interpreted in terms of the size biasing characterization of the Poisson distribution. Size biasing can also be naturally introduced into Stein’s method for the normal distribution, when the random variables under consideration are non-negative (Goldstein & Rinott, 1996). In Goldstein & Reinert (1997, 2005), a new biasing construction, known as zero biasing, was introduced. The construction associates a zero biased distribution to any distribution with zero mean and finite variance, in a way that fits very neatly with the Stein operators for normal approximation in one and higher dimensions.

An area in which Stein’s method has proved extremely fruitful is that of random geometrical graphs, in which points in space are taken to be neighbours if they are close in an appropriate sense. Such models arise very naturally in many branches of statistics. Starting from a paper of Avram & Bertsimas (1993), itself based on Baldi & Rinott (1989), Stein’s method has been used in a systematic way to prove many theorems about normal approximation to the characteristics of such graphs; see, for example, Penrose & Yukich (2005) for more details.

There has recently been an explosion of interest in the application of Stein’s method to functionals of Gaussian processes and fields. Using the method, combined with Malliavin calculus, it has been possible to generalize, refine and unify many central and non-central limit theorems for multiple Wiener–Itô integrals. A first step in this direction is to be found in Nourdin & Peccati (2009a), sharpened and refined in Nourdin & Peccati (2009b) to obtain one term Edgeworth expansions. Their approach has a wide range of application, and it is too soon for a complete picture of the possible developments to be discerned.

Stein’s (1986) monograph contains a chapter on large deviations, which can be seen as work in progress. This project has now been successfully developed by Chatterjee (2007a), who is able to prove strong concentration of measure inequalities (to distinguish them from bounds on the concentration function) using the Stein approach, in a variety of interesting

settings. This is but one of his recent impressive results. In Chatterjee (2009), he presents a rather general technique for proving central limit approximations for the distributions of linear statistics derived from high dimensional random matrices, using Stein’s method and a notion of second order Poincaré inequalities. One of the classes of matrices that he studies is that of random Gaussian Toeplitz matrices; for these matrices, a central limit theorem involving the spectrum is proved, even though the limiting formula for the associated variance is not known. For real-valued functionals $f(X)$ of independent random variables $X = (X_1, X_2, \dots, X_n)$, an abstract theorem bounding the distance between the distribution of $f(X)$ and the appropriate normal distribution is proved in Chatterjee (2008). The main part of the bound, a discrete version of his second-order Poincaré inequality, is expressed in terms of the variance of a random variable T , which is constructed using an independent copy X' of X . In the case of ‘local’ functions f , this variance is shown to be relatively accessible, and the approximation to have a number of interesting applications. He has also proved an approximation, with error bounds, to the distribution of the local field in the high temperature phase of the Sherrington–Kirkpatrick spin glass model: the approximation is by means of a two part mixture of normal distributions (Chatterjee, 2010).

It is not only the method itself that is broad in its scope: the range of application of Stein’s method is equally impressive. Stein (1986) gives applications to counting Latin rectangles, to random allocations, to the binary expansion of a random integer and to isolated trees in a Bernoulli random graph. Other applications include the analysis of molecular sequences (Arratia, Gordon & Waterman 1990; Neuhauser 1994), extreme value theory (Smith 1988), reliability theory (Godbole 1993), random fields (Takahata 1983), card shuffling (Fulman 2005), the eigenfunctions of the Laplacian on a manifold (Meckes 2009), logarithmic combinatorial structures (Arratia, Barbour & Tavaré 2003) and scan statistics (Glaz, Naus, Roos & Wallenstein 1994), among many others.

Stein’s method has emerged as a flexible and powerful tool for proving probability approximations, and has stimulated research in many different directions. Diaconis & Holmes observe that ‘for all these virtues, it still seems impossible to give a brief, understandable explanation of the essence of Stein’s method’. We entirely agree.

References.

- R. Arratia, A. D. Barbour & S. Tavaré (2003) *Logarithmic combinatorial structures*. European Math. Soc. Press, Zürich.
- R. Arratia, L. Goldstein & L. Gordon (1989) Two moments suffice for Poisson approximations: the Chen–Stein method. *Ann. Probab.*, **17**, 9–25.
- R. Arratia, L. Goldstein & L. Gordon (1990) Poisson approximation and the Chen–Stein method. *Stat. Science* **5**, 403–434.
- R. Arratia, L. Gordon & M. S. Waterman (1990) The Erdős–Rényi law in distribution, for coin tossing and sequence matching. *Annals of Statistics*, **18**, 539–70.
- F. Avram & D. Bertsimas (1993) On central limit theorems in geometrical probability. *Ann. Appl. Probab.* **3**, 1033–1046.
- P. Baldi & Y. Rinott (1989) On normal approximations of distributions in terms of dependency graphs. *Ann. Probab.* **17**, 1646–1650.
- P. Baldi, Y. Rinott & C. Stein (1989) A normal approximation for the number of local maxima of a random function on a graph. In: *Probability, Statistics and Mathematics*, Eds T. W. Anderson, K. B. Athreya and D. L. Iglehart. Academic Press, New York, pp. 59–81.
- A. D. Barbour (1988) Stein’s method and Poisson process convergence. *J. Appl. Probab.* **25(A)**, 175–184.
- A. D. Barbour (1990) Stein’s method for diffusion approximations. *Prob. Theory Rel. Fields* **84**, 297–322.
- A. D. Barbour (2001) Topics in Poisson approximation. In: *Handbook of Statistics* **19**, Eds D. N. Shanbhag and C. R. Rao, pp. 79–115.
- A. D. Barbour and T. C. Brown (1992) Stein’s method and point process approximation. *Stoch. Procs Applics* **43**, 9–31.
- A. D. Barbour & L. H. Y. Chen (Eds.) (2005a) *An introduction to Stein’s method*. IMS Lecture Note Series Volume 4, World Scientific Press, Singapore.
- A. D. Barbour & L. H. Y. Chen (Eds.) (2005b) *Stein’s method and applications*. IMS Lecture Note Series Volume 5, World Scientific Press, Singapore.
- A. D. Barbour & L. H. Y. Chen (2005c) The permutation distribution of matrix correlation statistics. In: *Stein’s method and applications*, Eds. A. D. Barbour & L. H. Y. Chen, IMS Lecture Note Series Volume 5, World Scientific Press, Singapore, pp. 223–246.

- A. D. Barbour, L. H. Y. Chen and W.-L. Loh (1992) Compound Poisson approximation for nonnegative random variables via Stein’s method. *Ann. Probab.* **20**, 1843–1866.
- A. D. Barbour & O. Chryssaphinou (2001) Compound Poisson approximation: a user’s guide. *Ann. Appl. Probab.* **11**, 964–1002.
- A. D. Barbour & P. Hall(1984) On the rate of Poisson convergence. *Math. Proc. Cam. Phil. Soc.* **95**, 473–480.
- A. D. Barbour, L. Holst and S. Janson (1992) *Poisson Approximation*. Oxford University Press.
- A. D. Barbour & S. Utev (1998) Solving the Stein Equation in compound Poisson approximation. *Adv. Appl. Probab.* **30**, 449–475.
- A. D. Barbour & S. Utev (1999) Compound Poisson approximation in total variation. *Stoch. Procs. Applics.* **82**, 89–125.
- A. D. Barbour & A. Xia (1999) Poisson perturbations. *ESAIM: P&S* **3**, 131–150.
- E. Bolthausen (1984) An estimate of the remainder in a combinatorial central limit theorem. *Z. Wahrscheinlichkeitstheorie verw. Geb.* **66**, 379–386.
- T. C. Brown & A. Xia (2001) Stein’s method and birth-death processes. *Ann. Probab.* **29**, 1373–1403.
- S. Chatterjee (2007) Stein’s method for concentration inequalities. *Probab. Theory Rel. Fields* **138**, 305–321.
- S. Chatterjee (2008) A new method of normal approximation. *Ann. Probab.* **36**, 1584–1610.
- S. Chatterjee (2009) Fluctuations of eigenvalues and second order Poincaré inequalities. *Probab. Theory Rel. Fields* **143**, 1–40.
- S. Chatterjee (2010) Spin glasses and Stein’s method. *Probab. Theory Rel. Fields* (to appear)
- L. H. Y. Chen (1975) Poisson approximation for dependent trials. *Ann. Probab.* **3**, 534–545.
- L. H. Y. Chen (1986) The rate of convergence in a central limit theorem for dependent random variables with arbitrary index set. IMA Preprint Series #243, Univ. Minnesota.
- L. H. Y. Chen & Q.-M. Shao (2001) A non-uniform Berry–Esseen bound via Stein’s method. *Prob. Theory Rel. Fields* **120**, 236–254.
- L. H. Y. Chen & Q.-M. Shao (2004) Normal approximation under local dependence. *Ann. Probab.* **32**, 1985–2028.

- L. H. Y. Chen & Q.-M. Shao (2007) Normal approximation for nonlinear statistics using a concentration inequality approach. *Bernoulli* **13**, 581–599.
- P. Diaconis (2004) Stein’s method for Markov chains: first examples. In: *Stein’s method: expository lectures and applications*, Eds P. Diaconis & S. Holmes, IMS Lecture Notes **46**, pp. 27–44.
- P. Diaconis & S. Holmes (Eds.) (2004) *Stein’s method: expository lectures and applications*. IMS Lecture Notes Vol. 46, Beachwood, Ohio.
- R. V. Erickson (1974) L_1 bounds for asymptotic normality of m -dependent sums using Stein’s technique. *Ann. Probab.* **2**, 522–529.
- J. Fulman (2004) Stein’s method and non-reversible Markov chains. In: *Stein’s method: expository lectures and applications*, Eds P. Diaconis & S. Holmes, IMS Lecture Notes **46**, pp. 69–78.
- J. Fulman (2005) Stein’s method and descents after riffle shuffles. *Electron. J. Probab.* **10**, 901–924.
- J. Fulman (2009) Stein’s method and characters of compact Lie groups. *Comm. Math. Phys.* **288**, 1181–1201.
- J. Glaz, J. Naus, M. Roos & S. Wallenstein (1994) Poisson approximations for distribution and moments of ordered m -spacings. *J. Appl. Probab.* **31A**, 271–281.
- A. P. Godbole (1993) Approximate reliabilities of m -consecutive- k -out-of- n : Failure systems. *Statist. Sinica* **3**, 321–328.
- F. Götze (1991) On the rate of convergence in the multivariate CLT. *Ann. Probab.* **19**, 724–739.
- F. Götze & A. N. Tikhomirov (2005) Limit theorems for spectra of random matrices with martingale structure. In: *Stein’s method and applications*, Eds. A. D. Barbour & L. H. Y. Chen, IMS Lecture Note Series Volume 5, World Scientific Press, Singapore, pp. 181–194.
- L. Goldstein & Y. Rinott (1996) Multivariate normal approximations by Stein’s method and size bias couplings. *J. Appl. Probab.* **33**, 1–17.
- L. Goldstein & G. Reinert (1997) Stein’s method and the zero-bias transformation with application to simple random sampling. *Ann. Appl. Probab.* **7**, 935–952.
- L. Goldstein & G. Reinert (2005) Zero biasing in one and higher dimensions, and applications. In: *Stein’s method and applications*, Eds. A. D. Barbour & L. H. Y. Chen, IMS Lecture Note Series Volume 5, World Scientific Press, Singapore, pp. 1–18.

- P. Hall and A. D. Barbour (1984) On the rate of Poisson convergence. *Math. Proc. Cam. Phil. Soc.* **95**, 473–480.
- S. T. Ho & L. H. Y. Chen (1978) An L_p bound for the remainder in a combinatorial central limit theorem. *Ann. Probab.* **6**, 231–249.
- S. Holmes & G. Reinert (2004) Stein’s method for the bootstrap. In: *Stein’s method: expository lectures and applications*, Eds P. Diaconis & S. Holmes, IMS Lecture Notes **46**, pp. 95–133.
- W.–L. Loh (1992) Stein’s method and multinomial approximation. *Ann. Appl. Probab.* **2**, 536–554.
- H. M. Luk (1994) *Stein’s method for the gamma distribution and related statistical applications*. Ph. D. Thesis, Univ. Southern California.
- E. Meckes (2008) Linear functions on the classical matrix groups. *Trans. Amer. Math. Soc.* **360**, 5355–5366.
- E. Meckes (2009) On the approximate normality of eigenfunctions of the Laplacian. *Trans. Amer. Math. Soc.* **361**, 5377–5399.
- C. Neuhauser (1994) A Poisson approximation theorem for sequence comparisons with insertions and deletions. *Ann. Statist.* **22**, 1603–1629.
- I. Nourdin & G. Peccati (2009a) Stein’s method on Wiener chaos. *Probab. Theory Rel. Fields* **145**, 75–118.
- I. Nourdin & G. Peccati (2009b) Stein’s method and exact Berry-Esseen asymptotics for functionals of Gaussian fields. *Ann. Probab.* **37**, 2231–2261.
- E. Peköz (1996) Stein’s method for geometric approximation. *J. Appl. Probab.* **33**, 707–713.
- M. D. Penrose & J. E. Yukich (2005) Normal approximation in geometric probability. In: *Stein’s method and applications*, Eds. A. D. Barbour & L. H. Y. Chen, IMS Lecture Note Series Volume 5, World Scientific Press, Singapore, pp. 37–58.
- Ju. V. Prohorov (1953) Asymptotic behaviour of the binomial distribution. *Uspekhi Math. Nauk* **83**, 135–43.
- G. Reinert (1995) A weak law of large numbers for empirical measures via Stein’s method. *Ann. Probab.* **23**, 334–354.
- G. Reinert & A. Röllin (2009) Multivariate normal approximation with Stein’s method of exchangeable pairs under a general linearity condition. *Ann. Probab.* **37**, 2150–2173.

- Y. Rinott & V. Rotar (1996) A multivariate CLT for local dependence with $n^{-1/2} \log n$ rate and applications to multivariate graph related statistics. *J. Multivariate Anal.* **56**, 333–350.
- Y. Rinott & V. Rotar (1997) On coupling constructions and rates in the CLT for dependent summands with applications to the antivoter model and weighted U -statistics. *Ann. Appl. Probab.* **7**, 1080–1105.
- A. Röllin (2008) A note on the exchangeability condition in Stein’s method. *Statist. Probab. Lett.* **78**, 1800–1806.
- M. Roos (1994) Stein’s method for compound Poisson approximation: the local approach. *Ann. Appl. Probab.* **4**, 1177–1187.
- R. L. Smith (1988) Extreme value theory for dependent sequences via the Stein–Chen method of Poisson approximation. *Stoch. Procs. Applics* **30**, 317–327.
- C. Stein (1972) A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. *Proc. Sixth Berkeley Symp. Math. Statist. Probab.* **2**, 583–602.
- C. Stein (1986) *Approximate computation of expectations*. IMS Lecture Notes Vol. 7, Hayward, Calif.
- C. Stein (1992) A way of using auxiliary randomization. In: *Probability Theory*, Eds L. H. Y. Chen, K. P. Choi, K. Hu and J.–H. Lou, W. de Gruyter, Berlin, pp. 159–180.
- C. Stein, P. Diaconis, S. Holmes & G. Reinert (2004) Use of exchangeable pairs in the analysis of simulations. In: *Stein’s method: expository lectures and applications*, Eds P. Diaconis & S. Holmes, IMS Lecture Notes **46**, pp. 1–26.
- H. Takahata (1983) On the rates in the central limit theorem for weakly dependent random fields. *Z. Wahrscheinlichkeitstheorie verw. Geb.* **64**, 445–456.