

Lecture Notes in  
STEIN'S METHOD

Martin Raič

Department of Statistics and Applied Probability  
National University of Singapore

Last change: April 18, 2020

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	Motivation and outline . . . . .	4
1.2	Basic examples of Stein operators . . . . .	6
<b>2</b>	<b>Poisson approximation</b>	<b>11</b>
2.1	Independent trials . . . . .	11
2.2	Solution to the Stein equation . . . . .	13
2.3	Abstract results . . . . .	16
2.3.1	Coupling . . . . .	16
2.3.2	Local dependence and decompositions . . . . .	18
2.4	Applications . . . . .	19
2.4.1	Pattern matching . . . . .	19
2.4.2	Birthday problems . . . . .	20
2.4.3	Random permutations . . . . .	21
2.4.4	Occupancy problems . . . . .	25
<b>3</b>	<b>Normal approximation</b>	<b>32</b>
3.1	Decomposable random variables . . . . .	32
3.2	Solution to the Stein equation . . . . .	38
3.3	Applications . . . . .	43
3.3.1	Local dependence and $U$ -statistics . . . . .	43
3.3.2	Random permutations . . . . .	46
3.4	The Berry–Esseen theorem . . . . .	51
<b>A</b>	<b>Convergence of probability measures</b>	<b>59</b>

A.1	Metrics based on test functions . . . . .	59
A.2	The Wasserstein metric . . . . .	61
A.3	More on the Kolmogorov metric . . . . .	64
A.4	Weak topologies . . . . .	65
A.5	Change of class of test functions . . . . .	70
<b>B</b>	<b>Some real analysis</b>	<b>75</b>
B.1	Differentiation of absolutely continuous functions . . . . .	75
B.2	Functions with bounded variation . . . . .	76
B.3	The Riemann–Stieltjes integral . . . . .	78
B.4	Signed measures . . . . .	81
<b>C</b>	<b>On the Mills ratio</b>	<b>84</b>
C.1	Basic properties . . . . .	84
C.2	Repeated integrals of the Gaussian density . . . . .	86

# Chapter 1

## Introduction

### 1.1 Motivation and outline

Stein's method serves as a powerful tool for estimating the error in approximating complicated probability distributions by more tractable ones. It was introduced in 1970 by Charles Stein (1920–2016) [31], first for assessing the error in the central limit theorem. Recall that, roughly speaking, if  $X_1, X_2, \dots, X_n$  are independent and identically distributed random variables with mean  $\mu_1$  and variance  $\sigma_1^2$ , then:

$$X_1 + X_2 + \dots + X_n \approx \mathcal{N}(n\mu_1, n\sigma_1^2)$$

or equivalently,

$$\frac{X_1 + X_2 + \dots + X_n - n\mu_1}{\sigma_1\sqrt{n}} \approx \mathcal{N}(0, 1).$$

Keeping  $\mu_1$  and  $\sigma_1^2$  fixed, it is known that the larger  $n$ , the smaller is the error in the preceding approximation. Although the error was estimated well before Stein, most prominently by Berry [7] and Esseen [16], Stein's method allows for numerous generalizations where classical approaches, such as characteristic functions, seem not to work. In particular, Stein's method succeeds to cope successfully with many sorts of dependence.

Stein's method has been adapted to numerous other approximations. The first such modification seems to be due to Louis H Y Chen [9], who adapted it to Poisson approximation. Recall that if  $\lambda > 0$  and  $X_1, X_2, \dots$  are independent and identically distributed random variables following the Bernoulli distribution  $\text{Be}(\lambda/n) = \begin{pmatrix} 0 & 1 \\ 1 - \lambda/n & \lambda/n \end{pmatrix}$ , then:

$$X_1 + X_2 + \dots + X_n \approx \text{Po}(\lambda).$$

There are several surveys of Stein's method. Stein summarizes his work in his book [32]. For more modern surveys, see Barbour [2], and Barbour and Chen [3, 4].

The idea behind Stein's method could be viewed as follows: let  $\mathcal{M}$  be a vector space of certain signed measures containing a tractable probability measure  $\nu$ . To approximate

other (presumably complicated) measures  $\mu \in \mathcal{M}$  by  $\nu$ , we first find a linear operator  $\mathcal{B}: \mathcal{M} \rightarrow \mathcal{Y}$  which annihilates  $\nu$ , i. e.,

$$\mathcal{B}\nu = 0;$$

here,  $\mathcal{Y}$  is another vector space. We strive to measure the size of the error  $\mu - \nu$  in terms of the size of  $\mathcal{B}\mu$ .

The equality  $\mathcal{B}\nu = 0$  can be rewritten as  $\nu \in \ker \mathcal{B}$ . If  $\mu - \nu$  can be estimated in terms of  $\mathcal{B}\mu$ , then  $\mathcal{B}$  should separate probability measures sufficiently well, so that it is reasonable to assume that  $\ker \mathcal{B} = \text{span}(\{\nu\})$ . Denoting

$$\langle h, \rho \rangle := \int h \, d\rho \tag{1.1.1}$$

(provided that  $h \in L^1(|\rho|)$ ), observe that the operator

$$\mathcal{Q}\rho := \rho - \langle 1, \rho \rangle \nu$$

has the very same kernel. Therefore, looking algebraically, there exists a linear operator  $\mathcal{T}: \mathcal{Y} \rightarrow \mathcal{M}$ , such that  $\mathcal{T}\mathcal{B} = \mathcal{Q}$ :

$$\begin{array}{ccc} \mathcal{M} & \xrightarrow{\mathcal{Q}} & \mathcal{M} \\ \mathcal{B} \downarrow & \nearrow \mathcal{T} & \\ \mathcal{Y} & & \end{array}$$

**Remark 1.1.1.** The operator  $\mathcal{Q}\rho := \rho - \langle 1, \rho \rangle \nu$  is not the only option, but in the present notes, we shall not consider other possibilities.

If  $\mu$  is a probability measure, then  $\mathcal{Q}\mu = \mathcal{T}\mathcal{B}\mu = \mu - \nu$ . Therefore, in order to bound the size of  $\mu - \nu$ , we can first bound the size of  $\mathcal{B}\mu$  and then the size of  $\mathcal{T}\mathcal{B}\mu$ .

By size, we here mean a suitable seminorm. Seminorms are often based on a certain class of test functions  $\mathcal{H}_1$ :

$$\|\rho\| = \sup_{h \in \mathcal{H}_1} |\langle h, \rho \rangle|$$

(provided that  $h \in L^1(|\rho|)$  for all  $h \in \mathcal{H}_1$ ): see Appendix A.

Suppose that  $\mathcal{B}$  is dual to some operator  $\mathcal{A}$ . That is, suppose that:

- $\mathcal{H}_1 \subseteq \mathcal{H}$ , where  $\mathcal{H}$  is another vector space, such that  $h \in L^1(|\rho|)$  for all  $h \in \mathcal{H}$ ;
- there is another vector space  $\mathcal{F}$  and a bilinear function from  $\mathcal{F} \times \mathcal{Y}$  to  $\mathbb{R}$  again denoted by  $\langle \cdot, \cdot \rangle$ ;
- $\mathcal{A}: \mathcal{F} \rightarrow \mathcal{H}$  is such that  $\langle f, \mathcal{B}\rho \rangle = \langle \mathcal{A}f, \rho \rangle$  for all  $f \in \mathcal{F}$  and  $\rho \in \mathcal{M}$ .

**Remark 1.1.2.** In practice,  $\mathcal{F}$  is a space of functions, though it could in general be any vector space.

Then, of course,

$$\langle \mathcal{A}f, \nu \rangle = \langle f, \mathcal{B}\nu \rangle = 0.$$

Seeking  $\mathcal{T}$  as being dual to some operator  $\mathcal{S}: \mathcal{H} \rightarrow \mathcal{F}$ , observe that

$$\langle h, \mu - \nu \rangle = \langle h, \mathcal{T}\mathcal{B}\mu \rangle = \langle \mathcal{A}\mathcal{S}h, \mu \rangle$$

and, on the other hand,

$$\langle h, \mu - \nu \rangle = \langle h - \langle h, \nu \rangle 1, \mu \rangle.$$

Putting all together, we obtain

**Outline of Stein's method.** To estimate  $\langle h, \mu - \nu \rangle$  for all  $h \in \mathcal{H}$ :

1. Find a vector space  $\mathcal{F}$  and a linear operator  $\mathcal{A}: \mathcal{F} \rightarrow \mathcal{H}$ , such that  $\langle \mathcal{A}f, \nu \rangle = 0$  for all  $f \in \mathcal{F}$ . The operator  $\mathcal{A}$  is called *Stein operator*.
2. For each  $h \in \mathcal{H}$ , find  $f \in \mathcal{F}$  solving the *Stein equation*

$$\mathcal{A}f = h - \langle h, \nu \rangle 1.$$

Then we have

$$\langle h, \mu - \nu \rangle = \langle \mathcal{A}f, \mu \rangle.$$

3. Efficiently bound the right hand side of the preceding equation, which is called *Stein expectation*.

**Remark 1.1.3.** Once we have found  $\mathcal{A}$ , we can drop the assumption that it is dual to some operator  $\mathcal{B}$ . This allows us to extend the space  $\mathcal{M}$ , i. e., the family of measures which can be approximated.

## 1.2 Basic examples of Stein operators

There is no unique way of finding a Stein operator. Indeed, there might be different Stein operators for the same approximating distribution  $\nu$ . Below, we give two basic methods:

1. We derive the operator from the *approximating distribution*  $\nu$  (provided that we already have one in mind). Typically we first derive the operator  $\mathcal{B}$  by comparing consequent point probabilities (for discrete distributions) or by differentiating the density (for continuous distributions). Then observe that  $\mathcal{B}$  is dual to some  $\mathcal{A}$ .
2. We derive the operator from the probability distribution  $\mu$ , which is *to be approximated*. Let  $W$  be a random variable with distribution  $\mu$ . Typically, we perturb  $W$  to  $W'$ , which is defined on the same probability space and follows the same distribution  $\mu$ . Often,  $W$  and  $W'$  are exchangeable, but in general, this is not necessary. Next, if  $\mathcal{G}$  is a linear operator, such that

$$\mathbb{E}[f(W') - f(W) \mid W] = \mathcal{G}f(W),$$

then  $\mathbb{E}[\mathcal{G}f(W)] = 0$ . If  $\mathcal{G}$  is tractable, we can seek  $\mathcal{A} \approx \mathcal{G}$ . We often rescale and seek  $\mathcal{A} \approx \kappa\mathcal{G}$ , where  $\kappa$  is a constant factor (if we consider a sequence of random variables  $W_n$ , the order of the underlying operators  $\mathcal{G}_n$  often depends on  $n$ , whereas we wish the operators  $\mathcal{A}_n$  to converge to a certain operator). If  $\mathcal{G}$  is not tractable, we can try

$$\mathbb{E}[\kappa(f(W') - f(W)) - R_f \mid W] = \mathcal{A}f(W),$$

where  $R_f$  is intended to be of smaller order than  $\kappa(f(W') - f(W))$ , as we seek  $\mathcal{A}$  close to  $\kappa\mathcal{G}$ .

It is also possible to consider several perturbations  $W'_\alpha$ . If they are indexed over a countable set, consider factors  $\kappa_\alpha$  and let

$$\mathbb{E}\left[\sum_{\alpha} \kappa_{\alpha}(f(W') - f(W)) - R_f \mid W\right] = \mathcal{A}f(W). \quad (1.2.1)$$

The perturbations can even be indexed over a measurable space. In this case, take  $\kappa$  to be a measure on this space and let

$$\mathbb{E}\left[\int (f(W'_\alpha) - f(W)) \kappa(d\alpha) - R_f \mid W\right] = \mathcal{A}f(W). \quad (1.2.2)$$

The approximating distribution  $\nu$  should be an annihilator of the image of  $\mathcal{A}$ , i. e.,  $\langle \mathcal{A}f, \nu \rangle = 0$  for all  $f \in \mathcal{F}$  (by this method,  $\mathcal{F}$  must be a space of functions). This method can also help solve the Stein equation. Details will be given later: see the beginning of Chapter 2.

Below we give three examples of Stein operators obtained by the first method. Applications of the second method will be given later.

1. *Poisson approximation.* Let  $\lambda > 0$  and let  $\nu = \text{Po}(\lambda)$  be the Poisson distribution. This is a probability measure on  $\mathbb{N}_0 := \{0, 1, 2, \dots\}$ . Signed measures on  $\mathbb{N}_0$  can be identified with their mass functions, so that we simply write

$$\nu(k) = \frac{\lambda^k e^{-\lambda}}{k!}.$$

Accordingly, for each signed measure  $\rho$  on  $\mathbb{N}_0$  and each function  $h \in L^1(|\rho|)$ , we have

$$\langle h, \rho \rangle = \sum_{k=0}^{\infty} h(k) \rho(k).$$

Now observe that

$$\frac{\nu(k-1)}{\nu(k)} = \frac{k}{\lambda}$$

for all  $k \in \mathbb{N} := \{1, 2, 3, \dots\}$ . This can be formulated in terms of linear operators: let  $\mathcal{B}$  be the operator mapping from the space  $\mathcal{M}$  of signed measures  $\rho$  on  $\mathbb{N}_0$  with  $\sum_{k=0}^{\infty} k |\rho(k)| < \infty$  to the space  $\mathcal{Y}$  of all signed measures on  $\mathbb{N}$  defined by

$$\mathcal{B}\rho(k) := \lambda \rho(k-1) - k \rho(k).$$

Then  $\mathcal{B}\nu = 0$ .

Now let  $\mathcal{F}$  be the space of all bounded functions on  $\mathbb{N}$ . Observe that for each  $f \in \mathcal{F}$ ,

$$\begin{aligned} \langle f, \mathcal{B}\rho \rangle &= \sum_{k=1}^{\infty} \lambda f(k) \rho(k-1) - \sum_{k=1}^{\infty} k f(k) \rho(k) \\ &= \sum_{w=0}^{\infty} \lambda f(w+1) \rho(w) - \sum_{w=1}^{\infty} w f(w) \rho(w) \end{aligned}$$

(notice that all sums converge due to the definitions of  $\mathcal{F}$  and  $\mathcal{Y}$ ). Let  $\mathcal{H}$  be the space of all functions on  $\mathbb{N}_0$ , such that there exists  $M$ , such that  $|h(w)| \leq M w$  for all  $w \in \mathbb{N}_0$ . Then  $\mathcal{B}$  is dual to the operator  $\mathcal{A}: \mathcal{F} \rightarrow \mathcal{H}$  defined by

$$\mathcal{A}f(w) := \begin{cases} \lambda f(w+1) - w f(w) & ; w = 1, 2, 3, \dots \\ \lambda f(1) & ; w = 0. \end{cases} \quad (1.2.3)$$

This is the usual Stein operator for the Poisson distribution. To simplify computations, we shall usually assume that  $f$  is also defined at 0 in an arbitrary way.

2. *Binomial approximation.* Let  $n \in \mathbb{N}$ , let  $0 < \theta < 1$  and let  $\nu = \text{Bin}(n, \theta)$  be the binomial distribution. This is a probability measure on  $\{0, 1, 2, \dots, n\}$  given by

$$\nu(k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}.$$

Similarly as for the Poisson distribution, observe that

$$\frac{\nu(k-1)}{\nu(k)} = \frac{k}{n-k+1} \frac{1-\theta}{\theta}$$

for all  $k \in \{1, 2, \dots, n\}$ . Thus, if  $\mathcal{B}$  be the operator mapping from the space  $\mathcal{M}$  of signed measures (or, equivalently, functions) on  $\{0, 1, \dots, n\}$  to the space  $\mathcal{Y}$  of all signed measures on  $\{1, 2, \dots, n\}$  defined by

$$\mathcal{B}\rho(k) := (n - k + 1)\theta \rho(k-1) - k(1 - \theta) \rho(k),$$

we have  $\mathcal{B}\nu = 0$ .

Letting  $\mathcal{F}$  be the space of all functions on  $\{1, 2, \dots, n\}$ , observe that for each  $f \in \mathcal{F}$ ,

$$\begin{aligned} \langle f, \mathcal{B}\rho \rangle &= \sum_{k=1}^n (n - k + 1)\theta f(k) \rho(k-1) - \sum_{k=1}^n k(1 - \theta) f(k) \rho(k) \\ &= \sum_{w=0}^{n-1} (n - w)\theta f(w+1) \rho(w) - \sum_{w=1}^n w(1 - \theta) f(w) \rho(w). \end{aligned}$$



Letting  $\mathcal{H}$  be the space of all functions on  $\{0, 1, 2, \dots, n\}$ ,  $\mathcal{B}$  is dual to the operator  $\mathcal{A}: \mathcal{F} \rightarrow \mathcal{H}$  defined by

$$\mathcal{A}f(w) := \begin{cases} (n-w)\theta f(w+1) - w(1-\theta)f(w) & ; w = 1, 2, 3, \dots, n-1 \\ (n-w)\theta f(1) & ; w = 0 \\ -n(1-\theta)f(w) & ; w = n. \end{cases}$$

This is one of possible choices for the Stein operator for the binomial distribution. Again, to simplify computations, we can assume that  $f$  is also defined at 0 and  $n+1$  in an arbitrary way.

3. *Normal and other continuous approximations.* Let  $\nu = \mathcal{N}(0, 1)$  be the standard normal distribution (the generalization to arbitrary mean and variance is straightforward). This is a continuous distribution, and we shall identify continuous distributions with their densities (more precisely, with equivalence classes of functions, where two functions are equivalent if they differ on a set with measure zero). Thus, we have

$$\nu(x) = e^{-x^2/2},$$

which is differentiable with

$$\nu'(x) = -x e^{-x^2/2}.$$

Thus, if, let's say,  $\mathcal{M}$  is the space of all continuously differentiable functions on the real line with exponentially decaying derivative,  $\mathcal{Y}$  is the space of all continuous exponentially decaying functions, and  $\mathcal{B}: \mathcal{M} \rightarrow \mathcal{Y}$  is defined by

$$\mathcal{B}\rho(x) := \rho'(x) + x\rho(x),$$

then  $\mathcal{B}\nu = 0$ . Let  $\mathcal{F}$  be the space of all continuously differentiable functions on the real line of polynomial growth. For each  $f \in \mathcal{F}$  and  $\rho \in \mathcal{M}$ , integration by parts gives

$$\begin{aligned} \langle f, \mathcal{B}\rho \rangle &= \int_{-\infty}^{\infty} f(x) \rho'(x) dx + \int_{-\infty}^{\infty} x f(x) \rho(x) dx \\ &= - \int_{-\infty}^{\infty} f'(x) \rho(x) dx + \int_{-\infty}^{\infty} x f(x) \rho(x) dx \end{aligned}$$

(polynomial growth of  $f$  and exponential decay of  $\rho$  ensure that  $\lim_{x \rightarrow \pm\infty} f(x) \rho(x) = 0$ ). Letting  $\mathcal{H}$  be the space of all functions on the real line of polynomial growth, we find that  $\mathcal{B}$  is dual to the operator  $\mathcal{A}: \mathcal{F} \rightarrow \mathcal{H}$  defined by

$$\mathcal{A}f(w) := -f'(x) + x f(x).$$

However, usually, we take  $\mathcal{A}$  with the opposite sign, i. e.,  $\mathcal{A}f(w) = f'(x) - x f(x)$ .

**Remark 1.2.1.** In view of Remark 1.1.3, the preceding derivation of the Stein operator for the normal distribution does not imply that normal approximation is restricted to measures with continuously differentiable densities. Indeed, the space  $\mathcal{M}$  can be extended to quite a general class of measures, but the extension can depend on the class  $\mathcal{H}$  of test functions.

**Remark 1.2.2.** Under suitable conditions, for a probability density  $\nu$  with continuously differentiable derivative and with  $\nu'(x) = \psi(x)\nu(x)$ , we obtain a Stein operator  $\mathcal{A}f(x) = f'(x) - \psi(x)f(x)$ . We omit the details.

# Chapter 2

## Poisson approximation

### 2.1 Independent trials

Let  $X_i \sim \text{Be}(p_i)$ ,  $i \in \mathcal{I}$ , be independent Bernoulli random variables, i. e.,  $\mathbb{P}(X_i = 1) = p_i$  and  $\mathbb{P}(X_i = 0) = 1 - p_i$  for all  $i \in \mathcal{I}$ . Assume that  $\sum_{i \in \mathcal{I}} p_i = \lambda < \infty$ . Notice that due to the latter condition, we can assume without loss of generality that  $\mathcal{I}$  is countable (but may be infinite). By a well known result, the sum  $W = \sum_{i \in \mathcal{I}} X_i$  converges almost surely.

It is known that if, roughly speaking, the  $p_i$ 's are small, then  $W$  approximately follows the Poisson distribution. However, we shall here pretend that we do not know this fact and try to find a Stein operator by the second method mentioned at the beginning of Section 1.2. Actually, at some point, we use the fact that if  $\mathcal{I} = \{1, 2, \dots, n\}$ ,  $p_i = \lambda/n$  for all  $i$  and  $n$  tends to the infinity, then  $W$  weakly converges to some distribution, but we will not need to know to which one.

Let the collection  $(X'_i)_{i \in \mathcal{I}}$  be an independent copy of the collection  $(X_i)_{i \in \mathcal{I}}$ . Denote  $W_i := W - X_i = \sum_{j \in \mathcal{I} \setminus \{i\}} X_j$ . Then  $W'_i := W_i + X'_i$  follows the same distribution as  $W$ . In fact,  $W$  and  $W'$  even form an exchangeable pair, but we shall not need this fact.

Define a relation  $U \sim V$  iff  $\mathbb{E}(U | W) = \mathbb{E}(V | W)$ . In addition, let  $\Delta f(w) := f(w + 1) - f(w)$  be the forward difference. Now take a function  $g: \mathbb{N}_0 \rightarrow \mathbb{R}$  and observe that

$$\begin{aligned} g(W'_i) - g(W) &= (g(W'_i) - g(W_i)) - (g(W) - g(W_i)) \\ &= \Delta g(W_i) X'_i - \Delta g(W - 1) X_i \\ &\sim p_i \Delta g(W_i) - \Delta g(W - 1) X_i \\ &= p_i \Delta g(W) - p_i \Delta^2 g(W - 1) X_i - \Delta g(W - 1) X_i \end{aligned}$$

with the following two supplements:

- The calculation is valid provided that  $\Delta^2 g$  is bounded: if  $|\Delta^2 g(w)| \leq M$  for all  $w \in \mathbb{N}_0$ , then  $|\Delta g(W_i) - \Delta g(W)| \leq |\Delta^2 g(W)| \leq M$ . Thus,  $\mathbb{E}(|\Delta g(W_i)| | W) \leq |\Delta g(W)| + M$  and  $\mathbb{E}(|\Delta g(W_i)| X_i | W) \leq p_i (|\Delta g(W)| + M)$ .
- For the sake of simplicity, we here assume that  $g$  is also defined at  $-1$ ; the values of the expressions are independent of the choice of  $g(-1)$ .

Now, in view of (1.2.1), sum over  $i$ :

$$\begin{aligned} \sum_{i \in \mathcal{I}} (g(W'_i) - g(W)) &\sim \lambda \Delta g(W) - \sum_{i \in \mathcal{I}} p_i \Delta^2 g(W-1) X_i - \Delta g(W-1)W \\ &= \tilde{\mathcal{A}}g(W) - \sum_{i \in \mathcal{I}} p_i \Delta^2 g(W-1) X_i, \end{aligned}$$

where  $\tilde{\mathcal{A}}g(w) := \lambda \Delta g(w) - \Delta g(w-1)w$ , noting that if  $\Delta^2 g$  is bounded, the summation is valid as well (in particular, all expressions almost surely exist). As a result, we have

$$\mathbb{E}[\tilde{\mathcal{A}}g(W)] = \sum_{i \in \mathcal{I}} p_i \mathbb{E}[\Delta^2 g(W-1) X_i],$$

From Section 1.2, recall that it makes sense to choose  $\tilde{\mathcal{A}}$  so that  $\sum_{i \in \mathcal{I}} p_i \Delta^2 g(W-1) X_i$  is of smaller order than  $\tilde{\mathcal{A}}g(W)$ . In our case, if  $|\Delta^2 g(w)| \leq M$  for all  $w \in \mathbb{N}_0$  (but not necessarily  $w = -1$ ), we can estimate:

$$\sum_{i \in \mathcal{I}} p_i \mathbb{E}|\Delta^2 g(W-1) X_i| \leq M \sum_{i \in \mathcal{I}} p_i^2.$$

Now if  $\mathcal{I} = \{1, 2, \dots, n\}$ ,  $p_i = \lambda/n$  for all  $i$  and  $n$  tends to the infinity, then recall that  $W$  weakly converges to some distribution, or, say, random variable. Then it is plausible that  $\tilde{\mathcal{A}}g(W)$  also converges in distribution to some random variable, while  $\sum_{i \in \mathcal{I}} p_i^2 = \lambda^2/n$  tends to zero. Thus, correct terms have been stored into  $\tilde{\mathcal{A}}g(W)$ .

We should now seek a probability measure  $\nu$  which annihilates the image of  $\tilde{\mathcal{A}}$ . However, we can just observe that  $\tilde{\mathcal{A}}g = \mathcal{A}f$ , where  $f(w) = \Delta g(w-1)$  and where  $\mathcal{A}$  is defined as in (1.2.3), that is

$$\mathcal{A}f(w) = \lambda f(w+1) - w f(w)$$

(again, we assume that  $f(0)$  is defined, but for  $w \in \mathbb{N}_0$ , the value of  $\mathcal{A}f(w)$  is independent of  $f(0)$ ). In Section 1.2, we have shown that the Poisson distribution  $\text{Po}(\lambda)$ , that is,

$$\text{Po}(\lambda)(k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

annihilates  $\mathcal{A}f$  (that is,  $\langle \mathcal{A}f, \text{Po}(\lambda) \rangle = 0$ ), provided that  $f$  is bounded. However, it is straightforward to check that this remains true if  $\Delta f$  is bounded. Therefore,  $\langle \tilde{\mathcal{A}}g, \text{Po}(\lambda) \rangle = 0$  if  $\Delta^2 g$  is bounded.

For each  $f: \mathbb{N}_0 \rightarrow \mathbb{R}$ , such that  $|\Delta f(w)| \leq M$  for all  $w \in \mathbb{N}$  (but not necessarily  $w = 0$ ), there exists  $g: \{-1, 0, 1, \dots\} \rightarrow \mathbb{R}$ , such that  $f(w) = \Delta g(w-1)$  for all  $w \in \mathbb{N}_0$ . Therefore, for each  $f$  satisfying the above-mentioned condition, we have

$$\mathbb{E}[\mathcal{A}f(W)] = \sum_{i \in \mathcal{I}} p_i \Delta f(W) X_i$$

(and it is also straightforward to verify this identity directly). If  $f$  is a solution to the Stein equation

$$\mathcal{A}f(w) = h(w) - \langle h, \text{Po}(\lambda) \rangle, \quad (2.1.1)$$

this gives a bound on the error in the Poisson approximation, as summarized in the assertion below.

**Proposition 2.1.1.** *If  $f$  is a solution to (2.1.1) with  $|\Delta f(w)| \leq M$  for all  $w \in \mathbb{N}$  (but not necessarily  $w = 0$ ), then we have*

$$|\mathbb{E}[h(W)] - \langle h, \text{Po}(\lambda) \rangle| \leq M \sum_{i \in \mathcal{I}} p_i^2. \quad (2.1.2)$$

□

## 2.2 Solution to the Stein equation

In the previous section, we have proved that we can efficiently estimate the error in the Poisson approximation of the sum of independent Bernoulli indicators with small success probabilities provided that we can bound  $\Delta f$  for some solution  $f$  to the Stein equation (2.1.1). This is what we shall do in this section. However, the bounds obtained here will not only be useful for sums of independent indicators, but also for indicators with certain dependence structure, which will be described in the sequel.

We shall here bound  $\Delta f$  for the case where  $h$  is an indicator of a set  $A \subseteq \mathbb{N}_0$ . This will allow us to bound the difference  $\mathbb{P}(W \in A) - \text{Po}(\lambda)(A)$ .

First, we shall find a solution  $f_j$  for the indicator of the singleton  $\{j\}$ , i. e., the solution to the equation

$$\lambda f_j(w+1) - w f_j(w) = \mathbf{1}(w=j) - \frac{\lambda^j}{j!} e^{-\lambda}. \quad (2.2.1)$$

Seeking  $f_j$  to be in the form

$$f_j(w) = \frac{(w-1)!}{\lambda^w} \psi_j(w) \quad \text{for all } w \in \mathbb{N},$$

equation (2.2.1) reduces to

$$\psi_j(w+1) - \psi_j(w) = \frac{\lambda^w}{w!} \left( \mathbf{1}(w=j) - \frac{\lambda^j}{j!} e^{-\lambda} \right) \quad \text{for all } w \in \mathbb{N}_0, \quad (2.2.2)$$

setting  $\psi_j(0) := 0$ . For  $w = 1, 2, \dots, j$ , we obtain

$$\psi_j(w) = \sum_{k=0}^{w-1} (\psi_j(k+1) - \psi_j(k)) = -\frac{\lambda^j}{j!} e^{-\lambda} \sum_{k=0}^{w-1} \frac{\lambda^k}{k!}. \quad (2.2.3)$$

Next, if we wish  $\Delta f_j$  to be bounded, we must have  $\lim_{w \rightarrow \infty} \psi_j(w) = 0$ . But then we have

$$\psi_j(w) = -\sum_{k=w}^{\infty} (\psi_j(k+1) - \psi_j(k)) = \frac{\lambda^j}{j!} e^{-\lambda} \sum_{k=w}^{\infty} \frac{\lambda^k}{k!} \quad (2.2.4)$$

for all  $w = j + 1, j + 2, \dots$ . The function  $\psi_j$  is now completely determined and satisfies equation (2.2.2) for all  $w \notin j$ . For  $w = j$ , we obtain

$$\begin{aligned} \psi_j(j+1) - \psi_j(j) &= \frac{\lambda^j}{j!} e^{-\lambda} \left( \sum_{k=0}^{j-1} \frac{\lambda^k}{k!} + \sum_{k=j+1}^{\infty} \frac{\lambda^k}{k!} \right) \\ &= \frac{\lambda^j}{j!} e^{-\lambda} \left( \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} - \frac{\lambda^j}{j!} \right) \\ &= \frac{\lambda^j}{j!} e^{-\lambda} \left( e^\lambda - \frac{\lambda^j}{j!} \right) \\ &= \frac{\lambda^j}{j!} \left( 1 - \frac{\lambda^j}{j!} e^{-\lambda} \right) \end{aligned}$$

and (2.2.2) holds true in this case, too.

Now we examine  $\Delta f_j(w)$ . We distinguish three cases.

*Case 1:*  $w = 1, 2, \dots, j - 1$ . Write

$$\begin{aligned} f_j(w+1) &= -\frac{\lambda^j}{j!} e^{-\lambda} \frac{w!}{\lambda^{w+1}} \sum_{l=0}^w \frac{\lambda^l}{l!}, \\ f_j(w) &= -\frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^w} \sum_{k=0}^{w-1} \frac{\lambda^k}{k!} \\ &= -\frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^w} \sum_{l=1}^w \frac{\lambda^{l-1}}{(l-1)!} \\ &= -\frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^{w+1}} \sum_{l=0}^w \frac{l \lambda^l}{l!}, \\ \Delta f_j(w) &= -\frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^{w+1}} \sum_{l=0}^w \frac{(w-l)\lambda^l}{l!} \\ &\leq 0. \end{aligned}$$

Case 2:  $w = j + 1, j + 2, \dots$ . Write

$$\begin{aligned} f_j(w+1) &= \frac{\lambda^j}{j!} e^{-\lambda} \frac{w!}{\lambda^{w+1}} \sum_{l=w+1}^{\infty} \frac{\lambda^l}{l!}, \\ f_j(w) &= \frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^w} \sum_{k=w}^{\infty} \frac{\lambda^k}{k!} \\ &= \frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^w} \sum_{l=w+1}^{\infty} \frac{\lambda^{l-1}}{(l-1)!} \\ &= \frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^{w+1}} \sum_{l=w+1}^{\infty} \frac{l \lambda^l}{l!}, \\ \Delta f_j(w) &= \frac{\lambda^j}{j!} e^{-\lambda} \frac{(w-1)!}{\lambda^{w+1}} \sum_{l=w+1}^{\infty} \frac{(w-l) \lambda^l}{l!} \\ &\leq 0. \end{aligned}$$

Case 3:  $w = j$ . Write

$$\begin{aligned} f_j(j) &= -\frac{\lambda^j}{j!} e^{-\lambda} \frac{(j-1)!}{\lambda^j} \sum_{k=0}^{j-1} \frac{\lambda^k}{k!} = -\frac{e^{-\lambda}}{j} \sum_{l=1}^j \frac{\lambda^{l-1}}{(l-1)!}, \\ f_j(j+1) &= \frac{\lambda^j}{j!} e^{-\lambda} \frac{j!}{\lambda^{j+1}} \sum_{k=j+1}^{\infty} \frac{\lambda^k}{k!} = \frac{e^{-\lambda}}{\lambda} \sum_{l=j+1}^{\infty} \frac{\lambda^l}{l!}, \\ \Delta f_j(j) &= \frac{e^{-\lambda}}{\lambda} \left( \frac{1}{j} \sum_{l=1}^j \frac{\lambda^l}{(l-1)!} + \sum_{l=j+1}^{\infty} \frac{\lambda^l}{l!} \right) \leq \frac{e^{-\lambda}}{\lambda} \sum_{l=1}^{\infty} \frac{\lambda^l}{l!} = \frac{1 - e^{-\lambda}}{\lambda}. \end{aligned}$$

Therefore,  $\Delta f_j(w) \leq \frac{1 - e^{-\lambda}}{\lambda}$  for all  $w \in \mathbb{N}$ .

For convenience, set  $f_j(0) := 0$  for all  $j \in \mathbb{N}_0$ . Now take  $A \subseteq \mathbb{N}_0$  and define  $f_A(w) := \sum_{j \in A} f_j(w)$ . We will show that the series converges pointwise and that  $f_A$  solves the Stein equation

$$\lambda f_A(w+1) - w f_A(w) = \mathbf{1}(w \in A) - \text{Po}(\lambda)(A). \quad (2.2.5)$$

Clearly,  $\mathbf{1}(w \in A) - \text{Po}(\lambda)(A) = \sum_{j \in A} (\mathbf{1}(w = j) - \text{Po}(\lambda)(\{j\}))$ . As to the left hand side, first deduce from (2.2.3) and (2.2.4) that  $|\psi_j(w)| \leq \frac{\lambda^j}{j!} e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = \frac{\lambda^j}{j!}$  for all  $w \in \mathbb{N}_0$ . Therefore,

$$\sum_{j=0}^{\infty} |f_j(w)| \leq \frac{(w-1)!}{\lambda^w} \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} = \frac{(w-1)!}{\lambda^w} e^{\lambda} < \infty.$$

As a result, the series  $\sum_{j \in A} f_j(w)$  converges pointwise and, moreover,  $\lambda f_A(w+1) - w f_A(w) = \sum_{j \in A} (\lambda f_j(w+1) - w f_j(w))$ . Combining all together, the proof of (2.2.5) is complete.

We have now shown that  $\Delta f_A(w) \leq \frac{1 - e^{-\lambda}}{\lambda}$  for all  $w \in \mathbb{N}$ . However,  $f_{\mathbb{N}_0}(w) = f_A + f_{A^c}$  solves  $\lambda f_{\mathbb{N}_0}(w+1) - w f_{\mathbb{N}_0}(w) = 0$ , which implies  $f_{\mathbb{N}_0}(w) = 0$  for all  $w \in \mathbb{N}$ . Therefore,  $\Delta f_A(w) = -\Delta f_{A^c}(w) \geq \frac{1 - e^{-\lambda}}{\lambda}$ . We have now proved the following assertion:

**Proposition 2.2.1.** *For each  $A \subseteq \mathbb{N}_0$ , there exists a solution  $f_A$  to the equation (2.2.5). This solution is uniquely determined on  $\mathbb{N}$  and can be arbitrary at 0. Finally, for all  $w \in \mathbb{N}$ , we have*

$$|\Delta f_A(w)| \leq \frac{1 - e^{-\lambda}}{\lambda}.$$

Combining with Proposition 2.1.1, we obtain the following neat result:

**Proposition 2.2.2.** *For independent random variables  $X_i \sim \text{Be}(p_i)$ ,  $i \in \mathcal{I}$ , with  $\lambda = \sum_{i \in \mathcal{I}} p_i < \infty$  and  $W = \sum_{i \in \mathcal{I}} X_i$ , we have*

$$|\mathbb{P}(W \in A) - \text{Po}(\lambda)(A)| \leq \frac{1 - e^{-\lambda}}{\lambda} \sum_{i \in \mathcal{I}} p_i^2 \quad (2.2.6)$$

for all  $A \subseteq \mathbb{N}_0$ .

**Remark 2.2.3.** We have  $\frac{1 - e^{-\lambda}}{\lambda} = \int_0^1 e^{-\lambda t} dt \leq 1$  and of course  $\frac{1 - e^{-\lambda}}{\lambda} \leq \frac{1}{\lambda}$ . Estimate (2.2.6) is essentially due to Le Cam [23], who proved the estimates

$$|\mathbb{P}(W \in A) - \text{Po}(\lambda)(A)| \leq \sum_{i \in \mathcal{I}} p_i^2 \quad \text{and} \quad |\mathbb{P}(W \in A) - \text{Po}(\lambda)(A)| \leq \frac{8}{\lambda} \sum_{i \in \mathcal{I}} p_i^2$$

(the latter provided that  $\max_i p_i \leq 1/4$ ), using an entirely different argument. However, as we shall see in the sequel, Stein's method can go far beyond independence. This is in contrast to the convolution argument used by Le Cam.

## 2.3 Abstract results

### 2.3.1 Coupling

Coupling of two probability distributions means constructing a joint probability distribution with marginals being the two given distributions. In other words, suppose that we have two random variables  $X$  and  $Y$ , which may be defined on different probability spaces. Coupling means constructing random variables  $X'$  and  $Y'$  defined on the *same* probability space, such that  $X'$  follows the same distribution as  $X$  and  $Y'$  follows the same distribution as  $Y$ .

In applications of Stein's method, coupling is extremely useful. In particular, if  $W$  is the random variable whose distribution is to be approximated and  $X$  is another random variable related to  $W$ , it is often useful to couple the (unconditional) distribution of  $W$  with conditional distributions of  $W$  given  $X = x$ . Without loss of generality, one can assume that there are random variables  $W_x$  defined on the same probability space as  $W$ , such that for each  $x$ , the distribution of  $W_x$  agrees with the conditional distribution of  $W$  given  $X = x$ .

**Example 2.3.1.** In Section 1.2, we mentioned that a Stein operator can be constructed by perturbing the original random variable  $W$  to a random variable  $W'$  with the same



distribution. Now let  $W$ ,  $X$  and  $W_x$  be as above and defined on the same probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . Suppose that the map  $(\omega, \omega') \mapsto W_\omega(\omega')$  is measurable with respect to  $\mathcal{F} \otimes \mathcal{F}$ . If there is another random variable  $X'$  which follows the same distribution of  $X$  and is independent of all other random variables, then the random variable  $W'(\omega) := W_{X(\omega)}(\omega)$  follows the same distribution as  $W$ .

Now let  $X_i \sim \text{Be}(p_i)$ ,  $i \in \mathcal{I}$ , be again a Bernoulli random variables with  $\lambda = \sum_{i \in \mathcal{I}} p_i < \infty$  and  $W = \sum_{i \in \mathcal{I}} X_i$ , but now we drop the assumption that they are independent. Instead, suppose that for each  $i \in \mathcal{I}$ , there exists a random variable  $\tilde{W}_i$ , such that the (unconditional) distribution of  $\tilde{W}_i + 1$  agrees with the conditional distribution of  $W$  given  $X_i = 1$ , or, equivalently, the (unconditional) distribution of  $\tilde{W}_i$  agrees with the conditional distribution of  $\sum_{j \in \mathcal{I} \setminus \{i\}} X_j$  given  $X_i = 1$ . Notice that if the random variables  $X_i$  are independent, we can simply take  $\tilde{W}_i = W - X_i$ . Now observe that

$$\begin{aligned} \mathbb{E}[\lambda f(W+1) - f(W)W] &= \sum_{i \in \mathcal{I}} \mathbb{E}[p_i f(W+1) - f(W)X_i] \\ &= \sum_{i \in \mathcal{I}} p_i \mathbb{E}[f(W+1) - f(\tilde{W}_i+1)]. \end{aligned}$$

If  $|\Delta f(w)| \leq M$  for all  $w \in \mathbb{N}$ , then

$$|\mathbb{E}[\lambda f(W+1) - f(W)W]| \leq M \sum_{i \in \mathcal{I}} p_i \mathbb{E}|W - \tilde{W}_i|.$$

Combining with Proposition 2.2.1, we find that

$$|\mathbb{P}(W \in A) - \text{Po}(\lambda)(A)| \leq \frac{1 - e^{-\lambda}}{\lambda} \sum_{i \in \mathcal{I}} p_i \mathbb{E}|W - \tilde{W}_i|$$

for all  $A \subseteq \mathbb{N}_0$ . This can be expressed in terms of the *total variation distance* between two probability measures on  $\mathbb{N}_0$  (see Example A.1.2):

$$d_{\text{TV}}(\mu, \nu) := \sup_{A \subseteq \mathbb{N}_0} |\mu(A) - \nu(A)|.$$

Denoting by  $\mathcal{L}(W)$  the distribution of  $W$ , that is,  $\mathcal{L}(W)(A) := \mathbb{P}(W \in A)$ , this leads to the following result:

**Theorem 2.3.2.** *For  $X_i$ ,  $W$  and  $\tilde{W}_i$  being as above, we have*

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \sum_{i \in \mathcal{I}} p_i \mathbb{E}|W - \tilde{W}_i|. \quad (2.3.1)$$

□

**Remark 2.3.3.** We have only needed to consider conditional distributions of  $W$  given  $X_i = 1$ , not given  $X_i = 0$ .

**Remark 2.3.4.** For independent indicators, the preceding theorem reduces to Proposition 2.2.2.

### 2.3.2 Local dependence and decompositions

Coupling is not the only approach to dependence where Stein's method can be applied. An important alternative concept is *local dependence*. In the context of Stein's method, the latter was introduced by Chen [9] and refined by Arratia, Goldstein and Gordon [1].

Again, take a sum  $W = \sum_{i \in \mathcal{I}} X_i$  of Bernoulli random variables  $X_i \sim \text{Be}(p_i)$  and assume that for each  $i \in \mathcal{I}$ , there exists a so called *dependence neighbourhood*  $\mathcal{N}_i$ , such that  $X_i$  is independent of the subfamily  $X_j, j \in \mathcal{I} \setminus (\{i\} \cup \mathcal{N}_i)$ . Under this assumption,  $\mathcal{N}_i$  can actually be a *punctured neighbourhood* of  $i$ .

More generally, suppose that for each  $i$ , there is a decomposition

$$W = X_i + Y_i + Z_i, \quad (2.3.2)$$

where  $Z_i$  is independent of  $X_i$ . In this case, write

$$\mathbb{E}[\lambda f(W+1) - f(W)W] = \sum_{i \in \mathcal{I}} \mathbb{E}[p_i f(X_i + Y_i + Z_i + 1) - f(Y_i + Z_i + 1)X_i].$$

Since  $X_i$  is independent of  $Z_i$ , we have

$$\mathbb{E}[p_i f(Z_i + 1) - f(Z_i + 1)X_i] = 0$$

and therefore

$$\begin{aligned} \mathbb{E}[\lambda f(W+1) - f(W)W] &= \sum_{i \in \mathcal{I}} p_i \mathbb{E}[f(X_i + Y_i + Z_i + 1) - f(Z_i + 1)] \\ &\quad - \sum_{i \in \mathcal{I}} \mathbb{E}[(f(Y_i + Z_i + 1) - f(Z_i + 1))X_i]. \end{aligned}$$

Now if  $|\Delta f(w)| \leq M$  for all  $w \in \mathbb{N}$ , we can estimate

$$\left| \mathbb{E}[\lambda f(W+1) - f(W)W] \right| \leq M \sum_{i \in \mathcal{I}} \left( p_i \mathbb{E}|X_i + Y_i| + \mathbb{E}[X_i |Y_i|] \right).$$

Combining with Proposition 2.2.2, we obtain the following result:

**Theorem 2.3.5.** *For  $W$  decomposed as in (2.3.2) and  $\tilde{W}_i$  being as above, we have*

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \sum_{i \in \mathcal{I}} \left( p_i \mathbb{E}|X_i + Y_i| + \mathbb{E}[X_i |Y_i|] \right). \quad (2.3.3)$$

□

**Remark 2.3.6.** For independent indicators, the preceding theorem again reduces to Proposition 2.2.2.

## 2.4 Applications

### 2.4.1 Pattern matching

Let  $\xi_i$ ,  $i \in \mathbb{N}_0$ , be a sequence of independent random variables taking values in a finite set  $L$ . The elements of  $L$  will be called *letters*. Fix a *pattern*  $l_0 l_1 \cdots l_{r-1}$ . We are interested in the number of occurrences of the given pattern in a certain subsequence of the random variables  $\xi_i$ . More formally, for  $i \in \mathcal{I} := \{0, 1, \dots, n-1\}$ , let

$$X_i := \mathbf{1}(\xi_i = l_0, \xi_{i+1} = l_1, \dots, \xi_{i+r-1} = l_{r-1}).$$

Then  $W := \sum_{i=0}^{n-1} X_i$  is the desired number of occurrences.

This problem is an obvious case of local dependence, as we can take dependence neighbourhoods

$$\mathcal{N}_i := \{i-r+1, i-r+2, \dots, i-1, i+1, i+2, \dots, i+r-1\} \cap \mathcal{I}$$

or, equivalently, take  $Y_i := X_{i-r+1} + X_{i-r+2} + \cdots + X_{i-1} + X_{i+1} + X_{i+2} + \cdots + X_{i+r-1}$ , letting  $X_j = 0$  for  $j \notin \mathcal{I}$ .

In the sequel, we shall only consider the case where  $L = \{0, 1\}$  and patterns of type  $11 \cdots 1$  (so called *r-runs*). In addition, we assume that  $\mathbb{P}(\xi_i = 1) = p$  for all  $i$ . In this case, we have  $\mathbb{E} X_i = p^r$  for all  $i$ , so that  $\lambda = \mathbb{E} W = np^r$ . Clearly, we have  $\mathbb{E} |X_i + Y_i| \leq (2r-1)p^r$ . Noting that

$$\mathbb{E}[X_i X_{i+k}] = \begin{cases} p^{r+k} & ; i, i+k \in \mathcal{I} \\ 0 & ; \text{otherwise} \end{cases}$$

for  $k = 0, 1, \dots, r-1$ , observe that  $\mathbb{E}[X_i | Y_i] \leq 2 \sum_{k=1}^{r-1} p^{r+k} = \frac{2(p^{r+1} - p^{2r})}{1-p}$ . Summing up, we find that

$$p^r \mathbb{E} |X_i + Y_i| + \mathbb{E}[X_i | Y_i] \leq (2r-3)p^{2r} + \frac{2p^{r+1}}{1-p}$$

and Theorem 2.3.5 yields

$$\begin{aligned} d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) &\leq \frac{1 - e^{-\lambda}}{\lambda} n \left( (2r-3)p^{2r} + \frac{2p^{r+1}}{1-p} \right) \\ &\leq \left( (2r-3)p^r + \frac{2p}{1-p} \right) \min\{1, np^r\}. \end{aligned} \tag{2.4.1}$$

However, it is also easy to construct explicit couplings: observe that a sequence of random letters following the conditional distribution given  $X_i = 1$  can be obtained from a sequence following the unconditional distribution simply by setting  $\xi_i, \xi_{i+1}, \dots, \xi_{i+r-1}$  to 1. Denote the new sequence by  $\tilde{\xi}_{i,0}, \tilde{\xi}_{i,1}, \dots$ . We may write  $\tilde{W}_i = \sum_{j \in \mathcal{I} \setminus \{i\}} \tilde{X}_{ij}$ , where  $\tilde{X}_{ij} = \mathbf{1}(\tilde{\xi}_{i,j} = \tilde{\xi}_{i,j+1} = \cdots = \tilde{\xi}_{i,j+r-1} = 1)$  (and  $\tilde{\xi}_{ij} = 0$  for  $j \notin \{0, 1, \dots, n-1\}$ ). Clearly,  $\tilde{X}_{ij} = X_j$  for  $|i-j| \geq r$ . Otherwise, observe that  $\tilde{X}_{ij}$  differs from  $X_j$  only if all appropriate  $i-j$  original letters equal 1. Therefore, we have  $\mathbb{E} |X_j - \tilde{X}_{ij}| \leq p^{|i-j|}$  for all  $j \neq i$ , so that

$$\mathbb{E} |W - \tilde{W}_i| \leq \sum_{j: 1 \leq |i-j| \leq r-1} p^{|i-j|} + p^r \leq 2 \sum_{k=1}^{\infty} p^k \leq \frac{2p}{1-p}$$

and Theorem 2.3.2 yields

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \frac{2np^{r+1}}{1-p} \leq \frac{2p}{1-p} \min\{1, np^r\},$$

which is better than (2.4.1).

## 2.4.2 Birthday problems

Suppose that there are  $n$  people and  $d$  days in the year. Each person has his/her birthday on a specific day with probability  $1/d$ , independently of other persons. Let  $W$  be the number of unordered pairs (i. e., sets with exactly two elements)  $\{i, j\}$ , such that the  $i$ -th and the  $j$ -th person have the same birthday.

First, we apply the concept of local dependence. Denoting by  $\mathcal{S}$  the set of all unordered pairs on  $\{1, 2, \dots, n\}$ , we have  $W = \sum_{\{i,j\} \in \mathcal{S}} X_{\{i,j\}}$ , where  $X_{\{i,j\}}$  is the indicator of the event that the  $i$ -th and the  $j$ -th person have the same birthday. Letting

$$Y_{\{i,j\}} := \sum_{\substack{1 \leq k \leq n \\ k \neq i,j}} (X_{\{i,k\}} + X_{\{k,j\}}),$$

observe that  $X_{\{i,j\}}$  is independent of  $Z_{\{i,j\}} := W - X_{\{i,j\}} - Y_{\{i,j\}}$ .

**Remark 2.4.1.** If  $i, j$  and  $k$  are distinct, then  $X_{\{i,j\}}$  is independent of  $X_{\{i,k\}}$  as well as of  $X_{\{k,j\}}$ . However, it is not *jointly* independent, i. e., it is not independent of the pair  $(X_{\{i,k\}}, X_{\{k,j\}})$ .

Noting that  $\mathbb{E} X_{\{i,j\}} = 1/d$  and  $\mathbb{E}[X_{\{i,j\}} X_{\{i,k\}}] = \mathbb{E}[X_{\{i,j\}} X_{\{k,j\}}] = 1/d^2$  for all distinct  $i, j, k$ , we have  $\lambda = \mathbb{E} W = \frac{n(n-1)}{2d}$  and Theorem 2.3.5 yields

$$\begin{aligned} d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) &\leq \frac{1 - e^{-\lambda}}{\lambda} \frac{n(n-1)(4n-7)}{2d^2} \\ &\leq \min \left\{ \frac{4n-7}{d}, \frac{n(n-1)(4n-7)}{2d^2} \right\}. \end{aligned} \quad (2.4.2)$$

Alternatively, one can also construct explicit couplings. Fixing the  $i$ -th and the  $j$ -th person, observe that a random pair of birthdays following the conditional distribution given that they both have the same birthday can be obtained from a pair following the unconditional distribution by picking one of them independently and uniformly at random and assigning the other the birthday of the chosen one. Since the birthdays are independent, we can leave the other persons unchanged. Accordingly, we construct  $\tilde{W}_{\{i,j\}}$ . We may write  $\tilde{W}_{\{i,j\}} = \sum_{\{k,l\} \in \mathcal{S} \setminus \{\{i,j\}\}} \tilde{X}_{\{i,j\},\{k,l\}}$ , where  $\tilde{X}_{\{i,j\},\{k,l\}}$  is the indicator of the event that  $k$  and  $l$  have the same birthday after the above-mentioned reassignment. Clearly,  $\tilde{X}_{\{i,j\},\{k,l\}} = X_{\{k,l\}}$  if  $\{i, j\} \cap \{k, l\} = \emptyset$ . Otherwise, we may assume without loss of generality that  $k = i$  (and  $l \neq i, j$ ). In this case,  $\tilde{X}_{\{i,j\},\{i,l\}}$  differs from  $X_{\{i,j\}}$  if, first,  $i$

is to be assigned the birthday of  $j$ , and second, if among the pairs  $\{i, l\}$  and  $\{j, l\}$ , there is exactly one of them such that both persons have the same birthday. Therefore,

$$\mathbb{E}|X_{\{i,j\}} - \tilde{X}_{\{i,j\},\{i,l\}}| = \frac{d-1}{d^2}$$

and

$$\mathbb{E}|W - \tilde{W}_{\{i,j\}}| \leq \frac{2(n-2)(d-1)}{d^2} + \frac{1}{d} = \frac{2nd - 3d - 2n + 4}{d^2}.$$

COmpared to (2.4.2), Theorem 2.3.2 yields a better bound

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \frac{n(n-1)(2nd - 3d - 2n + 4)}{2d^3}. \quad (2.4.3)$$

The exact probability  $\mathbb{P}(W = 0)$  is easy to compute – we have

$$\mathbb{P}(W = 0) = \frac{(d-1)(d-2) \cdots (d-n+1)}{d^{n-1}}$$

and for  $d = 365, n = 23$ , we have  $\mathbb{P}(W = 0) \doteq 0.4927028$  (closest to  $1/2$  for  $d = 365$ ). Poisson approximation yields  $\mathbb{P}(W = 0) \approx 0.4999982$ . From (2.4.3), we obtain an error bound 0.0587, which is quite larger than the actual error.

**Remark 2.4.2.** This example has many possible extensions. Among others, the birthday probabilities need not be equal and we may only consider pairs which form an edge in a certain given graph (i. e., we consider coloured graphs). Both extensions can be handled by Stein's method, see Barbour, Holst and Janson [5].

### 2.4.3 Random permutations

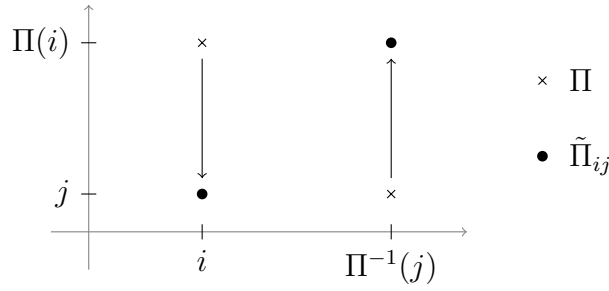
Let  $\Pi$  be a uniformly distributed random permutation of  $\{1, 2, \dots, n\}$ . Take  $C \subseteq \{1, 2, \dots, n\}^2$  and consider the statistic

$$W := \sum_{i=1}^n \mathbf{1}((i, \Pi(i)) \in C).$$

The indicators  $\mathbf{1}((i, \Pi(i)) \in C)$  do not exhibit any useful local dependence: any two of them are dependent. However, it is easy to construct a proper coupling along with a bound on the error in the Poisson approximation. One can rewrite  $W$  as

$$W = \sum_{(i,j) \in C} X_{ij}, \quad (2.4.4)$$

where  $X_{ij} := \mathbf{1}(\Pi(i) = j)$ . Given  $\Pi(i) = j$ ,  $\Pi$  is uniformly distributed over all permutations  $\tilde{\pi}$  with  $\tilde{\pi}(i) = j$ . Now define a new random permutation  $\tilde{\Pi}_{ij} := \tau_{\Pi(i),j} \circ \Pi$ , where  $\tau_{r,s}$  is the transposition exchanging  $r$  and  $s$  (and identity if  $r = s$ ):



We claim that  $\tilde{\Pi}_{ij}$  is uniformly distributed over all permutations  $\tilde{\pi}$  with  $\tilde{\pi}(i) = j$ . If  $\pi$  is a fixed permutation and  $\tilde{\pi} = \tau_{\pi(i),j} \circ \pi$ , then, clearly,  $\tilde{\pi}(i) = j$ . Now fix a permutation  $\tilde{\pi}$  with  $\tilde{\pi}(i) = j$  and examine all permutations  $\pi$  with  $\tilde{\pi} = \tau_{\pi(i),j} \circ \pi$ . As  $\pi = \tau_{j,\pi(i)} \circ \tilde{\pi}$ ,  $\pi$  must be of the form  $\pi = \tau_{j,l} \circ \tilde{\pi}$  for some  $l \in \{1, 2, \dots, n\}$ . Conversely, if  $\pi = \tau_{j,l} \circ \tilde{\pi}$ , then  $\pi(i) = l$  and  $\tau_{\pi(i),j} \circ \pi = \tau_{l,j} \circ \tau_{j,l} \circ \tilde{\pi} = \tilde{\pi}$ . Therefore, for each  $\tilde{\pi}$  with  $\tilde{\pi}(i) = j$ , there are exactly  $n$  permutations  $\pi$ , such that  $\tilde{\pi} = \tau_{\pi(i),j} \circ \pi$ . Therefore,  $\tilde{\Pi}_{ij}$  is indeed uniformly distributed over all permutations  $\tilde{\pi}$  with  $\tilde{\pi}(i) = j$ . Letting  $\tilde{X}_{ijkl} := \mathbf{1}(\tilde{\Pi}_{ij}(k) = l)$  and

$$\tilde{W}_{ij} = \sum_{(k,l) \in C} \tilde{X}_{ijkl} - 1,$$

the (unconditional) distribution  $\tilde{W} + 1$  agrees with the conditional distribution of  $W$  given  $X_{ij} = 1$ .

Now observe that  $\tilde{\Pi}_{ij}$  agrees with  $\Pi$  except on  $i$  and  $\Pi^{-1}(j)$ . More precisely, if  $\Pi(i) = j$ , then  $\tilde{X}_{ijkl} = X_{kl}$  for all  $(k, l) \in C$ . If  $\Pi(i) \neq j$ , then  $\tilde{X}_{ijkl}$  differs from  $X_{kl}$  in the following four disjoint cases:

- If  $k = i$  and  $l = j$ , then  $X_{kl} = 0$  and  $\tilde{X}_{ijkl} = 1$ .
- If  $k = i$  and  $l = \Pi(i)$ , then  $X_{kl} = 1$  and  $\tilde{X}_{ijkl} = 0$ .
- If  $\Pi(k) = j$  and  $l = j$ , then  $X_{kl} = 1$  and  $\tilde{X}_{ijkl} = 0$ .
- If  $\Pi(k) = j$  and  $l = \Pi(i)$ , then  $X_{kl} = 0$  and  $\tilde{X}_{ijkl} = 1$ .

Therefore, if  $(i, j) \in C$ , then

$$\begin{aligned} W - \tilde{W}_{ij} &= \sum_{(k,l) \in C} (X_{kl} - \tilde{X}_{ijkl}) + 1 \\ &= \sum_{l; (i,l) \in C} \mathbf{1}(\Pi(i) = l) + \sum_{k; (k,j) \in C} \mathbf{1}(\Pi(k) = j) - \sum_{(k,l) \in C} \mathbf{1}(\Pi(i) = l, \Pi(k) = j). \end{aligned}$$

Next, observe that the preceding identity remains true if  $\Pi(i) = j$ . Noting that the event  $\{\Pi(i) = l, \Pi(k) = j\}$  is only possible if either  $k = i$  and  $l = j$  or  $k \neq i$  and  $l \neq j$ , observe that

$$\begin{aligned} W - \tilde{W}_{ij} &= \mathbf{1}(\Pi(i) = j) + \sum_{\substack{l; (i,l) \in C \\ l \neq j}} \mathbf{1}(\Pi(i) = l) + \sum_{\substack{k; (k,j) \in C \\ k \neq i}} \mathbf{1}(\Pi(k) = j) \\ &\quad - \sum_{\substack{(k,l) \in C \\ k \neq i, l \neq j}} \mathbf{1}(\Pi(i) = l, \Pi(k) = j). \end{aligned}$$

Noting that  $\mathbb{P}(\Pi(i) = j) = \mathbb{P}(\Pi(k) = j) = \frac{1}{n}$  and  $\mathbb{P}(\Pi(i) = l, \Pi(k) = j) = \frac{1}{n(n-1)}$  for  $k \neq i, l \neq j$ , we estimate

$$\begin{aligned} \mathbb{E}|W - \tilde{W}_{ij}| &\leq \frac{1}{n} + \frac{1}{n} |\{l; (i, l) \in C, l \neq j\}| + \frac{1}{n} |\{k; (k, j) \in C, k \neq i\}| \\ &\quad + \frac{1}{n(n-1)} |\{(k, l) \in C; k \neq i, l \neq j\}| \\ &= \frac{n-2}{n(n-1)} + \frac{n-2}{n(n-1)} |\{l; (i, l) \in C, l \neq j\}| \\ &\quad + \frac{n-2}{n(n-1)} |\{k; (k, j) \in C, k \neq i\}| + \frac{1}{n(n-1)} |C| \\ &= -\frac{n-2}{n(n-1)} + \frac{n-2}{n(n-1)} |\{l; (i, l) \in C\}| \\ &\quad + \frac{n-2}{n(n-1)} |\{k; (k, j) \in C\}| + \frac{1}{n(n-1)} |C| \end{aligned}$$

For each  $(i, j) \in C$ , we have  $\mathbb{P}(X_{ij} = 1) = \frac{1}{n}$ . Summing up, we obtain

$$\begin{aligned} &\sum_{(i,j) \in C} \mathbb{P}(X_{ij} = 1) \mathbb{E}|W - \tilde{W}_{ij}| \\ &\leq -\frac{n-2}{n^2(n-1)} |C| + \frac{n-2}{n^2(n-1)} |\{(i, j, l); (i, j) \in C, (i, l) \in C\}| \\ &\quad + \frac{n-2}{n^2(n-1)} |\{(i, k, l) \in C; (i, j) \in C, (k, j) \in C\}| \\ &\quad + \frac{1}{n^2(n-1)} |C|^2. \end{aligned}$$

Introducing

$$p_i := \frac{|\{j; (i, j) \in C\}|}{n}, \quad q_j := \frac{|\{i; (i, j) \in C\}|}{n}, \quad \lambda := \frac{|C|}{n} = \mathbb{E}W \quad (2.4.5)$$

and applying Theorem 2.3.2, we obtain

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \left( \frac{n-2}{n-1} \sum_{i=1}^n p_i^2 + \frac{n-2}{n-1} \sum_{j=1}^n q_j^2 + \frac{\lambda^2}{n-1} - \frac{(n-2)\lambda}{n(n-1)} \right). \quad (2.4.6)$$

**Remark 2.4.3.** If a constant factor in the error bound does not matter, the bound in (2.4.6) can be simplified. First, by the inequality between the arithmetic and geometric mean, we can estimate

$$\lambda^2 = \sum_{i=1}^n p_i \sum_{j=1}^n q_j \leq \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (p_i^2 + q_j^2) = \frac{n}{2} \left( \sum_{i=1}^n p_i^2 + \sum_{j=1}^n q_j^2 \right),$$

leading to

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \left( \frac{3n-4}{2(n-1)} \sum_{i=1}^n p_i^2 + \frac{3n-4}{2(n-1)} \sum_{j=1}^n q_j^2 - \frac{(n-2)\lambda}{n(n-1)} \right).$$

In addition, since  $p_i \leq np_i^2$ , we can estimate

$$-\frac{(n-2)\lambda}{n(n-1)} = -\frac{\lambda}{n} + \frac{1}{n(n-1)} \sum_{i=1}^n p_i \leq -\frac{\lambda}{n} + \frac{1}{n-1} \sum_{i=1}^n p_i^2$$

combining all together, we obtain a simplified bound

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1-e^{-\lambda}}{\lambda} \left( \frac{3}{2} \sum_{i=1}^n p_i^2 + \frac{3}{2} \sum_{j=1}^n q_j^2 - \frac{\lambda}{n} \right).$$

Below we give two examples.

**Example 2.4.4.** The case of independent indicators can be regarded as a limiting case of random permutations. For the sake of simplicity, we assume that we have finitely many independent indicators  $X_1, X_2, \dots, X_r$  with success probabilities  $p_1, p_2, \dots, p_r$  all divisible by the same number  $s \in \mathbb{N}$ . For each  $m \in \mathbb{N}$ , define

$$C^{(m)} := \bigcup_{i=1}^r \{i\} \times \{1, 2, \dots, mp_i s\}.$$

For  $n := ms \geq r$ , let  $\Pi^{(m)}$  be a uniformly distributed random permutation of  $\{1, 2, \dots, n\}$ . Define

$$W^{(m)} := \sum_{i=1}^n X_i^{(m)} = \sum_{i=1}^r X_i^{(m)},$$

where  $X_i^{(m)} = \mathbf{1}((i, \Pi^{(m)}(i)) \in C)$ . Observe that  $X_i^{(m)} \sim \text{Be}(p_i)$ . It is plausible that in the limit as  $m$  tends to the infinity, the random variables  $X_1^{(m)}, \dots, X_r^{(m)}$  are independent. Next, observe that  $p_i$  and  $\lambda$  match the underlying quantities defined in (2.4.5). Letting  $q_j^{(m)}$  be the counterpart of  $q_j$  in (2.4.5), observe that  $q_j^{(m)} \leq r/n$ . Therefore, by (2.4.6), we have

$$d_{\text{TV}}(\mathcal{L}(W^{(m)}), \text{Po}(\lambda)) \leq \frac{1-e^{-\lambda}}{\lambda} \left( \frac{n-2}{n-1} \sum_{i=1}^n p_i^2 + \frac{n-2}{n-1} \frac{r^2}{n} + \frac{\lambda^2}{n-1} - \frac{(n-2)\lambda}{n(n-1)} \right)$$

and the right hand side tends to the bound in (2.2.6).

**Example 2.4.5. Matching problem.** Consider the setting of  $n = rd$  cards,  $d$  of which have face value  $i$ ,  $i = 1, 2, \dots, r$ , and draw  $r$  of them uniformly at random. A match occurs if a card with face  $i$  appears in the  $i$ -th drawing. Denoting by  $W$  the number of matchings, observe that this can be represented in form (2.4.4) if we take

$$C := \bigcup_{i=1}^r \{i\} \times \{d(i-1) + 1, d(i-1) + 2, \dots, d(i-1) + d\}.$$

Next, we have

$$p_i = \frac{1}{r} \mathbf{1}(i \leq r), \quad q_j = \frac{1}{n}, \quad \lambda = 1$$



and the bound (2.4.6) reduces after some calculation to

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq (1 - e^{-1}) \left( \frac{n-2}{n-1} \frac{1}{r} + \frac{1}{n-1} \right).$$

A matching problem was investigated in 1708 in de Montmort's book [14] on games of chance for the game 'Treize', where we have an ordinary deck of  $n = 52$  cards taking  $r = 13$  possible face values. The original problem was to find the probability that there was no match. The approximating probability equals  $\text{Po}(1)(\{0\}) = 1/e \doteq 0.3679$ . The exact probability can be computed by the inclusion–exclusion principle and equals 0.3569 (and was computed in 1711 by Nicolas Bernoulli, see page 324 of de Montmort [14]). The difference equals 0.0109, while the error bound equals 0.06 (all up to rounding errors).

### 2.4.4 Occupancy problems

Let  $r$  balls be thrown independently into a family of bins indexed by  $\mathcal{S}$ , with probability  $q_i$  of hitting the  $i$ -th bin. Suppose first that  $\mathcal{S}$  is finite and define  $W$  to be the number of empty bins. We can write  $W = \sum_{i \in \mathcal{S}} X_i$ , where  $X_i$  is the indicator of the event that the  $i$ -th bin is empty. Clearly,  $X_i \sim \text{Be}(p_i)$ , where

$$p_i = (1 - q_i)^r.$$

The event  $\{X_i = 1\}$  is the same as the event that the balls have been only thrown into bins in  $\mathcal{S} \setminus \{i\}$ . Conditioning on this event is the same as modifying the probability  $q_j$  to  $q_j/(1 - q_i)$  for  $j \neq i$  and  $q_i$  to 0. This is the same as first throwing the balls as initially and then relocating each ball that has landed in the  $i$ -th bin into the  $j$ -th bin,  $j \neq i$ , with probability  $q_j/(1 - q_i)$ . Taking  $W$  to be the number of empty bins before the rearrangement and  $\tilde{W}_i$  the number of empty bins except for the  $i$ -th bin after the rearrangement, we obtain the desired coupling.

We can write  $\tilde{W}_i = \sum_{j \in \mathcal{S} \setminus \{i\}} \tilde{X}_{ij}$ , where  $\tilde{X}_{ij}$  is the indicator of the event that the  $j$ -th bin is empty after the rearrangement. Write

$$W - \tilde{W}_i = X_i + \sum_{j \in \mathcal{S} \setminus \{i\}} (X_j - \tilde{X}_{ij}).$$

If the  $j$ -th bin is empty after the rearrangement, it must have been empty before the rearrangement, too. Therefore,  $X_j \geq \tilde{X}_{ij}$  and

$$\begin{aligned} \mathbb{E} |W - \tilde{W}_i| &= \mathbb{E} X_i + \sum_{j \in \mathcal{S} \setminus \{i\}} (\mathbb{E} X_j - \mathbb{E} \tilde{X}_{ij}) \\ &= (1 - q_i)^r + \sum_{j \in \mathcal{S} \setminus \{i\}} \left[ (1 - q_j)^r - \left( 1 - \frac{q_j}{1 - q_i} \right)^r \right]. \end{aligned}$$

We can further estimate this quantity by rewriting

$$\begin{aligned} (1 - q_j)^r - \left(1 - \frac{q_j}{1 - q_i}\right)^r &= (1 - q_j)^r - \left(1 - q_j - \frac{q_i q_j}{1 - q_i}\right)^r \\ &= (1 - q_j)^r \left[1 - \left(1 - \frac{q_i q_j}{(1 - q_i)(1 - q_j)}\right)^r\right]. \end{aligned}$$

Applying the inequality  $(1 - x)^r \geq 1 - rx$ , which holds true for all  $x \in [0, 1]$ , estimate

$$(1 - q_j)^r - \left(1 - \frac{q_j}{1 - q_i}\right)^r \leq r(1 - q_j)^r \frac{q_i q_j}{(1 - q_i)(1 - q_j)} = \frac{r p_j q_i q_j}{(1 - q_i)(1 - q_j)}.$$

Letting  $\lambda = \sum_{i \in \mathcal{I}} p_i$ , combining all together and applying Theorem 2.3.2, we obtain the following result:

**Proposition 2.4.6.** *The distribution of the number  $W$  of empty bins defined as above satisfies*

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \sum_{i \in \mathcal{I}} \left( p_i^2 + r \sum_{j \in \mathcal{I} \setminus \{i\}} \frac{p_i p_j q_i q_j}{(1 - q_i)(1 - q_j)} \right). \quad (2.4.7)$$

□

Now suppose that all bins are hit with equal probabilities, that is,  $q_i = \frac{1}{n}$ , where  $n$  denotes the number of bins. In this case, (2.4.7) reduces to

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq n \frac{1 - e^{-\lambda}}{\lambda} \left(1 - \frac{1}{n}\right)^{2r} \left(1 + \frac{r}{n-1}\right),$$

where  $\lambda = n\left(1 - \frac{1}{n}\right)^r$ . Letting  $a := r/n$ , observe that  $p_i \leq e^{-a}$  and consequently  $\lambda \leq n e^{-a}$ . If  $a$  is not too large, the latter bound is of correct order. In this case, by Remark 2.2.3,  $\frac{1 - e^{-\lambda}}{\lambda}$  is of order  $\min\left\{1, \frac{1}{n} e^a\right\}$ . As it turns out, this case is an upper bound.

**Lemma 2.4.7.** *There exists a constant  $B$ , such that, letting  $\lambda = n\left(1 - \frac{1}{n}\right)^{an}$ , we have  $\frac{1 - e^{-\lambda}}{\lambda} \leq \min\left\{1, \frac{B e^a}{n}\right\}$  for all  $n \geq 2$  and  $a \geq 0$ .*

**Corollary 2.4.8.** *Let  $W$  be as above. There exists a constant  $C$ , such that*

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq C \left( (1 + a) e^{-a} \min\{1, n e^{-a}\} \right).$$

for all  $n \geq 2$  and  $r \geq 0$ , where  $a = r/n$ . In particular, the total variation error in the Poisson approximation tends to zero uniformly in  $n$  as  $a \rightarrow \infty$ . □

**PROOF OF LEMMA 2.4.7.** Applying Taylor expansion, we find that

$$an \log \left(1 - \frac{1}{n}\right) = -a - \frac{a}{2n\left(1 - \frac{\theta}{n}\right)^2} \geq -a - \frac{2a}{n} \quad (2.4.8)$$

for some  $\theta \in [0, 1]$ . Therefore, for  $a \leq n/2$ , we have  $\lambda \geq n e^{-a-1}$ . By Remark 2.2.3, it follows that  $\frac{1 - e^{-\lambda}}{\lambda} \leq \min\left\{1, \frac{1}{\lambda}\right\} \leq \min\left\{1, \frac{e^{a+1}}{n}\right\}$ . On the other hand, for  $a \geq n/2$ , observe that  $\frac{e^{a+1}}{n} \geq \frac{e^{a+1}}{2a} \geq \frac{e^2}{2} \geq 1$ , so that  $\frac{1 - e^{-\lambda}}{\lambda} \leq 1 = \min\left\{1, \frac{e^{a+1}}{n}\right\}$ . This proves the result with  $B = e$ . □

Now consider the more general case considering the number of bins containing exactly  $m$  balls; denote this number by  $W^{(m)}$ . Again,  $W^{(m)} = \sum_{i \in \mathcal{I}} X_i^{(m)}$ , where  $X_i^{(m)}$  is the indicator of the event that the  $i$ -th bin contains exactly  $m$  balls. Clearly,  $X_i^{(m)} \sim \text{Be}(p_i^{(m)})$ , where

$$p_i^{(m)} = \binom{r}{m} q_i^m (1 - q_i)^{r-m}.$$

We have already constructed a coupling of the (unconditional) distribution of  $W$  with its conditional distribution given that the  $i$ -th bin is empty. Continuing from the balls rearranged in the latter way, further pick  $m$  balls uniformly at random and relocate them into the  $i$ -th bin. This makes a coupling with the conditional distribution given  $X_i^{(m)} = 1$ .

**Remark 2.4.9.** Notice that the latter coupling is not entirely optimal. A more refined coupling would go as follows: consider the number of balls in the  $i$ -th bin. If there are exactly  $m$  balls, leave it as it is. If there are more, suitably relocate the excess balls into other bins. If they are less, pick a suitable number of balls from other bins uniformly at random and relocate them into the  $i$ -th bin. However, in this case, the right hand side of (2.3.1) is more complicated to estimate and we do not benefit on the rate of convergence. Therefore, we keep the afore-mentioned two-step coupling.

Let  $\tilde{X}_{ij}^{(m)}$  denote the indicator of the event that the  $j$ -th bin contains exactly  $m$  balls after relocating all balls from the  $i$ -th bin. Next, let  $\tilde{\tilde{X}}_{ij}^{(m)}$  denote the indicator of the same event after relocating  $m$  balls into the  $i$ -th bin. As usual, let  $\tilde{W}_i^{(m)} := \sum_{j \in \mathcal{I} \setminus \{i\}} \tilde{X}_{ij}^{(m)}$  and  $\tilde{\tilde{W}}_i^{(m)} := \sum_{j \in \mathcal{I} \setminus \{i\}} \tilde{\tilde{X}}_{ij}^{(m)}$ . The (unconditional) distribution of  $1 + \tilde{\tilde{W}}_i^{(m)}$  agrees with the conditional distribution of  $W$  given  $X_i = 1$ .

Now estimate

$$\begin{aligned} |W^{(m)} - \tilde{\tilde{W}}_i^{(m)}| &\leq |W^{(m)} - \tilde{W}_i^{(m)}| + |\tilde{W}_i^{(m)} - \tilde{\tilde{W}}_i^{(m)}| \\ &\leq X_i + \sum_{j \in \mathcal{I} \setminus \{i\}} \left[ |X_j^{(m)} - \tilde{X}_{ij}^{(m)}| + |\tilde{X}_{ij}^{(m)} - \tilde{\tilde{X}}_{ij}^{(m)}| \right] \\ &= X_i + \sum_{j \in \mathcal{I} \setminus \{i\}} \left[ (X_j^{(m)} - X_j^{(m)} \tilde{X}_{ij}^{(m)}) + (\tilde{X}_{ij}^{(m)} - X_j^{(m)} \tilde{X}_{ij}^{(m)}) \right. \\ &\quad \left. + (\tilde{X}_{ij}^{(m)} - \tilde{X}_{ij}^{(m)} \tilde{\tilde{X}}_{ij}^{(m)}) + (\tilde{\tilde{X}}_{ij}^{(m)} - \tilde{X}_{ij}^{(m)} \tilde{\tilde{X}}_{ij}^{(m)}) \right] \end{aligned}$$

(notice that  $X_j^{(m)} \geq X_j^{(m)} \tilde{X}_{ij}^{(m)}$  and similarly for the other differences).

Now compute the expectations of all random variables, i. e., the probabilities of all underlying events.

- Recall that  $\mathbb{E} X_j^{(m)} = \mathbb{E} p_i^{(m)} = \binom{r}{m} q_i^m (1 - q_i)^{r-m}$ .
- After the first rearrangement, each ball is in the  $j$ -th bin with probability  $q_j / (1 - q_i)$ , where the balls are independent. Therefore,

$$\mathbb{E} \tilde{\tilde{X}}_{ij}^{(m)} = \binom{r}{m} \left( \frac{q_j}{1 - q_i} \right)^m \left( 1 - \frac{q_j}{1 - q_i} \right)^{r-m}.$$

- In the second rearrangement,  $m$  balls are selected to be relocated to the  $i$ -th bin. Given the choice of these  $m$  balls, each one of the rest is in the  $j$ -th bin with probability  $q_j/(1 - q_i)$ , where the balls are again independent. Therefore,

$$\mathbb{E} \tilde{X}_{ij}^{(m)} = \binom{r-m}{m} \left( \frac{q_j}{1-q_i} \right)^m \left( 1 - \frac{q_j}{1-q_i} \right)^{r-2m}.$$

- The product  $X_j^{(m)} \tilde{X}_{ij}^{(m)}$  is the indicator of the event that in the  $j$ -th bin, there are exactly  $m$  balls before and after the first rearrangement. There are  $\binom{r}{m}$  choices of balls that first land in the  $j$ -th bin, and for each of them, this occurs with probability  $q_j$ . Given the choice of these  $m$  balls, each one of the remaining  $r - m$  balls lands in the  $j$ -th bin with probability  $q_j/(1 - q_i)$  after the first rearrangement, and the balls are independent. Therefore,

$$\mathbb{E}[X_j^{(m)} \tilde{X}_{ij}^{(m)}] = \binom{r}{m} q_j^m \left( 1 - \frac{q_j}{1-q_i} \right)^{r-m}.$$

- The product  $\tilde{X}_{ij}^{(m)} \tilde{\tilde{X}}_{ij}^{(m)}$  is the indicator of the event that in the  $j$ -th bin, there are exactly  $m$  balls before and after the second rearrangement. Recall that in the second rearrangement,  $m$  balls are selected to be relocated to the  $i$ -th bin. We can choose these balls first. None of these balls should be in the  $j$ -th bin after the first rearrangement. Given the selection of these  $m$  balls, exactly  $m$  among the remaining  $r - m$  balls should be in the  $j$ -th bin after the first (and the second) rearrangement. Since the selection of the balls to be relocated in the second turn is independent of the location of the balls after the first rearrangement, we find that

$$\mathbb{E}[\tilde{X}_{ij}^{(m)} \tilde{\tilde{X}}_{ij}^{(m)}] = \binom{r-m}{m} \left( \frac{q_j}{1-q_i} \right)^m \left( 1 - \frac{q_j}{1-q_i} \right)^{r-m}.$$

Now estimate

$$\begin{aligned} \mathbb{E}(X_j^{(m)} - X_j^{(m)} \tilde{X}_{ij}^{(m)}) &= \binom{r}{m} q_j^m \left[ (1 - q_j)^{r-m} - \left( 1 - \frac{q_j}{1-q_i} \right)^{r-m} \right] \\ &= \binom{r}{m} q_j^m \left[ (1 - q_j)^{r-m} - \left( 1 - q_j - \frac{q_i q_j}{1-q_i} \right)^{r-m} \right] \\ &= \binom{r}{m} q_j^m (1 - q_j)^{r-m} \left[ 1 - \left( 1 - \frac{q_i q_j}{(1-q_i)(1-q_j)} \right)^{r-m} \right] \\ &\leq (r-m) \frac{q_i q_j}{(1-q_i)(1-q_j)} \binom{r}{m} q_j^m (1 - q_j)^{r-m}, \end{aligned}$$

$$\begin{aligned}
\mathbb{E}(\tilde{X}_{ij}^{(m)} - X_j^{(m)} \tilde{X}_{ij}^{(m)}) &= \binom{r}{m} \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m} \left[ \left(\frac{q_j}{1 - q_i}\right)^m - q_j^m \right] \\
&= \binom{r}{m} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m} [1 - (1 - q_i)^m] \\
&\leq m q_i \binom{r}{m} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m},
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}(\tilde{X}_{ij}^{(m)} - \tilde{X}_{ij}^{(m)} \tilde{X}_{ij}^{(m)}) &= \left[ \binom{r}{m} - \binom{r-m}{m} \right] \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m} \\
&= \sum_{k=r-m}^{r-1} \binom{k}{m-1} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m} \\
&\leq m \binom{r-1}{m-1} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m} \\
&= \frac{m^2}{r} \binom{r}{m} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-m}.
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}(\tilde{X}_{ij}^{(m)} - \tilde{X}_{ij}^{(m)} \tilde{X}_{ij}^{(m)}) &= \binom{r-m}{m} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-2m} \left[1 - \left(1 - \frac{q_j}{1 - q_i}\right)^m\right] \\
&\leq m \frac{q_j}{1 - q_i} \binom{r-m}{m} \left(\frac{q_j}{1 - q_i}\right)^m \left(1 - \frac{q_j}{1 - q_i}\right)^{r-2m}.
\end{aligned}$$

Collecting all together, we obtain (provided that  $m \geq 1$ )

$$\mathbb{E}|W^{(m)} - \tilde{W}_i^{(m)}| \leq p_i^{(m)} + \sum_{j \in \mathcal{J} \setminus \{i\}} p_{ij}^{(m)} \left[ \frac{(r-m)q_i q_j}{1 - q_j} + m q_i + m \frac{q_j}{1 - q_i} + \frac{m^2}{r} \right],$$

where

$$p_{ij}^{(m)} := \binom{r}{m} \left(\frac{q_j}{1 - q_i}\right)^m (1 - q_j)^{r-2m}.$$

Letting  $\lambda_m = \sum_{i \in \mathcal{J}} p_i^{(m)}$  and applying Theorem 2.3.2, we obtain the following result:

**Proposition 2.4.10.** *For  $m \geq 1$ , the distribution of the number  $W^{(m)}$  of the bins with exactly  $m$  balls satisfies*

$$\begin{aligned}
d_{\text{TV}}(\mathcal{L}(W^{(m)}), \text{Po}(\lambda_m)) &\leq \frac{1 - e^{-\lambda_m}}{\lambda_m} \sum_{i \in \mathcal{J}} \left( (p_i^{(m)})^2 \right. \\
&\quad \left. + \sum_{j \in \mathcal{J} \setminus \{i\}} p_i^{(m)} p_{ij}^{(m)} \left[ \frac{(r-m)q_i q_j}{1 - q_j} + m q_i + m \frac{q_j}{1 - q_i} + \frac{m^2}{r} \right] \right). \tag{2.4.9}
\end{aligned}$$

□

Now consider again the case where all bins are hit with equal probabilities, that is,  $q_i = \frac{1}{n}$ , where  $n$  denotes the number of bins. In this case, (2.4.9) reduces to

$$d_{\text{TV}}\left(\mathcal{L}(W^{(m)}), \text{Po}(\lambda_m)\right) \leq \frac{1 - e^{-\lambda_m}}{\lambda_m} \left[ n (p_1^{(m)})^2 + \left( r + 2m(n-1) + \frac{m^2 n(n-1)}{r} \right) p_1^{(m)} p_{12}^{(m)} \right],$$

where

$$p_1^{(m)} = \binom{r}{m} \frac{1}{n^m} \left(1 - \frac{1}{n}\right)^{r-m}, \quad p_{12}^{(m)} = \binom{r}{m} \frac{1}{(n-1)^m} \left(1 - \frac{1}{n}\right)^{r-2m}, \quad \lambda_m = n p_1^{(m)}.$$

Suppose that  $n \geq 2$  and  $r \geq m$ . Letting  $a := r/n$ , observe that

$$p_i^{(m)} \leq \frac{r^m}{m!} \frac{1}{n^m} e^{-a} \left(\frac{n}{n-1}\right)^m \leq \frac{(2a)^m}{m!} e^{-a}$$

and

$$p_{ij}^{(m)} \leq \frac{r^m}{m!} \frac{1}{n^m} e^{-a} \left(\frac{n}{n-1}\right)^{3m} \leq \frac{(8a)^m}{m!} e^{-a}$$

Thus,  $\lambda$  is of order at most  $n a^m e^{-a}$ . If this is the actual order, then, by Remark 2.2.3,  $\frac{1-e^{-\lambda}}{\lambda}$  is of order  $\min\{1, \frac{e^a}{n a^m}\}$ . Again, this order is an upper bound, as shown in the following generalization of Lemma 2.4.7:

**Lemma 2.4.11.** *For each  $m \in \mathbb{N}_0$ , there exists a constant  $B_m$ , such that, letting  $\lambda = \binom{an}{m} \frac{1}{n^{m-1}} \left(1 - \frac{1}{n}\right)^{an-m}$ , we have  $\frac{1-e^{-\lambda}}{\lambda} \leq \min\{1, \frac{B_m e^a}{n a^m}\}$  for all  $n \geq 2$  and  $a > 0$ .*

**Corollary 2.4.12.** *Let  $W^{(m)}$  and  $\lambda_m$  be as above. For each  $m \in \mathbb{N}_0$ , there exists a constant  $C_m$ , such that*

$$d_{\text{TV}}\left(\mathcal{L}(W^{(m)}), \text{Po}(\lambda_m)\right) \leq C_m \left( (a^{m+1} + a^{m-1}) e^{-a} \min\{1, n a^m e^{-a}\} \right).$$

for all  $n \geq 2$  and  $r \geq 0$ , where  $a = r/n$ . In particular, the total variation error in the Poisson approximation tends to zero uniformly in  $n$  as  $a \rightarrow \infty$ . For  $m \geq 2$ , this is also true as  $a \rightarrow 0$ .  $\square$

**Remark 2.4.13.** For  $m = 1$ , the error does not tend to zero uniformly in  $n$  as  $a \rightarrow 0$ . As a counterexample, consider the case where  $r \geq 1$  is constant, while  $n$  tends to the infinity. In this case, the number of bins with exactly one ball tends to the constant  $r$ . Since for  $r \geq 1$ , the total variation distance between the Dirac measure at  $r$  and any Poisson distribution is uniformly bounded away from zero, the total variation error in the Poisson approximation cannot tend to zero.

**PROOF OF LEMMA 2.4.11.** In view of Lemma 2.4.7, it suffices to prove the assertion for  $m \geq 1$ . First, assume that  $a \leq n/2$ . Recalling (2.4.8), we find that

$$\begin{aligned} \lambda &\geq n \frac{r^m}{m!} \left(1 - \frac{1}{r}\right) \left(1 - \frac{2}{r}\right) \cdots \left(1 - \frac{m-1}{r}\right) \frac{1}{n^m} e^{-a-2a/n} \\ &\geq n \frac{r^m}{m!} \left(1 - \frac{1}{m}\right) \left(1 - \frac{2}{m}\right) \cdots \left(1 - \frac{m-1}{m}\right) \frac{1}{n^m} e^{-a-1} \\ &\geq n \left(\frac{a}{m}\right)^m e^{-a-1}. \end{aligned}$$

Recalling Remark 2.2.3, this implies  $\frac{1-e^{-\lambda}}{\lambda} \leq \min\{1, \frac{e^{a+1}}{n} (\frac{m}{a})^m\}$ . On the other hand, for  $a \geq n/2$ , observe that  $\frac{e^{a+1}}{n} (\frac{m}{a})^m \geq \frac{m^m e^{a+1}}{2a^{m+1}} \geq \frac{m^m e^{m+2}}{2(m+1)^{m+1}} \geq \frac{e^{m+1}}{2(m+1)} \geq \frac{e}{2} \geq 1$ , so that  $\frac{1-e^{-\lambda}}{\lambda} \leq 1 = \min\{1, \frac{e^{a+1}}{n} (\frac{m}{a})^m\}$ . This proves the result with  $B_m = e m^m$ .  $\square$

# Chapter 3

## Normal approximation

### 3.1 Decomposable random variables

In Section 1.2, we derived that for a standard normal random variable  $W$ , the *Stein expectation*

$$\mathbb{E}[f'(W) - f(W)W] \tag{3.1.1}$$

vanishes for all continuously differentiable functions  $f: \mathbb{R} \rightarrow \mathbb{R}$  of polynomial growth. Here, we shall first show that the Stein expectation can be small for random variables  $W$  featuring a certain dependence structure. In the next section, we shall show that this implies proximity to the standard normal distribution in a certain metric.

In Subsection 2.3.2, we have shown that Poisson approximation by Stein's method works well for locally dependent random variables and, more generally, random variables which can be decomposed as in (2.3.2). Here, we adjust this approach to the case of normal approximation. This adjustment is due to Barbour, Karoński and Ruciński [6].

Consider a random variable  $W$  with  $\mathbb{E}W = 0$  and  $\text{var}(W) = 1$  (this, of course, means that  $\mathbb{E}(W^2) < \infty$ ). Suppose that

$$W = \sum_{i \in \mathcal{I}} X_i, \tag{3.1.2}$$

where  $\mathcal{I}$  is a countable set and where we assume that  $\sum_{i \in \mathcal{I}} \mathbb{E}|X_i| < \infty$ ,  $\sum_{i \in \mathcal{I}} \mathbb{E}(|X_i||W|) < \infty$  and  $\mathbb{E}X_i = 0$  for all  $i \in \mathcal{I}$  (notice that the first condition guarantees the almost sure existence of the random sum  $\sum_{i \in \mathcal{I}} X_i$ ). Then the variance can be expressed as

$$1 = \text{var}(W) = \mathbb{E}(W^2) = \sum_{i \in \mathcal{I}} \mathbb{E}(X_i W).$$

For a continuously differentiable function  $f$  with bounded derivative, this leads to the following expression of the Stein expectation:

$$\mathbb{E}[f'(W) - f(W)W] = \sum_{i \in \mathcal{I}} \mathbb{E}\left[f'(W)\mathbb{E}(X_i W) - f(W)X_i\right]$$



(notice that  $f$  is of at most linear growth, that is, there exist  $C_0$  and  $C_1$ , such that  $|f(w)| \leq C_0 + C_1|w|$  for all  $w \in \mathbb{R}$ ). Now suppose that for each  $i \in \mathcal{I}$ ,  $W$  can be decomposed as

$$W = W_i + R_i, \quad (3.1.3)$$

where  $W_i$  is independent of  $X_i$  and where  $\mathbb{E}(|X_i| |R_i|) < \infty$ . Then we have  $\mathbb{E}(X_i W) = \mathbb{E}(X_i R_i)$ . Next, by the fundamental theorem of calculus, we have

$$f(W) = f(W_i) + \int_{W_i}^W f'(t) dt = f(W_i) + \int_0^1 f'(W_i + tR_i) R_i dt.$$

Assuming in addition that  $\sum_{i \in \mathcal{I}} \mathbb{E}(|X_i| |R_i|) < \infty$  and combining all together, we express the Stein expectation as

$$\mathbb{E}[f'(W) - f(W)W] = \sum_{i \in \mathcal{I}} \mathbb{E} \left[ f'(W) \mathbb{E}(X_i R_i) - f(W_i) X_i - \int_0^1 f'(W_i + tR_i) X_i R_i dt \right].$$

By independence and since  $\mathbb{E} X_i = 0$ , the second term vanishes. Taking a random variable  $\theta_1$ , which is uniformly distributed over  $[0, 1]$  and independent of all other random variables, we can rewrite the Stein expectation as

$$\mathbb{E}[f'(W) - f(W)W] = \sum_{i \in \mathcal{I}} \mathbb{E} \left[ f'(W) \mathbb{E}(X_i R_i) - f'(W_i + \theta_1 R_i) X_i R_i \right].$$

Now suppose in addition that the random variables  $R_i$  can be expressed as sums

$$R_i = \sum_{j \in \mathcal{J}_i} X_{ij},$$

where we assume that  $\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}|) < \infty$ . Then we can further rewrite the Stein expectation as

$$\mathbb{E}[f'(W) - f(W)W] = \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ f'(W) \mathbb{E}(X_i X_{ij}) - f'(W_i + \theta_1 R_i) X_i X_{ij} \right].$$

Next, suppose that for each  $i \in \mathcal{I}$  and  $j \in \mathcal{J}_i$ ,  $W_i$  can be further decomposed as

$$W_i = W_{ij} + R_{ij}, \quad (3.1.4)$$

where  $W_{ij}$  is independent of the pair  $(X_i, X_{ij})$ . Assuming in addition that  $f$  is twice continuously differentiable with bounded second derivative and

$$\begin{aligned} \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|R_i + R_{ij}| &< \infty, \\ \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}| |\theta_1 R_i + R_{ij}|) &< \infty, \end{aligned}$$

we can further expand the Stein expectation as

$$\begin{aligned} & \mathbb{E}[f'(W) - f(W)W] \\ &= \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ f'(W_{ij}) \mathbb{E}(X_i X_{ij}) + f''(W_{ij} + \theta_2 R_{ij} + \theta_2 R_i)(R_i + R_{ij}) \mathbb{E}(X_i X_{ij}) \right. \\ & \quad \left. - f'(W_{ij}) X_i X_{ij} - f''(W_{ij} + \theta_2 R_{ij} + \theta_1 \theta_2 R_i) X_i X_{ij} (\theta_1 R_i + R_{ij}) \right], \end{aligned}$$

where  $\theta_2$  is another random variable which is uniformly distributed over  $[0, 1]$  and is independent of all other random variables. Because of independence, the first and the third term cancel, so that the Stein expectation can be expressed as

$$\begin{aligned} \mathbb{E}[f'(W) - f(W)W] &= \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ f''(W_{ij} + \theta_2 R_{ij} + \theta_2 R_i)(R_i + R_{ij}) \mathbb{E}(X_i X_{ij}) \right. \\ & \quad \left. - f''(W_{ij} + \theta_2 R_{ij} + \theta_1 \theta_2 R_i) X_i X_{ij} (\theta_1 R_i + R_{ij}) \right]. \end{aligned} \quad (3.1.5)$$

However, the assumptions implying the preceding identity can be relaxed to a certain extent. First, the assumption that  $f$  is twice continuously differentiable with bounded second derivative can be replaced by the assumption that  $f$  is differentiable and  $f'$  is Lipschitz (see Proposition B.1.9). This is because the fundamental theorem of calculus remains true for Lipschitz test functions (and more generally for the absolutely continuous functions): see Section B.1. This allows us to bound the Stein expectation in terms of

$$M_2(f) := \text{ess sup } |f''| = \sup_{x \neq y} \frac{|f'(x) - f'(y)|}{|x - y|},$$

where  $f'$  is the classical derivative of  $f$  and  $f''$  is an almost-everywhere derivative of  $f'$ .

In addition to the relaxed assumption on differentiability of  $f$ , one can also drop certain other assumptions. The following assertion makes it precise.

**Proposition 3.1.1.** *Let  $W$  be decomposed as follows:*

$$\begin{aligned} W &= \sum_{i \in \mathcal{I}} X_i, \\ W &= W_i + R_i, \text{ where } W_i \text{ is independent of } X_i, \\ R_i &= \sum_{j \in \mathcal{J}_i} X_{ij}, \\ W_i &= W_{ij} + R_{ij}, \text{ where } W_{ij} \text{ is independent of } (X_i, X_{ij}). \end{aligned}$$

Next, suppose that  $\mathbb{E} X_i = 0$  for all  $i \in \mathcal{I}$  and  $\text{var}(W) = 1$ , and that

$$\begin{aligned} \sum_{i \in \mathcal{I}} \mathbb{E} |X_i| < \infty, \quad \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} (|X_i| |X_{ij}|) < \infty, \quad \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} (|X_i| |X_{ij}| |R_{ij}|) < \infty, \\ \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} (|X_i| |X_{ij}| |R_i + R_{ij}|) < \infty, \quad \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} (|X_i| |X_{ij}|) \mathbb{E} |R_i + R_{ij}| < \infty. \end{aligned} \quad (3.1.6)$$

Finally, take a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with  $M_2(f) < \infty$ , where  $M_2(f)$  is defined as in (B.1.3). Then (3.1.5) remains true and we can estimate

$$\left| \mathbb{E}[f'(W) - f(W)W] \right| \leq M_2(f) \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ |X_i| |X_{ij}| \left( \mathbb{E} |R_i + R_{ij}| + \frac{1}{2} |R_{ij}| + \frac{1}{2} |R_i + R_{ij}| \right) \right]. \quad (3.1.7)$$

Before proving the preceding assertion, we formulate a couple of remarks. First, we no longer assume that  $f'$  is bounded. In particular, this allows us to apply Proposition 3.1.1 with the function  $f(w) = w|w|$ : this function is differentiable with derivative  $f'(w) = 2|w|$ , which is Lipschitz, but not bounded. This implies that  $W$  has finite third absolute moment. More precisely, with  $f$  as above, the following result is immediate from Proposition 3.1.1:

**Corollary 3.1.2.** *If  $W$  is as in Proposition 3.1.1, then*

$$\mathbb{E}(|W|^3) \leq 2 \mathbb{E} |W| + 2 \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ |X_i| |X_{ij}| \left( |R_i + R_{ij}| + \frac{1}{2} |R_{ij}| + \frac{1}{2} |R_i + R_{ij}| \right) \right].$$

□

**Remark 3.1.3.** If  $W$  is standard normal, then  $\mathbb{E} |W| = 2/\sqrt{2\pi}$  and  $\mathbb{E}(|W|^3) = 4/\sqrt{2\pi}$ , so that  $\mathbb{E}(|W|^3) = 2 \mathbb{E} |W|$ .

Now we turn to the proof of Proposition 3.1.1. As the first step, we formulate and prove the the following auxiliary result:

**Lemma 3.1.4.** *Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be absolutely continuous and let  $h: [0, \infty) \rightarrow [0, \infty)$  be non-decreasing. Suppose that  $|f'| \leq h$ . Finally, let a random variable  $W$  be as in Proposition 3.1.1, and let  $\theta_1$  be uniformly distributed over  $[0, 1]$  and independent of all other random variables. If*

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} [h(|W_i + \theta_1 R_i|) |X_i| |X_{ij}|] < \infty, \quad (3.1.8)$$

then  $\mathbb{E} |f(W)W| < \infty$  and

$$\mathbb{E}[f(W)W] = \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[f'(W_i + \theta_1 R_i) X_i X_{ij}]. \quad (3.1.9)$$

**PROOF.** Equation (3.1.9) is derived at the beginning of this section, but under stronger conditions. A closer look reveals that the calculations are still valid if  $W$  is as in Proposition 3.1.1 and  $f$  is Lipschitz (recall that in this case,  $f$  has an almost-everywhere derivative and satisfies the fundamental theorem of calculus – see Section B.1).

Now take any absolutely continuous function  $f$  with  $|f'| \leq h$ . For each  $n \in \mathbb{N}$ , define function  $\psi_n: [0, \infty) \rightarrow [0, 1]$  as  $\psi_n(t) := 1$  for  $t \leq n$ ,  $\psi_n(t) := 2 - t/n$  for  $n \leq t \leq 2n$  and  $\psi_n(t) := 0$  for  $t \geq 2n$ . Observe that  $\psi_n$  is well defined and that for each fixed  $n$ , the expression  $t \psi_n(t)$  is bounded in  $t$ . Clearly,  $\psi$  is absolutely continuous with  $\psi'_n(t) = 0$  for

$t < n$ ,  $\psi'_n(t) = -1/n$  for  $n < t < 2n$  and  $\psi'_n(t) = 0$  for  $t > 2n$ . Observe that  $\psi'$  can be extended to  $[0, \infty)$  so that the expression  $t |\psi'_n(t)|$  is uniformly bounded in  $t$  and  $n$ .

Now let  $f_n(w) := f(\psi_n(|w|)w)$ . Applying the chain rule (for details on the validity in the context of absolutely continuous functions, see Corollary 6.5.4 of [20]), we find that  $f_n$  is absolutely continuous with

$$f'_n(w) = f'(\psi_n(|w|)w) (\psi'_n(|w|)|w| + \psi_n(|w|)).$$

Since  $t |\psi'_n(t)|$  is uniformly bounded in  $t$  and  $n$  and since  $h$  is non-decreasing, there exists a constant  $C$ , such that  $|f'_n(w)| \leq C h(|w|)$  for all  $n$  and  $w$ .

For each fixed  $n$ , the expression  $\psi_n(|w|)w$  is bounded in  $w \in \mathbb{R}$ . Since  $f'$  is bounded on bounded sets,  $|f'_n(w)|$  is also bounded in  $w \in \mathbb{R}$ . Therefore, (3.1.9) applies with  $f_n$  in place of  $f$ .

Now observe that the functions  $f_n$  converge pointwise to  $f$  and that their derivatives  $f'_n$  converge pointwise to  $f'$  as well. Recalling (3.1.8) and applying the dominated convergence theorem, we obtain

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}[f_n(W)W] &= \lim_{n \rightarrow \infty} \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[f'_n(W_i + \theta_1 R_i) X_i X_{ij}] \\ &= \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[f'(W_i + \theta_1 R_i) X_i X_{ij}]. \end{aligned} \quad (3.1.10)$$

Now consider the function  $\tilde{f}: \mathbb{R} \rightarrow \mathbb{R}$  defined by  $\tilde{f}(w) := \int_0^w h(|t|) dt$ . Observe that  $\tilde{f}$  is absolutely continuous with  $\tilde{f}'(w) = h(|w|)$  and that  $\tilde{f}(w)w \geq 0$  for all  $w \in \mathbb{R}$ . Letting  $\tilde{f}_n(w) := \tilde{f}(\psi_n(|w|)w)$ , we also have  $\tilde{f}_n(w)w \geq 0$  for all  $w \in \mathbb{R}$ . In addition, (3.1.10) applies with  $\tilde{f}_n$  and  $\tilde{f}$  in place of  $f_n$  and  $f$ . Combining with Fatou's lemma, we find that

$$\begin{aligned} \mathbb{E}[\tilde{f}(W)W] &\leq \lim_{n \rightarrow \infty} \mathbb{E}[\tilde{f}_n(W)W] \\ &= \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[\tilde{f}'(W_i + \theta_1 R_i) X_i X_{ij}] \\ &\leq \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[h(|W_i + \theta_1 R_i|) |X_i| |X_{ij}|] \\ &< \infty. \end{aligned} \quad (3.1.11)$$

Now estimate

$$\begin{aligned} |f_n(w)| &= |f(\psi_n(|w|)w)| \leq |f(0)| + \int_0^{\psi_n(|w|)|w|} h(t) dt \leq |f(0)| + \int_0^{|w|} h(t) dt, \\ |f_n(w)w| &\leq |f(0)| |w| + |w| \int_0^{|w|} h(t) dt \leq |F(0)| |w| + \tilde{f}(w)w. \end{aligned}$$

Recalling (3.1.11), it follows that the sequence of random variables  $f_n(W)W$  is dominated by a non-negative random variable with finite expectation. Applying the dominated convergence theorem and combining with (3.1.10), finiteness of  $\mathbb{E}|f(W)W|$  along with (3.1.9) follows.  $\square$

PROOF OF PROPOSITION 3.1.1. We shall go from (3.1.5) backwards. Define

$$\begin{aligned} \rho := \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ f''(W_{ij} + \theta_2 R_{ij} + \theta_2 R_i)(R_i + R_{ij}) \mathbb{E}(X_i X_{ij}) \right. \\ \left. - f''(W_{ij} + \theta_1 \theta_2 R_i + \theta_2 R_{ij}) X_i X_{ij} (\theta_1 R_i + R_{ij}) \right]. \end{aligned} \quad (3.1.12)$$

Noting also that

$$\theta_1 R_i + R_{ij} = (1 - \theta_1) R_{ij} + \theta_1 (R_i + R_{ij}), \quad (3.1.13)$$

and applying (3.1.6), we find that the right hand side of (3.1.12) exists, along with the bound

$$|\rho| \leq M \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ |X_i| |X_{ij}| \left( \mathbb{E} |R_i + R_{ij}| + \frac{1}{2} |R_{ij}| + \frac{1}{2} |R_i + R_{ij}| \right) \right].$$

It remains to show that  $\mathbb{E}|f'(W) - f(W)W| < \infty$  and  $\rho = \mathbb{E}[f'(W) - f(W)W]$ . Observe first that

$$\begin{aligned} \mathbb{E}(|X_i| |X_{ij}| |W_{ij}|) &= \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|W_{ij}| = \\ &\leq \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|W| + \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|R_i + R_{ij}|. \end{aligned}$$

Applying (3.1.6), we find that

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}| |W_{ij}|) = \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|W_{ij}| < \infty.$$

Noting also that

$$|f'(w)| \leq |f'(0)| + M|w|, \quad (3.1.14)$$

we find that

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|f'(W_{ij})| |X_i| |X_{ij}|) < \infty$$

and

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[|f'(W_{ij})|] |\mathbb{E}(X_i X_{ij})| < \infty.$$

Moreover,

$$\begin{aligned} \rho &= \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ f'(W_{ij}) \mathbb{E}(X_i X_{ij}) + f''(W_{ij} + \theta_2 R_{ij} + \theta_2 R_i)(R_i + R_{ij}) \mathbb{E}(X_i X_{ij}) \right. \\ &\quad \left. - f'(W_{ij}) X_i X_{ij} - f''(W_i + \theta_2 R_{ij} + \theta_1 \theta_2 R_i) X_i X_{ij} (\theta_1 R_i + R_{ij}) \right] \end{aligned}$$

(and all expectations exist and the sum converges). Next, since the fundamental theorem of calculus holds for  $f'$  and  $f''$ , we have

$$\begin{aligned} &\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[|f'(W)|] |\mathbb{E}(X_i X_{ij})| \\ &\leq \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \left( \mathbb{E}[|f'(W_{ij})|] |\mathbb{E}(X_i X_{ij})| + M \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|R_i + R_{ij}| \right) < \infty. \end{aligned}$$

Similarly, recalling (3.1.13), we have

$$\begin{aligned} & \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[|f'(W_i + \theta_1 R_i)| |X_i| |X_{ij}|] \\ & \leq \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \left( \mathbb{E}[|f'(W_{ij})| |X_i| |X_{ij}|] + \frac{1}{2} M \mathbb{E}[|X_i| |X_{ij}| (|R_{ij}| + |R_i + R_{ij}|)] \right) < \infty. \end{aligned}$$

In particular, the choice  $f(w) = \frac{1}{2}w^2$  gives

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|W_i + \theta_1 R_i| |X_i| |X_{ij}|) < \infty. \quad (3.1.15)$$

Moreover,

$$\rho = \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ f'(W) \mathbb{E}(X_i X_{ij}) - f'(W_i + \theta_1 R_i) X_i X_{ij} \right].$$

Recalling (3.1.14) and applying Lemma 3.1.4 with  $h(t) = f(0) + Mt$ , making use of (3.1.15), we obtain that  $\mathbb{E}|f(W)W| < \infty$  and  $\mathbb{E}[f(W)W] = \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}[f'(W_i + \theta_1 R_i) X_i X_{ij}]$ . Finally, applying Lemma 3.1.4 with  $f(w) = w$ , we find that  $\mathbb{E}(W^2) = \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(X_i X_{ij})$ . Together with (3.1.14), this completes the proof.  $\square$

## 3.2 Solution to the Stein equation

As indicated in Section 1.2, the proximity to the standard normal distribution can be assessed as follows: for a test function  $h$ , find a function  $f$  solving the Stein equation

$$f'(w) - f(w)w = h(w) - \mathcal{N}h \quad (3.2.1)$$

where

$$\mathcal{N}h := \langle h, \mathcal{N}(0, 1) \rangle = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} h(x) e^{-x^2/2} dx.$$

Taking a random variable  $W$ , we then have

$$\mathbb{E}[f'(W) - f(W)W] = \mathbb{E}[h(W)] - \mathcal{N}h. \quad (3.2.2)$$

Estimating the left hand side for sufficiently many test functions  $h$ , we are able to bound the error in the normal approximation with respect to a suitable metric. In particular, if  $W$  is as in Proposition 3.1.1 and if there exists a class  $\mathcal{H}$  of test functions, such that for each  $h \in \mathcal{H}$ , there exists  $f$  which solves (3.2.1) and satisfies  $M_2(f) \leq 1$ , then we have a bound

$$d_{\mathcal{H}}(\mathcal{L}(W), \mathcal{N}(0, 1)) \leq \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ |X_i| |X_{ij}| \left( \mathbb{E}[|R_i + R_{ij}| + \frac{1}{2}|R_{ij}| + \frac{1}{2}|R_i + R_{ij}|] \right) \right],$$

where  $d_{\mathcal{H}}$  is defined as in (A.1.1).

Equation (3.2.1) is an ordinary differential equation of the first order. Such differential equations can be solved in two steps: first, we solve the homogeneous part, then we perform the variation of constant. A simple calculation shows that the solution to the homogeneous part

$$f'_H(w) - f_H(w)w = 0 \quad (3.2.3)$$

is  $f_H(w) = C e^{w^2/2}$ . The solution to the original equation (3.2.1) can be sought as  $f(w) = k(w) e^{w^2/2}$ , where  $k$  is now a function. Another simple calculation shows that  $k$  must be the indefinite integral

$$k(w) = \int (h(x) - \mathcal{N}h) e^{-x^2/2} dx.$$

However, unless  $\lim_{w \rightarrow \pm\infty} k(w) = 0$ ,  $f$  grows very rapidly, so that there is no hope for  $f'$  to be Lipschitz. Since  $\int_{-\infty}^{\infty} (h(x) - \mathcal{N}h) dx = 0$ ,  $\lim_{w \rightarrow -\infty} k(w) = 0$  implies  $\lim_{w \rightarrow \infty} k(w) = 0$  and the 'tame' solution to (3.2.1) can be expressed as

$$f(w) = e^{w^2/2} \int_{-\infty}^w (h(x) - \mathcal{N}h) e^{-x^2/2} dx = e^{w^2/2} \int_w^{\infty} (\mathcal{N}h - h(x)) e^{-x^2/2} dx. \quad (3.2.4)$$

We summarize the preceding calculations into the following statement:

**Proposition 3.2.1.** *For any measurable function  $h: \mathbb{R} \rightarrow \mathbb{R}$  with  $\mathcal{N}|h| < \infty$ , the function  $f$  defined by (3.2.4) is an almost-everywhere solution to the Stein equation (3.2.1), i. e.,  $f$  is absolutely continuous and the function  $w \mapsto f(w)w + h(w) - \mathcal{N}h$  is an almost-everywhere derivative of  $f$ . If  $h$  is continuous,  $f$  is a classical solution, i. e., continuously differentiable and (3.2.1) holds for all  $w \in \mathbb{R}$ .*

Denoting the standard normal density by

$$\phi(z) := \frac{1}{\sqrt{2\pi}} e^{-z^2/2},$$

we can also write

$$f(w) = \frac{1}{\phi(w)} \int_{-\infty}^w (h(x) - \mathcal{N}h) \phi(x) dx = \frac{1}{\phi(w)} \int_w^{\infty} (\mathcal{N}h - h(x)) \phi(x) dx.$$

Of course, we can also take affine combinations of the two forms, that is, for each  $a \in \mathbb{R}$ , we have

$$\begin{aligned} f(w) &= \frac{1}{\phi(w)} \left( (1-a) \int_{-\infty}^w (h(x) - \mathcal{N}h) \phi(x) dx + a \int_w^{\infty} (\mathcal{N}h - h(x)) \phi(x) dx \right) \\ &= \frac{1}{\phi(w)} \left( (1-a) \int_{-\infty}^w h(x) \phi(x) dx - a \int_w^{\infty} h(x) \phi(x) dx \right) \\ &\quad + [a(1 - \Phi(w)) - (1-a)\Phi(w)] \mathcal{N}h, \end{aligned} \quad (3.2.5)$$

where

$$\Phi(w) := \int_{-\infty}^w \phi(x) dx$$

is the standard normal cumulative distribution function. Choosing  $a = \Phi(w)$ , the second term in the final expression of (3.2.5) vanishes and we obtain the following form:

$$f(w) = \frac{1 - \Phi(w)}{\phi(w)} \int_{-\infty}^w h(x) \phi(x) dx - \frac{\Phi(w)}{\phi(w)} \int_w^{\infty} h(x) \phi(x) dx.$$

Noting that  $1 - \Phi(w) = \Phi(-w)$  and introducing the *Mills ratio*:

$$\psi(w) := \frac{\Phi(w)}{\phi(w)}$$

we can rewrite the solution as

$$f(w) = \psi(-w) \int_{-\infty}^w h(x) \phi(x) dx - \psi(w) \int_w^{\infty} h(x) \phi(x) dx.$$

In view of Proposition 3.1.1, it is beneficial to study the behaviour of the second derivative, including the general case where  $f'$  is absolutely continuous. Differentiating the preceding formula, we obtain

$$\begin{aligned} f'(w) &= -\psi'(-w) \int_{-\infty}^w h(x) \phi(x) dx + \psi(-w) h(w) \phi(w) \\ &\quad - \psi'(w) \int_w^{\infty} h(x) \phi(x) dx + \psi(w) h(w) \phi(w) \\ &= h(w) - \psi'(-w) \int_{-\infty}^w h(x) \phi(x) dx - \psi'(w) \int_w^{\infty} h(x) \phi(x) dx, \end{aligned} \tag{3.2.6}$$

where the last equality follows from the identity  $\Phi(w) + \Phi(-w) = 1$  (this is also Proposition C.2.3 for the case  $r = 0$ ).

**Remark 3.2.2.** Similarly as in Proposition 3.2.1, Formula 3.2.6 can be interpreted in two ways: if  $h$  is continuous,  $f'$  defined as in (3.2.6) is the classical derivative of  $f$ , whereas in the general case,  $f$  is absolutely continuous and  $f'$  is an almost-everywhere derivative of  $f$ .

Now assume that  $h$  is absolutely continuous (it must be if so is  $f'$ ) and that it is of polynomial growth. In this case, we can apply the integration by parts formula (see Theorem 6.4.6 of Heil [20]), where we differentiate  $h$  and integrate  $\phi$ . However, in the first integral of (3.2.6), we integrate  $\phi(x)$  to  $\Phi(x)$ , while in the second one, we integrate  $\phi(x)$  to  $-\Phi(-x)$ . Noting that the polynomial growth of  $h$  along with Corollary C.2.2 implies  $\lim_{w \rightarrow -\infty} h(w) \Phi(w) = 0$  and  $\lim_{w \rightarrow \infty} h(w) \Phi(-w) = 0$ , we obtain

$$\begin{aligned} f'(w) &= h(w) - \psi'(-w) h(w) \Phi(w) + \psi'(-w) \int_{-\infty}^w h'(x) \Phi(x) dx \\ &\quad - \psi'(w) h(w) \Phi(-w) - \psi'(w) \int_w^{\infty} h'(x) \Phi(-x) dx \\ &= \psi'(-w) \int_{-\infty}^w h'(x) \Phi(x) dx - \psi'(w) \int_w^{\infty} h'(x) \Phi(-x) dx, \end{aligned} \tag{3.2.7}$$



where the last inequality is due to Proposition C.2.3 for  $r = 1$ . Differentiating (3.2.7), we find that

$$\begin{aligned} f''(w) &= -\psi''(-w) \int_{-\infty}^w h'(x) \Phi(x) dx + \psi'(-w) h'(w) \Phi(w) dw \\ &\quad - \psi''(w) \int_w^{\infty} h'(x) \Phi(-x) dx + \psi'(w) h'(w) \Phi(-w) \\ &= h'(w) - \psi''(-w) \int_{-\infty}^w h'(x) \Phi(x) dx - \psi''(w) \int_w^{\infty} h'(x) \Phi(-x) dx, \end{aligned} \quad (3.2.8)$$

where the last inequality again follows from Proposition C.2.3 for  $r = 1$ . Again, if  $h'$  is continuous,  $f''$  is the classical derivative of  $f'$ , whereas in the general case,  $f'$  is absolutely continuous and  $f''$  is an almost-everywhere derivative of  $f'$ .

The preceding formula allows us to bound  $M_2(f)$  in terms of  $M_1(f)$ . For random variables decomposed according to Barbour, Karoński and Ruciński, this allows us to bound the error in the normal approximation in the Wasserstein metric.

**Theorem 3.2.3.** *For any function  $h: \mathbb{R} \rightarrow \mathbb{R}$  with  $M_1(h) < \infty$ , the function  $f$  defined by (3.2.4) is a classical solution to the Stein equation (3.2.1), which satisfies  $M_2(f) \leq M_1(h)$ .*

**Corollary 3.2.4.** *For a random variable  $W$  decomposed as in Proposition 3.1.1, we have*

$$d_W(\mathcal{L}(W), \mathcal{N}(0, 1)) \leq \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E} \left[ |X_i| |X_{ij}| \left( 2 \mathbb{E} |R_i + R_{ij}| + |R_{ij}| + |R_i + R_{ij}| \right) \right]. \quad (3.2.9)$$

**PROOF OF THEOREM 3.2.3.** Since  $M_1(h) < \infty$ ,  $h$  is absolutely continuous and of linear growth. Consequently,  $\mathcal{N}|h| < \infty$  and Proposition 3.2.1 applies. Moreover, since  $h$  is absolutely continuous and of linear growth, the derivation of (3.2.8) is valid. By Proposition C.1.5 and Proposition C.2.3 for  $r = 2$ , we can estimate

$$|f''(w)| \leq M_1(h) \left[ 1 + \psi''(-w) \int_{-\infty}^w \Phi(x) dx + \psi''(w) \int_w^{\infty} \Phi(-x) dx \right] = 2 M_1(f).$$

Taking the supremum over  $w$ , the proof is complete.  $\square$

**Example 3.2.5.** Let  $\xi_1, \xi_2, \dots$  be independent and identically distributed random variables with  $\mathbb{E} \xi_1 = 0$ ,  $\text{var}(\xi_1) = 1$  and  $\mathbb{E} |\xi_1|^3 < \infty$ . Then the rescaled sum

$$W^{(n)} := \frac{\xi_1 + \xi_2 + \dots + \xi_n}{\sqrt{n}}$$

satisfies  $\mathbb{E} W^{(n)} = 0$  and  $\text{var}(W^{(n)}) = 1$ , and can be trivially decomposed by setting  $\mathcal{I}^{(n)} := \{1, 2, \dots, n\}$ ,  $X_i^{(n)} := \xi_i / \sqrt{n}$ ,  $R_i^{(n)} := X_i^{(n)}$ ,  $\mathcal{J}_i^{(n)} := \{0\}$ ,  $X_{i0}^{(n)} := X_i^{(n)}$  and  $R_{i0}^{(n)} := 0$ . Corollary 3.2.4 yields

$$d_W(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)) \leq \sum_{i=1}^n \left( 2 \mathbb{E} (X_i^{(n)})^2 \mathbb{E} |X_i^{(n)}| + \mathbb{E} |X_i^{(n)}|^3 \right).$$

By Jensen's inequality, we have  $\mathbb{E}|X_i^{(n)}| \leq \left(\mathbb{E}|X_i^{(n)}|^3\right)^{1/3}$  and  $\mathbb{E}(X_i^{(n)})^2 \leq \left(\mathbb{E}|X_i^{(n)}|^3\right)^{2/3}$ , leading to

$$d_W(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)) \leq 3 \sum_{i=1}^n \mathbb{E}|X_i^{(n)}|^3 = \frac{3\mathbb{E}|\xi_1|^3}{\sqrt{n}}.$$

This is the typical rate of convergence in the central limit theorem and can in general not be improved. As an example, take  $\mathbb{P}(\xi_1 = 1) = \mathbb{P}(\xi_1 = -1) = 1/2$ . For even  $n$ ,  $W^{(n)}$  takes values in the set  $\{2k/\sqrt{n}; k \in \mathbb{Z}\}$ . Now define functions  $f_n: \mathbb{R} \rightarrow \mathbb{R}$  by  $f_n(w) := |w - \frac{2k}{\sqrt{n}}|$  for  $\frac{2k-1}{\sqrt{n}} \leq w \leq \frac{2k+1}{\sqrt{n}}$ , where  $k \in \mathbb{Z}$ . Then we have  $\mathbb{E}[f_n(W^{(n)})] = 0$ . Observe that

$$\begin{aligned} \mathcal{N}f_n &= \int_{-\infty}^{\infty} f_n(x) \phi(x) dx \\ &= \frac{1}{2} \int_{-\infty}^{\infty} \left[ f_n(x) \phi(x) dx + f_n\left(x + \frac{1}{\sqrt{n}}\right) \phi\left(x + \frac{1}{\sqrt{n}}\right) \right] dx \\ &= \frac{1}{2} \int_{-\infty}^{\infty} \left[ f_n(x) \phi(x) dx + f_n\left(x + \frac{1}{\sqrt{n}}\right) \phi\left(x + \frac{1}{\sqrt{n}}\right) \right] dx \\ &= \frac{1}{2} \int_{-\infty}^{\infty} \left[ f_n(x) + f_n\left(x + \frac{1}{\sqrt{n}}\right) \right] \phi(x) dx \\ &\quad + \frac{1}{2} \int_{-\infty}^{\infty} f_n\left(x + \frac{1}{\sqrt{n}}\right) \left[ \phi\left(x + \frac{1}{\sqrt{n}}\right) - \phi(x) \right] dx \\ &= \frac{1}{2\sqrt{n}} + \frac{1}{2\sqrt{n}} \int_{-\infty}^{\infty} f_n\left(x + \frac{1}{\sqrt{n}}\right) \int_0^1 \phi'\left(x + \frac{t}{\sqrt{n}}\right) dt dx. \end{aligned}$$

Noting that  $0 \leq f_n(x) \leq 1/\sqrt{x}$  for all  $x$  and that the function  $\phi'$  is bounded, we find that  $\lim_{n \rightarrow \infty} \mathcal{N}f_n \sqrt{n} = 1/2$ . Therefore, for each  $\varepsilon > 0$ , there exists  $n_0 \in \mathbb{N}$ , such that  $d_W(\mathcal{L}(W_n), \mathcal{N}(0, 1)) \geq \frac{1-\varepsilon}{2\sqrt{n}}$ . This proves that the rate of  $1/\sqrt{n}$  cannot be improved.

**Remark 3.2.6.** The bound in Theorem 3.2.3 is sharp: consider functions

$$h_n(w) := \begin{cases} w + \frac{2}{n} & ; w \leq -\frac{1}{n} \\ -w & ; -\frac{1}{n} \leq w \leq \frac{1}{n} \\ w - \frac{2}{n} & ; w \geq \frac{1}{n}. \end{cases}$$

Clearly,  $M_1(h_n) = 1$ . Next, observe that the underlying functions  $f_n''$  are continuous at the origin and we have

$$f_n''(0) = -1 - \psi''(0) \int_{-\infty}^0 h_n'(x) \Phi(x) dx - \psi''(0) \int_0^{\infty} h_n'(x) \Phi(-x) dx.$$

Therefore,

$$M_2(f_n) \geq \left| 1 + \psi''(0) \int_{-\infty}^0 h_n'(x) \Phi(x) dx + \psi''(0) \int_0^{\infty} h_n'(x) \Phi(-x) dx \right|.$$

The functions  $h_n'$  are uniformly bounded and converge pointwise to the constant 1 (except at 0, where they converge to  $-1$ ). By the dominated convergence theorem, the right hand

side converges to

$$\left| 1 + \psi''(0) \int_{-\infty}^0 \Phi(x) dx + \psi''(0) \int_0^{\infty} \Phi(-x) dx \right| = 2$$

by Proposition C.2.3 for  $r = 2$ . Therefore, for each  $\varepsilon > 0$ , there exists  $n$  such that  $M_2(f_n) > 2 - \varepsilon$ .

## 3.3 Applications

### 3.3.1 Local dependence and $U$ -statistics

In Subsection 2.3.2, we already considered locally dependent random variables. The concept of local dependence was expressed in terms of dependence neighbourhoods. Here, we shall take a stronger concept expressed in terms of the *dependence graph*.

For two vertices  $i$  and  $j$  of an undirected graph  $\Gamma$ , we shall denote  $i \sim j$  if they are equal or adjacent, and  $i \not\sim j$  otherwise. For a vertex  $i$  and a set of vertices  $\mathcal{J}$ , we shall denote  $i \sim \mathcal{J}$  if either  $i \in \mathcal{J}$  or there is an edge with one endpoint equal to  $i$  and the other in  $\mathcal{J}$ , and  $i \not\sim \mathcal{J}$  otherwise. Finally, for sets of vertices  $\mathcal{I}$  and  $\mathcal{J}$ , we shall denote  $\mathcal{I} \sim \mathcal{J}$  if either  $\mathcal{I} \cap \mathcal{J} \neq \emptyset$  or there is an edge with one endpoint in  $\mathcal{I}$  and the other in  $\mathcal{J}$ , and  $\mathcal{I} \not\sim \mathcal{J}$  otherwise.

**Definition 3.3.1.** Let  $(X_i)_{i \in \mathcal{I}}$  be a family of random variables and let  $\Gamma$  be an undirected graph with vertex set  $\mathcal{I}$ . The dependence structure of the family  $(X_i)_{i \in \mathcal{I}}$  fits  $\Gamma$  if for any sets  $\mathcal{J}, \mathcal{K} \subseteq \mathcal{I}$  with  $\mathcal{J} \not\sim \mathcal{K}$ , the subfamilies  $(X_j)_{j \in \mathcal{J}}$  and  $(X_k)_{k \in \mathcal{K}}$  are independent. We shall call  $\Gamma$  a *dependence graph* for the family  $(X_i)_{i \in \mathcal{I}}$ .

As in Section 3.1, consider a sum  $W = \sum_{i \in \mathcal{I}} X_i$ , such that  $\sum_{i \in \mathcal{I}} \mathbb{E} |X_i| < \infty$ ,  $\mathbb{E} X_i = 0$  for all  $i \in \mathcal{I}$  and  $\text{var}(W) = 1$ ; in addition, suppose that  $\sum_{i \in \mathcal{I}} \mathbb{E} |X_i|^3 < \infty$ . Next, let  $\Gamma$  be a dependence graph for the family  $(X_i)_{i \in \mathcal{I}}$ . Take  $D < \infty$  and suppose that for each  $i \in \mathcal{I}$ , there are no more than  $D$  vertices  $j$  with  $i \sim j$ . In other words, the degrees of all vertices are strictly less than  $D$ . We shall show that under these conditions,  $W$  can be reasonably decomposed as in Proposition 3.1.1. To achieve this, first set

$$\mathcal{J}_i := \{j \in \mathcal{I} ; i \not\sim j\}, \quad X_{ij} := X_j, \quad R_i := \sum_{j: i \sim j} X_j, \quad W_i := \sum_{j: i \not\sim j} X_j.$$

Since the dependence structure of the family  $(X_i)_{i \in \mathcal{I}}$  fits  $\Gamma$ ,  $W_i$  is independent of  $X_i$ . Next, set

$$R_{ij} := \sum_{k: i \not\sim k, j \sim k} X_k, \quad W_{ij} := \sum_{k: k \not\sim \{i, j\}} X_k$$

and again observe that  $W_{ij}$  is independent of the pair  $(X_i, X_j)$ . Now we bound the quantities appearing in (3.1.6) and (3.1.7). First, by the inequality between the arithmetic and the geometric mean, we have

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} |X_i| |X_{ij}| = \sum_{(i, j): i \sim j} |X_i| |X_j| \leq \frac{1}{2} \sum_{(i, j): i \sim j} (X_i^2 + X_j^2).$$

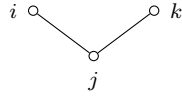
Since the vertex degrees are bounded from above by  $D$ , we have  $|\{j ; i \sim j\}| \leq D$  and  $|\{j ; i \sim j\}| \leq D$ , leading to

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} |X_i| |X_{ij}| \leq D \sum_{i \in \mathcal{I}} \mathbb{E}(X_i^2) \leq \frac{D}{2} \sum_{i \in \mathcal{I}} \left( \mathbb{E} |X_i| + \mathbb{E} |X_i|^3 \right) < \infty,$$

applying again the inequality between the arithmetic and the geometric mean in the second inequality. Next, observe that

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} |X_i| |X_{ij}| |R_{ij}| = \sum_{(i,j,k) \in \mathcal{T}_1} |X_i| |X_j| |X_k| \leq \frac{1}{3} \sum_{(i,j,k) \in \mathcal{T}_1} (|X_i|^3 + |X_j|^3 + |X_k|^3),$$

where  $\mathcal{T}_1 := \{(i, j, k) ; i \sim j, j \sim k, i \not\sim k\}$  as illustrated below:



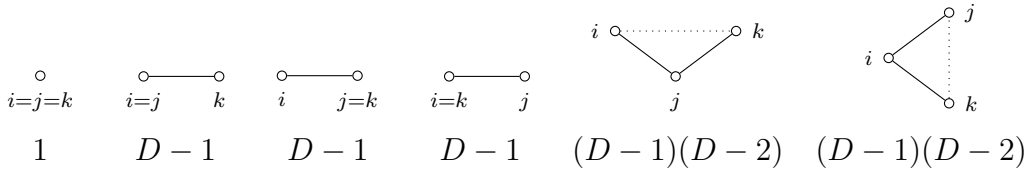
The elements  $i$ ,  $j$  and  $k$  must be distinct. Therefore,  $|\{(j, k) ; (i, j, k) \in \mathcal{T}_1\}| \leq (D-1)(D-2)$  for all  $i$ ,  $|\{(i, k) ; (i, j, k) \in \mathcal{T}_1\}| \leq (D-1)(D-2)$  for all  $j$  and  $|\{(i, j) ; (i, j, k) \in \mathcal{T}_1\}| \leq (D-1)(D-2)$  for all  $k$ , leading to

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}| |R_{ij}|) \leq (D-1)(D-2) \sum_{i \in \mathcal{I}} \mathbb{E} |X_i|^3.$$

Next, observe that

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} |X_i| |X_{ij}| |R_i + R_{ij}| = \sum_{(i,j,k) \in \mathcal{T}_2} |X_i| |X_j| |X_k| \leq \frac{1}{3} \sum_{(i,j,k) \in \mathcal{T}_2} (|X_i|^3 + |X_j|^3 + |X_k|^3),$$

where  $\mathcal{T}_2 := \{(i, j, k) ; i \sim j, i \sim \{j, k\}\}$ . Elements  $(i, j, k)$  of the set  $\mathcal{T}_2$  can be divided into the cases illustrated below, along with upper bounds on the number of pairs  $(j, k)$  with  $(i, j, k) \in \mathcal{T}_2$  for fixed  $i$ , which are also upper bounds on the number of pairs  $(i, j)$  with  $(i, j, k) \in \mathcal{T}_2$  for fixed  $j$  and upper bounds on the number of pairs  $(i, j)$  with  $(i, j, k) \in \mathcal{T}_2$  for fixed  $k$ :



Therefore,  $|\{(j, k) ; (i, j, k) \in \mathcal{T}_2\}| \leq 2D^2 - 3D + 2$  for all  $i$ ,  $|\{(i, k) ; (i, j, k) \in \mathcal{T}_2\}| \leq 2D^2 - 3D + 2$  for all  $j$  and  $|\{(i, j) ; (i, j, k) \in \mathcal{T}_2\}| \leq 2D^2 - 3D + 2$  for all  $k$ , leading to

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \mathbb{E}(|X_i| |X_{ij}| |R_i + R_{ij}|) \leq (2D^2 - 3D + 2) \sum_{i \in \mathcal{I}} \mathbb{E} |X_i|^3.$$

Writing  $\mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|R_i + R_{ij}| = \mathbb{E}(|X_i| |X_{ij}| |R'_i + R'_{ij}|)$ , where the pair  $(R'_i, R'_{ij})$  is an independent copy of the pair  $(R_i, R_{ij})$ , we can similarly estimate

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{I}_i} \mathbb{E}(|X_i| |X_{ij}|) \mathbb{E}|R_i + R_{ij}| \leq (2D^2 - 3D + 2) \sum_{i \in \mathcal{I}} \mathbb{E}|X_i|^3.$$

Collecting all together and applying Corollary 3.2.4, we obtain the following result:

**Theorem 3.3.2.** *Let  $(X_i)_{i \in \mathcal{I}}$  be a family of random variables with dependence structure which fits a graph  $\Gamma$ . Suppose that for each  $i \in \mathcal{I}$ , there are no more than  $D < \infty$  vertices  $j$  with  $i \sim j$ . Assume that  $\sum_{i \in \mathcal{I}} \mathbb{E}|X_i| < \infty$ ,  $\sum_{i \in \mathcal{I}} \mathbb{E}|X_i|^3 < \infty$  and  $\mathbb{E}X_i = 0$  for all  $i \in \mathcal{I}$ . Let  $W = \sum_{i \in \mathcal{I}} X_i$  and suppose that  $\text{var}(W) = 1$ . Then we have*

$$d_W(\mathcal{L}(W), \mathcal{N}(0, 1)) \leq (7D^2 - 12D + 8) \sum_{i \in \mathcal{I}} \mathbb{E}|X_i|^3.$$

□

**Example 3.3.3.** If the summands  $X_i$  are independent, we can set  $i \sim j$  if and only if  $i = j$ , leading to  $D = 1$  and the bound

$$d_W(\mathcal{L}(W), \mathcal{N}(0, 1)) \leq 3 \sum_{i \in \mathcal{I}} \mathbb{E}|X_i|^3,$$

which is the same as the bound in Example 3.2.5.

**Example 3.3.4.** Consider  $U$ -statistics: let  $\xi_1, \xi_2, \dots$  be independent and identically distributed random variables taking values in a measurable space  $(S, \mathcal{S})$ . Let  $F: S \times S \rightarrow \mathbb{R}$  be a symmetric product measurable function. Suppose that  $\mathbb{E}|F(\xi_1, \xi_2)|^3 < \infty$  and  $\mathbb{E}F(\xi_1, \xi_2) = 0$ . Consider the sum

$$U_n := \sum_{1 \leq i < j \leq n} F(\xi_i, \xi_j).$$

Now compute the variance:

$$\sigma_n^2 := \text{var}(U_n) = \sum_{1 \leq i < j \leq n} \sum_{1 \leq k < l \leq n} \text{cov}(F(\xi_i, \xi_j), F(\xi_k, \xi_l)).$$

Noting that

$$\text{cov}(F(\xi_i, \xi_j), F(\xi_k, \xi_l)) = \begin{cases} \text{var}(F(\xi_1, \xi_2)) & ; i = j, k = l \\ \text{cov}(F(\xi_1, \xi_2), F(\xi_1, \xi_3)) & ; |\{i, j\} \cap \{k, l\}| = 1 \\ 0 & ; \text{otherwise} \end{cases}$$

and letting  $V_1 := \text{var}(F(\xi_1, \xi_2))$  and  $V_2 := \text{cov}(F(\xi_1, \xi_2), F(\xi_1, \xi_3))$ , we obtain

$$\sigma_n^2 = \frac{n(n-1)}{2} V_1 + n(n-1)(n-2) V_2.$$

Assume that  $F$  is not almost everywhere zero (with respect to the joint distribution of  $(\xi_1, \xi_2)$ ), so that  $V_1 > 0$ . The covariance  $V_2$  can be decomposed as

$$V_2 = \mathbb{E} \left[ \text{cov}(F(\xi_1, \xi_2), F(\xi_1, \xi_3) \mid \xi_1) \right] + \text{cov} \left[ \mathbb{E}(F(\xi_1, \xi_2) \mid \xi_1), \mathbb{E}(F(\xi_1, \xi_3) \mid \xi_1) \right].$$

Since  $\xi_2$  and  $\xi_3$  are conditionally independent given  $\xi_1$ , the first term vanishes. Next,  $\mathbb{E}(F(\xi_1, \xi_2) \mid \xi_1) = \mathbb{E}(F(\xi_1, \xi_3) \mid \xi_1)$ . Therefore,

$$V_2 = \text{var} \left[ \mathbb{E}(F(\xi_1, \xi_2) \mid \xi_1) \right] \geq 0.$$

Thus, if  $\mathbb{E}(F(\xi_1, \xi_2) \mid \xi_1)$  is not almost surely constant (or equivalently not almost surely zero), we have  $V_2 > 0$ . In this case, the  $U$ -statistic is called *non-degenerate*.

The dependence structure of the family  $(F(\xi_i, \xi_j))_{1 \leq i < j \leq n}$  fits the graph on the vertex set  $\{(i, j) ; 1 \leq i < j \leq n\}$ , where vertices  $(i, j)$  and  $(k, l)$  are adjacent if  $\{i, j\} \cap \{k, l\} \neq \emptyset$ : taking families  $(i_\alpha, j_\alpha)_{\alpha \in A}$  and  $(k_\beta, l_\beta)_{\beta \in B}$ , such that  $\{i_\alpha, j_\alpha\} \cap \{k_\beta, l_\beta\} = \emptyset$  for all  $\alpha$  and  $\beta$ , the sets  $\{i_\alpha, j_\alpha ; \alpha \in A\}$  and  $\{k_\beta, l_\beta ; \beta \in B\}$  must be disjoint. Consequently, we may take  $D = 2n - 3$ .

Rescaling and applying Theorem 3.3.2, we obtain

$$d_W \left( \mathcal{L} \left( \frac{U_n}{\sigma_n} \right), \mathcal{N}(0, 1) \right) \leq \frac{n(n-1)(28n^2 - 48n + 59)}{2 \left( \frac{n(n-1)}{2} V_1 + n(n-1)(n-2) V_2 \right)^{3/2}} \mathbb{E} |F(\xi_1, \xi_2)|^3.$$

Notice that if the statistic is non-degenerate, then the rate of convergence is again  $1/\sqrt{n}$ , like for independent random variables (see Example 3.2.5).

If the statistic is degenerate, the distributions may not converge to the standard normal. A typical example is if  $F$  is of the form  $F(x, y) = G(x)G(y)$ . In this case, we have

$$U_n = \frac{1}{2} \left( \sum_{i=1}^n G(\xi_i) \right)^2 - \frac{1}{2} \sum_{i=1}^n (G(\xi_i))^2. \quad (3.3.1)$$

Observe that  $\mathbb{E}[F(\xi_1, \xi_2)] = (\mathbb{E}[G(\xi_1)])^2 = 0$ , so that  $\mathbb{E}[G(\xi_1)] = 0$  by our assumption. Letting  $\tau_1^2 := \text{var}[G(\xi_1)]$ , we find that the first term in the right hand side of (3.3.1) has expectation  $n\tau_1^2$  and variance of order  $n^2$ , while the second term has expectation  $n\tau_1^2$  and variance of order  $n$ . By the central limit theorem, the distribution of the first term divided by  $n$  approaches the chi squared distribution with one degree of freedom, scaled by a constant factor. Therefore, the distributions of  $U_n/\sqrt{n}$  converge to the centred version of that distribution.

### 3.3.2 Random permutations

Let  $\mathbf{A} = [a(i, j)]_{i, j}$  be a  $n \times n$  matrix of real numbers. Consider the statistic

$$W = \sum_{i=1}^n a(i, \Pi(i)) = \sum_{i=1}^n \sum_{j=1}^n a(i, j) \mathbf{1}(\Pi(i) = j),$$

where  $\Pi$  is a uniformly distributed random permutation of  $\{1, 2, \dots, n\}$ . Observe that  $W$  changes just for a constant if we add a constant value to all entries of a particular row or column. Subtracting the averages, we can make all rows to have sum zero. Doing the same with the columns, observe that the row sums still remain zero. Therefore, we can assume without for generality that

$$\sum_{j=1}^n a(i, j) = 0 \quad \text{for all } i \quad \text{and} \quad \sum_{i=1}^n a(i, j) = 0 \quad \text{for all } j. \quad (3.3.2)$$

In this case, of course,  $\mathbb{E}W = 0$ . Now compute

$$\text{var}(W) = \mathbb{E}(W^2) = \sum_{1 \leq i, j \leq n} \sum_{1 \leq i', j' \leq n} a(i, j) a(k, l) \mathbb{P}(\Pi(i) = j, \Pi(k) = l).$$

Noting that

$$\mathbb{P}(\Pi(i) = j, \Pi(k) = l) = \begin{cases} \frac{1}{n} & ; i = k, j = l \\ \frac{1}{n(n-1)} & ; i \neq k, j \neq l, \\ 0 & ; \text{otherwise,} \end{cases}$$

we further compute

$$\begin{aligned} \text{var}(W) &= \frac{1}{n} \sum_{1 \leq i, j \leq n} a(i, j)^2 + \frac{1}{n(n-1)} \sum_{1 \leq i, j \leq n} \sum_{1 \leq k, l \leq n; k \neq i, l \neq j} a(i, j) a(k, l) \\ &= \frac{1}{n} \sum_{1 \leq i, j \leq n} a(i, j)^2 + \frac{1}{n(n-1)} \sum_{1 \leq i, j \leq n} \sum_{1 \leq k, l \leq n} a(i, j) a(k, l) \\ &\quad - \frac{1}{n(n-1)} \sum_{1 \leq i, j \leq n} \sum_{1 \leq l \leq n} a(i, j) a(i, l) \\ &\quad - \frac{1}{n(n-1)} \sum_{1 \leq i, j \leq n} \sum_{1 \leq k \leq n} a(i, j) a(k, j) \\ &\quad + \frac{1}{n(n-1)} \sum_{1 \leq i, j \leq n} a(i, j)^2 \\ &= \frac{1}{n-1} \sum_{1 \leq i, j \leq n} a(i, j)^2. \end{aligned}$$

**In the sequel, we shall assume (3.3.2) and  $\text{var}(W) = 1$ .**

To construct decompositions from Proposition 3.1.1, we introduce the concept of simple random relocation.

**Definition 3.3.5.** Let  $A \subseteq M$  be finite sets. A *simple random relocation* of the set  $A$  within the set  $M$  is a random permutation  $T_A$  of the set  $M$ , which acts as follows:

- The elements of the set  $A$  are mapped to any elements of  $M$  uniformly at random.
- Given the latter and denoting by  $B$  the image of  $A$  under  $T_A$ , all elements of  $B \setminus A$  are mapped to the elements of  $A \setminus B$  uniformly at random.

- The other elements are left unchanged.

**Proposition 3.3.6.** *Let  $A \subseteq M$  be finite sets and let  $T_A$  be a simple random relocation of  $A$  within  $M$ . Take a uniformly distributed random permutation  $\Pi$  of  $M$ , independent of  $T_A$ . Then  $\Pi \circ T_A^{-1}$  is also uniformly distributed and is independent of the restriction of  $\Pi$  to  $A$ .*

PROOF. Let  $A = \{i_1, \dots, i_r\}$ . What we need to prove is that given  $\Pi(i_1) = j_1, \dots, \Pi(i_r) = j_r$ ,  $\Pi \circ T_A^{-1}$  is uniformly distributed. Suppose that  $T_A(i_1) = k_1, \dots, T_A(i_r) = k_r$ , and observe that the map  $\sigma \mapsto \sigma \circ T_A^{-1}$  is a one-to-one correspondence between the set of permutations  $\sigma$  with  $\sigma(i_1) = j_1, \dots, \sigma(i_r) = j_r$  and the set of permutations  $\sigma$  with  $\sigma(k_1) = j_1, \dots, \sigma(k_r) = j_r$ . Therefore, given  $T_A(i_1) = k_1, \dots, T_A(i_r) = k_r, \Pi(i_1) = j_1, \dots, \Pi(i_r) = j_r$ , the random permutation  $\sigma \circ T_A^{-1}$  is uniformly distributed over all permutations  $\sigma$  with  $\sigma(k_1) = j_1, \dots, \sigma(k_r) = j_r$ . Noting that the  $r$ -tuple  $(T_A(i_1), \dots, T_A(i_r))$  is uniformly distributed over all possible elements of  $M^r$  with distinct elements, the proof is complete.  $\square$

Now set

$$\mathcal{I}_i := \{1, 2, \dots, n\}, \quad X_i := a(i, \Pi(i)),$$

noting that  $\mathbb{E} X_i = 0$  for all  $i$ . Next, take a family of simple random relocations  $T_A$ ,  $A \subseteq \{1, 2, \dots, n\}$  of sets  $A$  within  $\{1, 2, \dots, n\}$ , which are all independent of  $\Pi$ . By Proposition 3.3.6, the random variable

$$W_i := \sum_{j=1}^N a(j, \Pi(T_{\mathcal{I}_i}^{-1}(j))) = \sum_{j=1}^N a(T_{\mathcal{I}_i}(j), \Pi(j)), \quad (3.3.3)$$

is independent of  $X_i$  for each  $i$ , so that we can set

$$\mathcal{I}_i := \{1, 2, \dots, n\}, \quad X_{ij} := a(j, \Pi(j)) - a(T_{\mathcal{I}_i}(j), \Pi(j)), \quad R_i := \sum_{j \in \mathcal{I}_i} X_{ij},$$

noting that  $R_i = W - W_i$ .

The pair  $(X_i, X_{ij})$  is uniquely determined by  $T_{\mathcal{I}_i}$  and the restriction of  $\Pi$  to  $\{i, j\}$ . Thus, take another family of simple random relocations  $T'_A$ ,  $A \subseteq \{1, 2, \dots, n\}$ , which is independent of all other random variables. Similarly as before, put

$$W_{ij} := \sum_{k=1}^N a(k, \Pi(T'_{\mathcal{I}_{i,j}}^{-1}(k))) = \sum_{j=1}^N a(T'_{\mathcal{I}_{i,j}}(k), \Pi(k)) \quad (3.3.4)$$

and by Proposition 3.3.6,  $W_{ij}$  is independent of the pair  $(X_i, X_{ij})$ .

To apply Corollary 3.2.4, we need to bound the sums

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}| |R_i + R_{ij}|], \quad \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}| |R_{ij}|], \\ \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}|] \mathbb{E}|R_i + R_{ij}|. \end{aligned} \quad (3.3.5)$$



As  $T_A$  preserves all points outside  $A \cup T_A(A) = A \cup T_A^{-1}(A)$  (and similarly  $T'_A$ ), we have

$$\begin{aligned}
|X_{ij}| &\leq \mathbf{1}\left[j \in \{i, T_{\{i\}}^{-1}(i)\}\right] \left(|a(j, \Pi(j))| + |a(T_{\{i\}}(j), \Pi(j))|\right), \\
|R_i + R_{ij}| &\leq \sum_{k=1}^n \mathbf{1}\left[k \in \{i, j, T_{\{i,j\}}^{-1}(i), T'_{\{i,j\}}^{-1}(j)\}\right] \left(|a(k, \Pi(k))| + |a(T'_{\{i,j\}}(k), \Pi(k))|\right), \\
|R_{ij}| &\leq \sum_{k=1}^n \mathbf{1}\left[k \in \{i, j, T_{\{i\}}^{-1}(i), T'_{\{i,j\}}^{-1}(i), T'_{\{i,j\}}^{-1}(j)\}\right] \\
&\quad \times \left(|a(T_{\{i\}}(k), \Pi(k))| + |a(T'_{\{i,j\}}(k), \Pi(k))|\right).
\end{aligned} \tag{3.3.6}$$

The resulting bounds on the sums in (3.3.5) can be expressed in terms of auxiliary random indices and permutations: let  $I$  be uniformly distributed over  $\{1, 2, \dots, n\}$  and independent of all other random elements. Letting

$$\begin{aligned}
J_1 &:= I, & J_2 &:= T_{\{I\}}^{-1}(I); \\
J_{r1} &:= J_r, & J_{r2} &:= T_{\{I\}}(J_r); \\
K_{r1} &:= I, & K_{r2} &:= J_r, & K_{r3} &:= T'_{\{I, J_r\}}^{-1}(I), & K_{r4} &:= T'_{\{I, J_r\}}^{-1}(J_r); \\
K_{rs1} &:= K_{rs}, & J_{rs2} &:= T_{\{I, J_r\}}(K_{rs}),
\end{aligned}$$

we find that

$$\begin{aligned}
&\sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}| |R_i + R_{ij}|] \\
&\leq n \mathbb{E}\left[|a(I, \Pi(I))| \sum_{r=1}^2 \sum_{u=1}^2 |a(J_{ru}, \Pi(J_r))| \sum_{s=1}^4 \sum_{v=1}^2 |a(K_{rsv}, \Pi(K_{rs}))|\right].
\end{aligned}$$

Applying the inequality between the arithmetic and the geometric mean, we obtain

$$\begin{aligned}
&\sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}| |R_i + R_{ij}|] \\
&\leq \frac{n}{3} \sum_{r=1}^2 \sum_{u=1}^2 \sum_{s=1}^4 \sum_{v=1}^2 \left[ \mathbb{E}|a(I, \Pi(I))|^3 + \mathbb{E}|a(J_{ru}, \Pi(J_r))|^3 + \mathbb{E}|a(K_{rsv}, \Pi(K_{rs}))|^3 \right].
\end{aligned}$$

Using independence, we find that each random index  $L$  being equal to  $I$ ,  $J_{ru}$  or  $K_{rsv}$  is uniformly distributed over  $\{1, 2, \dots, n\}$ . Moreover, since  $\Pi$  is independent of all these random indices,  $\Pi(L)$  is independent of  $L$  and also uniformly distributed over  $\{1, 2, \dots, n\}$ . This leads to the bound

$$\sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}| |R_i + R_{ij}|] \leq \frac{32}{n} \sum_{i=1}^n \sum_{j=1}^n |a(i, j)|^3.$$

Similarly, introducing  $K_{r5} := T'_{\{I\}}^{-1}(I)$ , we estimate

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}| |R_{ij}|] \\ & \leq n \mathbb{E} \left[ \left| a(I, \Pi(I)) \right| \sum_{r=1}^2 \sum_{u=1}^2 \left| a(J_{ru}, \Pi(J_r)) \right| \sum_{s=1}^5 \sum_{v=1}^2 \left| a(K_{rsv}, \Pi(K_{rs})) \right| \right] \\ & \leq \frac{40}{n} \sum_{i=1}^n \sum_{j=1}^n |a(i, j)|^3. \end{aligned}$$

Finally, to estimate the last sum in (3.3.5), we introduce another random permutation  $\Pi'$ , which is uniformly distributed and independent of all other random variables. Similarly as before, we obtain

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[|X_i| |X_{ij}|] \mathbb{E}|R_i + R_{ij}| \\ & \leq n \mathbb{E} \left[ \left| a(I, \Pi(I)) \right| \sum_{r=1}^2 \sum_{u=1}^2 \left| a(J_{ru}, \Pi(J_r)) \right| \sum_{s=1}^4 \sum_{v=1}^2 \left| a(K_{rsv}, \Pi'(K_{rs})) \right| \right] \\ & \leq \frac{32}{n} \sum_{i=1}^n \sum_{j=1}^n |a(i, j)|^3. \end{aligned}$$

Collecting all together and applying Corollary 3.2.4, we obtain the following result:

**Proposition 3.3.7.** *Let  $a(i, j)$ ,  $1 \leq i, j \leq n$ , be real numbers, such that*

$$\sum_{j=1}^n a(i, j) = 0 \quad \text{for all } i \quad \text{and} \quad \sum_{i=1}^n a(i, j) = 0 \quad \text{for all } j.$$

*Take a uniformly distributed random permutation  $\Pi$  of the set  $\{1, 2, \dots, n\}$  and let*

$$W := \sum_{i=1}^n a(i, \Pi(i)).$$

*If  $\text{var}(W) = 1$ , then we have*

$$d_W(\mathcal{L}(W), \mathcal{N}(0, 1)) \leq \frac{136}{n} \sum_{i=1}^n \sum_{j=1}^n |a(i, j)|^3.$$

□

**Remark 3.3.8.** More careful computations along with a more sophisticated version of Corollary 3.2.4 would yield a better constant. Moreover, as mentioned in Example 2.4.4, it is plausible that a sum of  $m$  independent random variables can be obtained as a limit of statistics  $W$  defined as above, where  $a(i, j) = 0$  for  $i > m$ ,  $m$  is fixed and  $n$  tends to infinity. It is possible to derive a bound in the normal approximation which approaches the bound for sums of independent random variables as stated in Example 3.3.3.

### 3.4 The Berry–Esseen theorem

The celebrated Berry–Esseen theorem bounds the error in the central limit theorem for independent and identically distributed random variables in terms of the *Kolmogorov distance*:

$$d_K(\mu, \nu) := \sum_{a \in \mathbb{R}} \left| \mu((-\infty, a]) - \nu((-\infty, a]) \right|.$$

**Theorem 3.4.1** (Berry [7], Esseen [16]). *Let  $\xi_1, \xi_2, \dots$  be independent and identically distributed random variables with  $\mathbb{E} \xi_1 = 0$ ,  $\text{var}(\xi_1) = 1$  and  $\mathbb{E} |\xi_1|^3 < \infty$ . Letting*

$$W^{(n)} := \frac{\xi_1 + \xi_2 + \dots + \xi_n}{\sqrt{n}},$$

we have

$$d_K(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)) \leq \frac{C \mathbb{E} |\xi_1|^3}{\sqrt{n}},$$

where  $C$  is a universal constant.

**Remark 3.4.2.** The calculation of the constant  $C$  has a long history. In 1941, Berry [7] claimed that the result holds with  $C = 1.88$ , but his calculations turned out not to be entirely correct (see Hsu [21]). In 1945, Esseen [16] proved the result with  $C = 7.59$ . Over the next decades, the constant was significantly improved. The best value obtained so far seem to be  $C = 0.4748$ , obtained in 2011 by Shevtsova [28] (the same author even claims  $C = 0.469$  in her paper [29], but provides no proof). Shevtsova's bound is not far from optimal: a lower bound  $\frac{\sqrt{10+3}}{6\sqrt{2\pi}} > 0.4097$  on the constant  $C$  was derived in 1956 by Esseen [17]. For more details on the history of calculation of  $C$ , see Korolev and Shevtsova [22].

**Remark 3.4.3.** Example 3.3.3 provides a bound in the Wasserstein distance:

$$d_W(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)) \leq \frac{3 \mathbb{E} |\xi_1|^3}{\sqrt{n}},$$

which, combined with (A.2.7), gives

$$d_W(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)) \leq \frac{\sqrt{6}}{\sqrt[4]{2\pi}} \frac{\sqrt{\mathbb{E} |\xi_1|^3}}{\sqrt[4]{n}}.$$

However, the latter bound does not preserve the rate of convergence.

Here, we prove Theorem 3.4.1 by Stein's method. Unfortunately, the constant will be far from optimal. However, the main advantage of Stein's method is that it can be readily applied to sums of dependent random variables, where most other methods do not seem work. In particular, most of the improvements of the constant in the Berry–Esseen theorem (including the improvement by Shevtsova [28]) are proved by the method of characteristic functions. For sums of independent random variables, characteristic functions satisfy the multiplication formula, which has no straightforward extension to, let's say, local dependence.

Thus, the proof of the Berry–Esseen theorem can serve as a guideline how to extend the result to sums of dependent random variables. Here, we shall not do the latter. In fact, bounds of correct order in terms of the Kolmogorov metric are much harder to obtain than for the Wasserstein metric. In particular, it seems to be very difficult to derive a result being as general as Corollary 3.2.4. Instead, more special results have been derived. Bolthausen [8] proves a Berry–Esseen type result for random permutations. Chen and Shao [12] prove a result for local dependence, but that result does not yield the counterpart of Theorem 3.3.2 for local dependence. One can do more under the assumption of higher moments: see Chen, Goldstein and Röllin [10] and references therein. The assumption of boundedness admits even more general and cleaner results: see Dembo and Rinott [15], Goldstein [19] and Raič [26].

In order to derive the Berry–Esseen theorem by Stein's method, it is beneficial to consider functions with bounded total variation in view of Section A.3.

**Theorem 3.4.4.** *Consider a function  $h: \mathbb{R} \rightarrow \mathbb{R}$ .*

- (1) *If  $h$  has bounded variation, then the function  $f$  defined by (3.2.4) is an almost-everywhere solution to the Stein equation (3.2.1), which satisfies  $V(f') \leq 2V(h)$ .*
- (2) *If  $h$  is absolutely continuous and  $h'$  has bounded variation, the function  $f$  defined by (3.2.4) is a classical solution to the Stein equation (3.2.1). Moreover,  $f'$  is absolutely continuous and  $V(f'') \leq 2V(h')$ .*

PROOF.

*Part (1).* First recall that  $h$  is bounded by Remark B.2.2 and measurable by Corollary B.2.5. Therefore,  $\mathcal{N}|h| < \infty$ . By Proposition 3.2.1, there exists an almost-everywhere solution  $f$  to the Stein equation, which means that  $f$  is absolutely continuous. Recalling (3.2.6), there is an almost-everywhere derivative of  $f$  given by

$$f'(w) = h(w) - \psi'(-w) \int_{-\infty}^w h(x) \phi(x) dx - \psi'(w) \int_w^{\infty} h(x) \phi(x) dx. \quad (3.4.1)$$

We shall derive a formula similar to (3.2.7), which was derived from (3.2.6) by integration by parts, integrating  $\phi$  and differentiating  $h$ . This can be done if  $h$  is absolutely continuous, which is not an assumption in our case. However, we can use the (improper) *Riemann–Stieltjes integral*. By assumption,  $h$  has bounded variation. Letting  $\Phi^-(y) := \Phi(-y)$ ,  $\Phi$  and  $\Phi^-$  are absolutely continuous. By Proposition B.3.8, we can then rewrite (3.4.1) as

$$f'(w) = h(w) - \psi'(-w) \int_{-\infty}^w h(x) d\Phi(x) + \psi'(w) \int_w^{\infty} h(x) d\Phi^-(x).$$

By the integration by parts formula (Proposition B.3.6) and noting that

$\lim_{w \rightarrow -\infty} h(w) \Phi(w) = \lim_{w \rightarrow \infty} h(w) \Phi(-w) = 0$ , we have

$$\begin{aligned} f'(w) &= h(w) - \psi'(-w) h(w) \Phi(w) + \psi'(-w) \int_{-\infty}^w \Phi(y) dh(y) \\ &\quad - \psi'(w) h(w) \Phi(-w) - \psi'(w) \int_w^{\infty} \Phi(-y) dh(y) \\ &= \psi'(-w) \int_{-\infty}^w \Phi(y) dh(y) - \psi'(w) \int_w^{\infty} \Phi(-y) dh(y), \end{aligned}$$

where the last equality is due to Proposition C.2.3 for  $r = 1$ .

Let  $\Lambda_h$  be the signed measure associated to  $h$  in view of Definition B.4.9. By part (2) of Proposition B.4.11, we have  $\int_a^w \Phi(y) dh(y) = \int_{(a,w]} \Phi d\Lambda_h$  and

$\int_w^b \Phi(-y) dh(y) = \int_{(w,b]} \Phi^- d\Lambda_h$  for all  $a \leq w \leq b$ . Since  $\Phi$  is bounded on  $(-\infty, w]$  and  $\Phi^-$  is bounded on  $(w, \infty)$ , we may take the limit in  $a$  and  $b$ , leading to  $\int_{-\infty}^w \Phi(y) dh(y) = \int_{(-\infty,w]} \Phi d\Lambda_h$  and  $\int_w^{\infty} \Phi(-y) dh(y) = \int_{(w,\infty)} \Phi^- d\Lambda_h$ . We may rewrite this as

$$f'(w) = \int_{-\infty}^{\infty} F(w, y) \Lambda_h(dy),$$

where

$$F(w, y) = \begin{cases} -\psi'(w) \Phi(-y) & ; w < y \\ \psi'(-w) \Phi(y) & ; w \geq y. \end{cases}$$

Letting  $F_y(w) := F(w, y)$ , it is left to the reader as an exercise to show that  $\sum_{k=1}^n |F_y(w_k) - F_y(w_{k-1})| \leq 2(\psi'(y) \Phi(-y) + \psi'(-y) \Phi(y))$  for every sequence  $w_0 \leq w_1 \leq \dots \leq w_n$  (use Proposition C.1.5). Therefore,  $V(F_y) \leq 2(\psi'(y) \Phi(-y) + \psi'(-y) \Phi(y))$ . Moreover, by Proposition C.1.6, we have  $\lim_{t \rightarrow \infty} [ |F_y(-t) - F_y(w - \frac{1}{t})| + |F_y(w - \frac{1}{t}) - F_t(y)| + |F_y(y) - F_y(t)| ] = 2(\psi'(y) \Phi(-y) + \psi'(-y) \Phi(y))$ . By Proposition C.2.3, we have  $\psi'(y) \Phi(-y) + \psi'(-y) \Phi(y) = 1$ , so that  $V(F_y) = 2$  for all  $y$ . Finally, applying part (2) of Proposition B.2.7 along with part (1) of Proposition B.4.11, we estimate  $V(f') \leq \int_{-\infty}^{\infty} V(F_y) |\Lambda_h|(dy) = 2 \|\Lambda_h\| = 2V(h)$ , proving part (1).

*Part (2).* We proceed similarly as in part (1). First recall that  $h'$  is bounded by Remark B.2.2. Therefore,  $h$  is of linear growth and  $\mathcal{N}|h| < \infty$ . By Proposition 3.2.1, there exists a classical solution  $f$  to the Stein equation, such that  $f'$  is absolutely continuous and, recalling (3.2.8),

$$f''(w) = h'(w) - \psi''(-w) \int_{-\infty}^w h'(y) \Phi(y) dy - \psi''(w) \int_w^{\infty} h'(y) \Phi(-y) dy. \quad (3.4.2)$$

Again, we rewrite this formula in terms of the Riemann–Stieltjes integral. By assumption,  $h'$  has finite total variation. The functions  $\Phi_2(y) := \int_{-\infty}^y \Phi(t) dt$   $\Phi_2^-(y) := \Phi_2(-y)$  are absolutely continuous. By Proposition B.3.8, we have

$$f''(w) = h'(w) - \psi''(-w) \int_{-\infty}^w h'(y) d\Phi_2(y) + \psi''(w) \int_w^{\infty} h'(y) d\Phi_2^-(y).$$

By the integration by parts formula (Proposition B.3.6) and noting that  $\lim_{w \rightarrow -\infty} h'(w) \Phi_2(w) = \lim_{w \rightarrow \infty} h'(w) \Phi_2(-w) = 0$  by Corollary C.2.2, we have

$$\begin{aligned} f''(w) &= h'(w) - \psi''(-w) h'(w) \Phi_2(w) + \psi''(-w) \int_{-\infty}^w \Phi_2(y) dh'(y) \\ &\quad - \psi''(w) h'(w) \Phi_2(-w) - \psi''(w) \int_w^{\infty} \Phi_2(-y) dh'(y) \\ &= \psi''(-w) \int_{-\infty}^w \Phi_2(y) dh'(y) - \psi''(w) \int_w^{\infty} \Phi_2(-y) dh'(y), \end{aligned}$$

where the last equality is due to Proposition C.2.3 for  $r = 2$ . Similarly as in the proof of part (1), we can rewrite this as

$$f''(w) = \int_{-\infty}^{\infty} G(w, y) \Lambda_{h'}(dy),$$

where

$$G(w, y) = \begin{cases} -\psi''(w) \Phi_2(-y) & ; w < y \\ \psi''(-w) \Phi_2(y) & ; w \geq y. \end{cases}$$

Letting  $G_y(w) := G(w, y)$ , we again find that  $V(G_y) = 2(\psi''(y) \Phi_2(-y) + \psi''(-y) \Phi_2(y)) = 2$ , with the last equality due to Proposition C.2.3 for  $r = 2$ . Finally, applying part (2) of Proposition B.2.7 along with part (1) of Proposition B.4.11, we estimate  $V(f'') \leq \int_{-\infty}^{\infty} V(G_y) |\Lambda_{h'}|(dy) = 2 \|\Lambda_{h'}\| = 2V(h')$ , completing the proof.  $\square$

**PROOF OF THEOREM 3.4.1.** Sums of independent random variables are clearly a very special case of decompositions of Barbour, Karoński and Ruciński introduced in Section 3.1. Letting

$$\begin{aligned} \mathcal{I} &:= \{1, 2, \dots, n\}, & \mathcal{I}_i &:= \{0\}, \\ X_i^{(n)} &:= R_i^{(n)} := X_{i0}^{(n)} := \frac{\xi_i}{\sqrt{n}}, & W_i^{(0)} &:= W_{i0}^{(n)} := W^{(n)} - X_i^{(n)}, & R_{i0}^{(n)} &:= 0, \end{aligned}$$

all conditions specified in Section 3.1 are fulfilled and (3.1.5) reduces to

$$\begin{aligned} \mathbb{E}[f'(W) - f(W)W] &= \sum_{i=1}^n \mathbb{E} \left[ f''(W_i^{(n)} + \theta_2 X_i^{(n)}) X_i^{(n)} \mathbb{E}((X_i^{(n)})^2) \right. \\ &\quad \left. - \theta_1 f''(W_i^{(n)} + \theta_1 \theta_2 X_i^{(n)}) (X_i^{(n)})^3 \right], \end{aligned} \tag{3.4.3}$$

where  $\theta_1$  and  $\theta_2$  are uniformly distributed over  $[0, 1]$  and independent of each other as well as of  $X_1^{(n)}, \dots, X_n^{(n)}$ . This is true for all functions  $f$  with  $M_2(f) < \infty$  (see Proposition 3.1.1).

Denoting the indicators of half-lines by

$$h_a(w) := \begin{cases} 1 & ; w \leq a \\ 0 & ; x > a, \end{cases}$$

the error in the normal approximation in the Kolmogorov metric is expressed as

$$d_K\left(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)\right) = \sup_{a \in \mathbb{R}} \left| \mathbb{E}[h_a(W^{(n)})] - \mathcal{N}h_a \right|.$$

Unfortunately, the derivative of the solution of the Stein equation (3.2.1) with  $h_a$  in place of  $h$  is not absolutely continuous (exercise). Therefore, similarly as in the proof of Propositions A.2.9 and A.4.11, we approximate the *step* functions  $h_a$  by the *slope* functions

$$h_{a,b}(x) := \begin{cases} 1 & ; x \leq a \\ \frac{b-x}{b-a} & ; a \leq x \leq b \\ 0 & ; x \geq b, \end{cases}$$

defined for  $a < b$ . Observe that

$$\begin{aligned} \mathcal{N}h_{a,b} - \mathcal{N}h_a &= \int_a^b \frac{b-x}{b-a} \phi(x) dx \leq \frac{1}{\sqrt{2\pi}} \int_a^b \frac{b-x}{b-a} dx = \frac{b-a}{2\sqrt{2\pi}}, \\ \mathcal{N}h_b - \mathcal{N}h_{a,b} &= \int_a^b \frac{x-a}{b-a} \phi(x) dx \leq \frac{1}{\sqrt{2\pi}} \int_a^b \frac{x-a}{b-a} dx = \frac{b-a}{2\sqrt{2\pi}}. \end{aligned}$$

Therefore, for any  $a \in \mathbb{R}$  and  $\varepsilon > 0$ ,

$$\begin{aligned} \mathbb{E}[h_a(W^{(n)})] - \mathcal{N}h_a &\leq \mathbb{E}[h_{a,a+\varepsilon}(W^{(n)})] - \mathcal{N}h_{a,a+\varepsilon} + \frac{\varepsilon}{2\sqrt{2\pi}}, \\ \mathcal{N}h_a - \mathbb{E}[h_a(W^{(n)})] &\leq \mathcal{N}h_{a-\varepsilon,a} - \mathbb{E}[h_{a-\varepsilon,a}(W^{(n)})] + \frac{\varepsilon}{2\sqrt{2\pi}} \end{aligned}$$

and consequently

$$d_K\left(\mathcal{L}(W^{(n)}), \mathcal{N}(0, 1)\right) \leq \sup_{a \in \mathbb{R}} \left| \mathbb{E}[h_{a,a+\varepsilon}(W^{(n)})] - \mathcal{N}h_{a,a+\varepsilon} \right| + \frac{\varepsilon}{2\sqrt{2\pi}}. \quad (3.4.4)$$

Let  $f_{a,b}$  be the solution to the Stein equation (3.2.1) defined by (3.2.4) with  $f_{a,b}$  in place of  $f$  and  $h_{a,b}$  in place of  $h$ . Since  $M_1(h_{a,b}) < \infty$ , we have  $M_1(f_{a,b}) < \infty$  by Theorem 3.2.3, so that (3.4.3) applies with  $f_{a,b}$  in place of  $f$ . Therefore,

$$\begin{aligned} \mathbb{E}[h_{a,a+\varepsilon}(W^{(n)})] - \mathcal{N}h_{a,a+\varepsilon} &= \sum_{i=1}^n \mathbb{E} \left[ f''_{a,a+\varepsilon}(W_i^{(n)} + \theta_2 X_i^{(n)}) X_i^{(n)} \mathbb{E}((X_i^{(n)})^2) \right. \\ &\quad \left. - \theta_1 f''_{a,a+\varepsilon}(W_i^{(n)} + \theta_1 \theta_2 X_i^{(n)}) (X_i^{(n)})^3 \right]. \end{aligned} \quad (3.4.5)$$

Now fix  $x \in \mathbb{R}$  and consider the expectation

$$\mathbb{E} \left[ f''_{a,a+\varepsilon}(W_i^{(n)} + x) \right] = \mathbb{E} \left[ f''_{a,a+\varepsilon} \left( \sqrt{\frac{n-1}{n}} W^{(n-1)} + x \right) \right]. \quad (3.4.6)$$

Denoting by  $\mathcal{N}(\mu, \sigma^2)$  the normal distribution with mean  $\mu$  and variance  $\sigma^2$ , write

$$\left\langle f''_{a,a+\varepsilon}, \mathcal{N} \left( x, \frac{n-1}{n} \right) \right\rangle = \sqrt{\frac{n}{n-1}} \int_{-\infty}^{\infty} f''_{a,a+\varepsilon}(w) \phi \left( (w-x) \sqrt{\frac{n}{n-1}} \right) dw.$$

However, since  $\int_{-\infty}^{\infty} f''_{a,a+\varepsilon}(w) dw = \lim_{w \rightarrow \infty} f'_{a,a+\varepsilon}(w) - \lim_{w \rightarrow -\infty} f'_{a,a+\varepsilon}(w) = 0$ , we also have

$$\left\langle f''_{a,a+\varepsilon}, \mathcal{N}\left(x, \frac{n-1}{n}\right) \right\rangle = \sqrt{\frac{n}{n-1}} \int_{-\infty}^{\infty} f''_{a,a+\varepsilon}(w) \left[ \phi\left((w-x)\sqrt{\frac{n}{n-1}}\right) - \frac{1}{2\sqrt{2\pi}} \right] dw.$$

Applying by Proposition B.2.6 and part (1) of Theorem 3.4.4, we estimate

$$\begin{aligned} \left| \left\langle f''_{a,a+\varepsilon}, \mathcal{N}\left(x, \frac{n-1}{n}\right) \right\rangle \right| &\leq \frac{1}{2\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} \int_{-\infty}^{\infty} |f''_{a,a+\varepsilon}(w)| dw \\ &= \frac{1}{2\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} V(f'_{a,a+\varepsilon}) \\ &\leq \frac{1}{\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} V(h_{a,a+\varepsilon}) \\ &= \frac{\sqrt{n}}{\sqrt{2\pi(n-1)}}. \end{aligned} \tag{3.4.7}$$

Next, observe that, by Corollary A.3.3 and part (2) of Theorem 3.4.4,

$$\begin{aligned} &\left| \mathbb{E} \left[ f''_{a,a+\varepsilon} \left( \sqrt{\frac{n-1}{n}} W^{(n-1)} + x \right) \right] - \left\langle f''_{a,a+\varepsilon}, \mathcal{N}\left(x, \frac{n-1}{n}\right) \right\rangle \right| \\ &\leq V(f''_{a,a+\varepsilon}) d_K \left( \mathcal{L} \left( \sqrt{\frac{n-1}{n}} W^{(n-1)} + x \right), \mathcal{N}\left(x, \frac{n-1}{n}\right) \right) \\ &= V(f''_{a,a+\varepsilon}) d_K \left( \mathcal{L}(W^{(n-1)}), \mathcal{N}(0, 1) \right) \\ &\leq 2 V(h'_{a,a+\varepsilon}) d_K \left( \mathcal{L}(W^{(n-1)}), \mathcal{N}(0, 1) \right) \\ &= \frac{4}{\varepsilon} d_K \left( \mathcal{L}(W^{(n-1)}), \mathcal{N}(0, 1) \right). \end{aligned} \tag{3.4.8}$$

Combining (3.4.6), (3.4.7) and (3.4.8), we obtain

$$\left| \mathbb{E} \left[ f''_{a,a+\varepsilon}(W_i^{(n)} + x) \right] \right| \leq \frac{\sqrt{n}}{\sqrt{2\pi(n-1)}} + \frac{4\delta_{n-1}}{\varepsilon},$$

where

$$\delta_n := d_K \left( \mathcal{L}(W^{(n)}), \mathcal{N}(0, 1) \right).$$

We also have

$$\begin{aligned} \left| \mathbb{E} \left[ f''_{a,a+\varepsilon}(W_i^{(n)} + \theta_2 X_i^{(n)}) \mid X_i^{(n)}, \theta_1, \theta_2 \right] \right| &\leq 2\sqrt{\frac{n}{n-1}} + \frac{4\delta_{n-1}}{\varepsilon}, \\ \left| \mathbb{E} \left[ f''_{a,a+\varepsilon}(W_i^{(n)} + \theta_1 \theta_2 X_i^{(n)}) \mid X_i^{(n)}, \theta_1, \theta_2 \right] \right| &\leq 2\sqrt{\frac{n}{n-1}} + \frac{4\delta_{n-1}}{\varepsilon}. \end{aligned}$$



Plugging into (3.4.5), we obtain

$$\begin{aligned} & \left| \mathbb{E}[h_{a,a+\varepsilon}(W^{(n)})] - \mathcal{N}h_{a,a+\varepsilon} \right| \\ & \leq \left( \frac{\sqrt{n}}{\sqrt{2\pi(n-1)}} + \frac{4\delta_{n-1}}{\varepsilon} \right) \sum_{i=1}^n \left( \mathbb{E}|X_i^{(n)}| \mathbb{E}(X_i^{(n)})^2 + \mathbb{E}(\theta_1|X_i^{(n)}|^3) \right). \end{aligned}$$

By Jensen's inequality, we have  $\mathbb{E}|X_i^{(n)}| \leq \left( \mathbb{E}|X_i^{(n)}|^3 \right)^{1/3}$  and  $\mathbb{E}(X_i^{(n)})^2 \leq \left( \mathbb{E}|X_i^{(n)}|^3 \right)^{2/3}$ , leading to

$$\begin{aligned} \left| \mathbb{E}[h_{a,a+\varepsilon}(W^{(n)})] - \mathcal{N}h_{a,a+\varepsilon} \right| & \leq \left( \frac{3}{2} \frac{\sqrt{n}}{\sqrt{2\pi(n-1)}} + \frac{6\delta_{n-1}}{\varepsilon} \right) n \mathbb{E}|X_1^{(n)}|^3 \\ & = \left( \frac{3\sqrt{n}}{\sqrt{8\pi(n-1)}} + \frac{6\delta_{n-1}}{\varepsilon} \right) \frac{\mathbb{E}|\xi_1|^3}{\sqrt{n}}. \end{aligned}$$

Combining with (3.4.4), we obtain

$$\delta_n \leq \left( \frac{3\sqrt{n}}{\sqrt{8\pi(n-1)}} + \frac{6\delta_{n-1}}{\varepsilon} \right) \frac{\mathbb{E}|\xi_1|^3}{\sqrt{n}} + \frac{\varepsilon}{2\sqrt{2\pi}}.$$

It is easy to check that  $\inf_{\varepsilon>0} \left( \frac{a}{\varepsilon} + b\varepsilon \right) = 2\sqrt{ab}$  for all  $a, b \geq 0$ . Therefore,

$$\delta_n \leq \frac{3\mathbb{E}|\xi_1|^3}{\sqrt{8\pi(n-1)}} + \sqrt{\frac{12\mathbb{E}|\xi_1|^3}{\sqrt{2\pi}} \frac{\delta_{n-1}}{\sqrt{n}}}.$$

Letting  $C_n := \delta_n \sqrt{n} / \mathbb{E}|\xi_1|^3$ , we rewrite this as

$$C_n \leq \frac{3\sqrt{n}}{\sqrt{8\pi(n-1)}} + \sqrt{\frac{12}{\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} C_{n-1}}.$$

Trivially,  $\delta_n \leq 1$  for all  $n$ . Next, by Jensen's inequality, we have  $\mathbb{E}|\xi_1|^3 \geq (\mathbb{E}\xi_1^2)^{3/2} = 1$ . Therefore,  $C_n \leq \sqrt{n}$  for all  $n$ . However, our goal is to bound  $C_n$  uniformly in  $n$ . Numerical calculations show that

$$\sqrt{n} < \frac{3\sqrt{n}}{\sqrt{8\pi(n-1)}} + \sqrt{\frac{12}{\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} \sqrt{n-1}}$$

for all  $n \leq 35$ . Therefore, in this case, we can merely say that  $C_n \leq \sqrt{n}$ . Letting

$$C_{35}^* := \sqrt{35}, \quad C_n^* := \frac{3\sqrt{n}}{\sqrt{8\pi(n-1)}} + \sqrt{\frac{12}{\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} C_{n-1}^*},$$

numerical calculations show

$n$	$C_n^*$
35	5.916080
36	5.966365
37	5.987821
38	5.996282
39	5.998944
40	5.999055
41	5.998073

Clearly,  $C_n \leq C_{35}^*$  for all  $n \leq 35$  and  $C_n \leq C_n^*$  for all  $n \geq 35$ . We have  $C_{35}^* < C_{36}^* < C_{37}^* < C_{38}^* < C_{39}^* < C_{40}^*$ , so that  $C_n \leq C_{40}^*$  for all  $n \leq 40$ . Now prove by induction that this remains true for  $n \geq 40$ . Indeed, for  $n \geq 41$ , observe that

$$\begin{aligned} C_n &\leq \frac{3\sqrt{n}}{\sqrt{8\pi(n-1)}} + \sqrt{\frac{12}{\sqrt{2\pi}} \sqrt{\frac{n}{n-1}} C_{n-1}} \\ &\leq \frac{3\sqrt{41}}{\sqrt{320\pi}} + \sqrt{\frac{12}{\sqrt{2\pi}} \sqrt{\frac{41}{40}} C_{40}^*} \\ &< 5.998074 < C_{40}^*. \end{aligned}$$

Thus, we have proved the result for  $C = C_{40}^* < 6$ . □

**Remark 3.4.5.** For large  $n$ , we can do slightly better: one can show that  $\limsup_{n \rightarrow \infty} C_n \leq C_*$ , where  $C_*$  is the unique solution to the equation

$$C^* = \frac{3}{\sqrt{8\pi}} + \sqrt{\frac{12}{\sqrt{2\pi}} C^*}.$$

The solution can be computed explicitly as  $C^* = \frac{15 + \sqrt{216}}{\sqrt{2\pi}} < 5.923683$ . Thus, the improvement is still very small comparable to the result of Shevtsova [28], which gives  $C = 0.4748$ .

As regards Stein's method, Chen and Shao [11] succeed to derive the Berry–Esseen inequality with  $C = 4.1$  (however, their constant also applies for sums of non-identically distributed random variables). Even this constant is far away from 0.4748. As already mentioned, the advantage of Stein's method is possibility of extension to dependent summands, but we shall not tackle this issue here.

# Appendix A

## Convergence of probability measures

**Note.** This appendix requires a basic knowledge of the theory of metric spaces and topology. For basic definitions related to the latter as well for a slightly deeper insight, the reader is referred to [24] or [30].

We are often interested whether a sequence of probability distributions  $(\mu_n)_{n \in \mathbb{N}}$  on a measurable space  $(S, \mathcal{S})$  converges to a given probability measure  $\mu$ . In order to make this precise, we have to endow the space  $\text{Pr}(S, \mathcal{S})$ , the set of all probability measures on  $(S, \mathcal{S})$ , with a topology. The latter will be based on *test functions*: the probability measures  $\mu$  and  $\nu$  are “close” if the integrals  $\langle f, \mu \rangle$  and  $\langle f, \nu \rangle$  are close for a suitable class of measurable functions  $f$ , where we recall from (1.1.1):

$$\langle f, \mu \rangle := \int f \, d\mu. \tag{A.0.1}$$

**Remark A.0.6.** If the test functions are unbounded, they cannot test all probability measures. In this case,  $\text{Pr}(S, \mathcal{S})$  should be replaced by a suitable subspace.

Throughout this appendix,  $(S, \mathcal{S})$  will denote a measurable space, while  $\mathcal{M}$  will denote a subspace of  $\text{Pr}(S, \mathcal{S})$ .

### A.1 Metrics based on test functions

One of the ways to construct a topology on  $\mathcal{M}$  from a class of test functions  $\mathcal{F}$  is to define the following metric:

$$d_{\mathcal{F}}(\mu, \nu) = \sup_{f \in \mathcal{F}} |\langle f, \nu \rangle - \langle f, \mu \rangle|. \tag{A.1.1}$$

Of course, we must check whether the right hand side is well defined and whether it represents a metric. Firstly,  $\langle |f|, \mu \rangle$  must be finite for all  $\mu \in \mathcal{M}$  and  $f \in \mathcal{F}$ . However, this does not imply that  $d(\mu, \nu) < \infty$ . Suppose that the latter is true. In this case, the symmetry is obvious and the triangle inequality is easy to check, too. It remains to check

that  $d(\mu, \nu) = 0$  implies  $\mu = \nu$ . This is obviously true if the class  $\mathcal{F}$  is rich enough. The proof of the following assertion is easy and is therefore left to the reader:

**Proposition A.1.1.** *Let  $d_{\mathcal{F}}$  be as in (A.1.1) and let  $d(\mu, \nu) < \infty$  for any  $\mu$  and  $\nu$ . Then  $d_{\mathcal{F}}$  is a metric if and only if the test functions from  $\mathcal{F}$  separate the probability measures from  $\mathcal{M}$ , that is, when for any two probability measures  $\mu \neq \nu \in \mathcal{M}$ , there exists a function  $f \in \mathcal{F}$ , such that  $\langle f, \mu \rangle \neq \langle f, \nu \rangle$ .  $\square$*

**Example A.1.2.** Letting  $\mathcal{F}$  be the class of indicators of all measurable sets, we obtain the *total variation metric*, which will be denoted by  $d_{\text{TV}}$ . Thus, for example, Proposition 2.2.2 can be rewritten as

$$d_{\text{TV}}(\mathcal{L}(W), \text{Po}(\lambda)) \leq \frac{1 - e^{-\lambda}}{\lambda} \sum_{i \in \mathcal{I}} p_i^2.$$

where  $\mathcal{L}(W)$  denotes the distribution (law) of  $W$ , that is,  $\mathcal{L}(W)(A) = \mathbb{P}(W \in A)$ . Clearly, we can take  $\mathcal{M} = \text{Pr}(S, \mathcal{S})$  and the class of indicators of all measurable sets separates all probability measures *by definition*, so that  $d_{\text{TV}}$  is a metric.

**Remark A.1.3.** The total variation metric typically works well in discrete spaces. Otherwise, it is usually too strong. As an example, consider a sequence  $x_1, x_2, \dots$  of points in a metric space, which converges to  $x$ . If all the points  $x_n$  are different from  $x$ , then  $d_{\text{TV}}(\delta_{x_n}, \delta_x) = 1$  for all  $n$ , so that the sequence of Dirac measures  $\delta_{x_1}, \delta_{x_2}, \dots$  does not converge to  $\delta_x$ .

**Example A.1.4.** Letting  $\mathcal{F}$  be the class of indicators of all half-lines  $(-\infty, a]$  on the real line endowed with the Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R})$ , we obtain the *Kolmogorov metric*, which will be denoted by  $d_{\text{K}}$ . Again, we can take  $\mathcal{M} = \text{Pr}(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . We claim that  $\mathcal{F}$  separates all probability measures. To show this, take Borel probability measures  $\mu$  and  $\nu$  which agree on  $\mathcal{F}$ , that is, on all all half-lines  $(-\infty, a]$ . The family of all half-lines is closed under intersections and generates the Borel  $\sigma$ -algebra. By Theorem A.1.6 below,  $\mu$  and  $\nu$  agree on all Borel sets, so that  $\mu = \nu$ . Therefore,  $\mathcal{F}$  separates all probability measures, so that  $d_{\text{K}}$  is also a metric. Clearly,  $d_{\text{K}} \leq d_{\text{TV}}$ , so that  $d_{\text{K}}$  is weaker than  $d_{\text{TV}}$ .

**Definition A.1.5.** A  $\pi$ -*system* is a collection of sets which is closed under finite intersections.

The following result is a consequence of Dynkin's  $\pi$ - $\lambda$  theorem (see, for example, Corollary 1.6.3 of Cohn [13]).

**Theorem A.1.6.** *If two probability measures agree on a  $\pi$ -system which generates a  $\sigma$ -algebra  $\mathcal{S}$ , then they agree on  $\mathcal{S}$ .  $\square$*

The fact that the set of all half-lines  $(-\infty, a]$  separates probability measures can be generalized in the following way:

**Proposition A.1.7.** *Let  $\mathcal{F}$  be a class of measurable functions on  $(S, \mathcal{S})$  and let  $\mathcal{P}$  be a  $\pi$ -system which generates  $\mathcal{S}$ . For each  $A \in \mathcal{P}$ , suppose that there exists a uniformly bounded sequence of functions  $f_n \in \text{span } \mathcal{F}$  which converges pointwise to the indicator  $\mathbf{1}_A$  of a set  $A$ . Then  $\mathcal{F}$  separates all probability measures from  $\text{Pr}(S, \mathcal{S})$ .*

PROOF. Let  $\mu, \nu \in \Pr(S, \mathcal{S})$  be such that  $\langle f, \mu \rangle = \langle f, \nu \rangle$  for all  $f \in \mathcal{F}$ . Thanks to linearity, this is also true for all  $f \in \text{span } \mathcal{F}$ . From the dominated convergence theorem, it immediately follows that  $\mu(A) = \nu(A)$  for all  $A \in \mathcal{P}$ . By Theorem A.1.6,  $\mu(A) = \nu(A)$  for all  $A \in \mathcal{S}$ , so that  $\mu = \nu$ .  $\square$

**Remark A.1.8.** Again, if  $x_1, x_2, \dots$  is a sequence of real numbers converging to  $x$  and all numbers  $x_n$  are different from  $x$ , we still have  $d_K(\delta_{x_n}, \delta_x) = 1$  for all  $n$ , so that the sequence of Dirac measures  $\delta_{x_1}, \delta_{x_2}, \dots$  does not converge to  $\delta_x$  in the Kolmogorov metric either.

Now consider another example: take  $0 < p < 1$  and let  $X_n \sim \text{Bin}(n, p)$ ,  $n \in \mathbb{N}$  be binomial random variables. Define

$$Y_n := \text{var}(X_n)^{-1/2}(X_n - \mathbb{E} X_n) = \frac{X_n - np}{\sqrt{np(1-p)}}. \quad (\text{A.1.2})$$

Then the sequence  $Y_n$ ,  $n \in \mathbb{N}$ , does not converge to the standard normal distribution  $\mathcal{N}(0, 1)$  in the total variation metric: letting  $A := \{k/\sqrt{np(1-p)}; k \in \mathbb{Z}, n \in \mathbb{N}\}$ , we have  $\mathbb{P}(Y_n \in A) = 1$  for all  $n \in \mathbb{N}$ , while  $\mathcal{N}(0, 1)\{A\} = 0$ . However, by the classical Laplace central limit theorem, it converges in  $d_K$ .

## A.2 The Wasserstein metric

As seen in the previous section, convergence of a sequence of points does not imply convergence of the underlying sequence of the Dirac measures in all metrics. Here, we define a metric for which this is true. However, this metric will not be defined on all probability measures.

**Definition A.2.1.** Let  $S$  be endowed with a metric  $d$  and let  $\mathcal{S}$  be the underlying Borel  $\sigma$ -algebra. A probability measure  $\mu \in \Pr(S, \mathcal{S})$  has *finite first absolute moment* with respect to the underlying metric  $d$  if

$$\int d(x, y) \mu(dy) < \infty \quad (\text{A.2.1})$$

for some (all)  $x \in S$ . The space of Borel probability measures with finite first absolute moment with respect to a metric  $d$  will be denoted by  $\Pr^{L^1}(S, d)$ .

**Remark A.2.2.** For each probability measure  $\mu \in \Pr^{L^1}(S, d)$  and each Lipschitz function  $f: S \rightarrow \mathbb{R}$ , we have  $\langle |f|, \mu \rangle < \infty$ .

**Definition A.2.3.** The *Wasserstein metric* (also known as Dudley, Fortet–Mourier or Kantorovich metric) on  $\Pr^{L^1}(S, d)$  is defined by

$$d_W(\mu, \nu) := \sup_{M_1(f) \leq 1} |\langle f, \nu \rangle - \langle f, \mu \rangle|, \quad (\text{A.2.2})$$

where

$$M_1(f) := \sup_{x \neq y} \frac{|f(x) - f(y)|}{d(x, y)} \quad (\text{A.2.3})$$

(notice that  $M_1(f)$  is defined differently in (B.1.2), but the definitions coincide by Proposition B.1.9).

**Proposition A.2.4.** *The right hand side of (A.2.2) is finite for any two probability measures  $\mu, \nu \in \text{Pr}^{L^1}(S, d)$ .*

PROOF. First, we check finiteness. Taking  $f$  with  $M_1(f) \leq 1$  and choosing arbitrary  $z \in S$ , observe that

$$\begin{aligned} |\langle f, \nu \rangle - \langle f, \mu \rangle| &= \left| \int_S \int_S (f(x) - f(y)) \mu(dx) \nu(dy) \right| \\ &\leq \int_S \int_S d(x, y) \mu(dx) \nu(dy) \\ &\leq \int_S \int_S (d(x, z) + d(y, z)) \mu(dx) \nu(dy) \\ &= \int_S d(x, z) \mu(dx) + \int_S d(y, z) \nu(dy) < \infty. \end{aligned} \tag{A.2.4}$$

From Proposition A.1.7, it follows that Lipschitz functions separate probability measures. Therefore, by Proposition A.1.1,  $d_W$  is indeed a metric. More precisely, the conditions of Proposition A.1.7 are fulfilled because the indicator of any closed set can be expressed as a limit of a uniformly bounded sequence of Lipschitz functions, while closed sets form a  $\pi$ -system generating the Borel  $\sigma$ -algebra  $\mathcal{S}$ . This completes the proof.  $\square$

There are several alternative definitions of the Wasserstein metric. A very important one, stated here as a theorem, is based on *couplings*.

**Theorem A.2.5.** *If  $S$  is separable, then for any two probability measures  $\mu, \nu \in \text{Pr}^{L^1}(S, d)$ , we have*

$$d_W(\mu, \nu) = \inf \mathbb{E}[d(X, Y)], \tag{A.2.5}$$

where the minimum runs over all pairs of random variables  $X$  and  $Y$  defined on the same probability space, where  $X$  follows the distribution  $\mu$  and  $Y$  follows the distribution  $\nu$ .

PARTIAL PROOF. For  $X \sim \mu$  and  $Y \sim \nu$  being defined on the same probability space and for  $f: S \rightarrow \mathbb{R}$  with  $M_1(f) \leq 1$ , observe that

$$|\langle f, \mu \rangle - \langle f, \nu \rangle| = |\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]| \leq \mathbb{E}|f(X) - f(Y)| \leq \mathbb{E}[d(X, Y)].$$

Taking the supremum over all  $f$  in the left hand side and the infimum over all appropriate pairs  $(X, Y)$  in the right hand side, we find that  $d_W(\mu, \nu) \leq \inf \mathbb{E}[d(X, Y)]$ . The proof of the opposite inequality is much more difficult and will be omitted here, but see Rachev [25].  $\square$

Clearly, the total variation distance is stronger than the Kolmogorov distance. However, these two metrics are in general uncomparable to the Wasserstein distance, although in certain special cases, comparison is possible.

**Example A.2.6.** On  $\mathbb{Z}$ , we have  $M_1(\mathbf{1}_A) \leq 1$  for all  $A \subseteq \mathbb{Z}$ . Therefore,  $d_{\text{TV}} \leq d_{\text{W}}$ , so that the Wasserstein metric is uniformly stronger than the total variation metric.

**Example A.2.7.** If  $S$  is bounded, the total variation metric is stronger than the Wasserstein metric – see Corollary A.5.9.

**Example A.2.8.** On an unbounded metric space with at least one accumulation point, the total variation and the Wasserstein metric are uncomparable: neither is stronger than the other. This is shown by the following two counterexamples: first, take an accumulation point  $x$ , so that there exist a sequence of points  $x_1, x_2, \dots$ , which are all different from  $x$  and converge to  $x$ . Clearly,  $d_{\text{TV}}(\delta_{x_n}, \delta_x) = 1$  and  $d_{\text{W}}(\delta_{x_n}, \delta_x) = d(x_n, x)$ . Therefore, the sequence  $\delta_{x_1}, \delta_{x_2}, \dots$  converges to  $\delta_x$  in the Wasserstein metric, but it does not converge in the total variation metric. Next, choose an arbitrary point  $y \in S$ . Since  $S$  is unbounded, there exists a sequence  $y_1, y_2, \dots$ , such that  $d(y_n, y) \geq n$  for all  $n$ . Now consider the sequence of probability measures  $\mu_1, \mu_2, \dots$  defined as

$$\mu_n := \begin{pmatrix} y & y_n \\ 1 - \frac{1}{n} & \frac{1}{n} \end{pmatrix} = \left(1 - \frac{1}{n}\right) \delta_y + \frac{1}{n} \delta_{y_n}. \quad (\text{A.2.6})$$

Since  $d_{\text{TV}}(\mu_n, \delta_y) = 1/n$ , this sequence converges to the Dirac measure  $\delta_y$  in the total variation metric. However, taking the test function  $f(x) := d(x, x_0)$ , we find that  $\langle f, \mu_n \rangle \geq 1$  for all  $n \geq 1$ , while  $\langle f, \delta_{x_0} \rangle = 0$ , so that the sequence  $\mu_n$  does not converge to  $\delta_{x_0}$  in the Wasserstein metric.

Considering the same two examples on the real line, we find that the Kolmogorov and the Wasserstein metric are uncomparable, too.

However, it turns out that for some probability measures  $\nu$ , any neighbourhood of  $\nu$  in the Kolmogorov metric restricted to  $\text{Pr}^{L^1}(S, d)$  is also a neighbourhood of  $\nu$  in the Wasserstein metric. Consequently, any sequence of probability measures in  $\text{Pr}^{L^1}(S, d)$  which converges to  $\nu$  in the Wasserstein metric also converges to  $\nu$  in the Kolmogorov metric.

**Proposition A.2.9.** *If  $\nu \in \text{Pr}^{L^1}(\mathbb{R})$  is a probability measure with density bounded from above by  $B$ , then we can estimate*

$$d_{\text{K}}(\mu, \nu) \leq \sqrt{2B d_{\text{W}}(\mu, \nu)} \quad (\text{A.2.7})$$

for all  $\mu \in \text{Pr}^{L^1}(\mathbb{R})$ .

**PROOF.** For any  $a \in \mathbb{R}$ , define  $f_a(x) := \mathbf{1}(x \leq a)$ , so that  $d_{\text{K}}(\mu, \nu) = \sup_{a \in \mathbb{R}} |\langle f_a, \mu \rangle - \langle f_a, \nu \rangle|$ . Next, for any  $a < b$ , define

$$f_{a,b}(x) := \begin{cases} 1 & ; x \leq a \\ \frac{b-x}{b-a} & ; a \leq x \leq b \\ 0 & ; x \geq b \end{cases}$$

and observe that  $M_1(f_{a,b}) = 1/(b-a)$ . Therefore,  $|\langle f_{a,b}, \mu \rangle - \langle f_{a,b}, \nu \rangle| \leq d_W(\mu, \nu)/(b-a)$ . Next, observe that

$$\begin{aligned}\langle f_{a,b}, \nu \rangle - \langle f_a, \nu \rangle &= \int_a^b \frac{b-x}{b-a} \nu(dx) \leq B \int_a^b \frac{b-x}{b-a} dx = \frac{B(b-a)}{2}, \\ \langle f_b, \nu \rangle - \langle f_{a,b}, \nu \rangle &= \int_a^b \frac{x-a}{b-a} \nu(dx) \leq B \int_a^b \frac{x-a}{b-a} dx = \frac{B(b-a)}{2}.\end{aligned}$$

Therefore, for any  $a \in \mathbb{R}$  and  $\varepsilon > 0$ ,

$$\begin{aligned}\langle f_a, \mu \rangle - \langle f_a, \nu \rangle &\leq \langle f_{a,a+\varepsilon}, \mu \rangle - \langle f_{a,a+\varepsilon}, \nu \rangle + \frac{B\varepsilon}{2}, \\ \langle f_a, \nu \rangle - \langle f_a, \mu \rangle &\leq \langle f_{a-\varepsilon,a}, \nu \rangle - \langle f_{a-\varepsilon,a}, \mu \rangle + \frac{B\varepsilon}{2}.\end{aligned}$$

Noting that the functions  $f_{a,a+\varepsilon}$  and  $f_{a-\varepsilon,a}$  have Lipschitz constant  $1/\varepsilon$  and combining both estimates with the bound in terms of the Wasserstein distance, we find that

$$|\langle f_a, \mu \rangle - \langle f_a, \nu \rangle| \leq \frac{d_W(\mu, \nu)}{\varepsilon} + \frac{B\varepsilon}{2}.$$

Taking the supremum over all  $a$ , we obtain

$$d_K(\mu, \nu) \leq \frac{d_W(\mu, \nu)}{\varepsilon} + \frac{B\varepsilon}{2}.$$

Optimization over  $\varepsilon$  completes the proof.  $\square$

### A.3 More on the Kolmogorov metric

The Kolmogorov metric measures the distance between two probability measures in terms of the indicators of half-lines. However, we often need to consider more general test functions. Here, we show that we can take functions with bounded total variation on the whole real line. We refer to the results listed in Section B.2.

In view of definition B.2.1, define the total variation of a function  $f$  on the whole real line as

$$V(f) := V(f; \mathbb{R}) = \sup \sum_{i=1}^n |f(x_i) - f(x_{i-1})|,$$

where the supremum runs over all possible finite sequences  $x_0 \leq x_1 \leq x_2 \leq \dots \leq x_n$ . A function  $f$  has *bounded variation* if its total variation is finite.

**Proposition A.3.1.** *For any two probability measures  $\mu$  and  $\nu$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , we have*

$$d_K(\mu, \nu) = \sup \{ |\langle f, \mu \rangle - \langle f, \nu \rangle| ; V(f) \leq 1 \}.$$

**Remark A.3.2.** If  $f$  has bounded variation,  $\langle f, \mu \rangle$  exists for all Borel measures  $\mu$  because  $f$  is bounded and Borel measurable.



**Corollary A.3.3.** *For any two probability measures  $\mu$  and  $\nu$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  and any function  $f$  with bounded variation, we have*

$$|\langle f, \mu \rangle - \langle f, \nu \rangle| \leq d_K(\mu, \nu) V(f).$$

PROOF OF PROPOSITION A.3.1. Letting  $\delta := d_K(\mu, \nu)$ , it suffices to prove that

$$|\langle f, \mu \rangle - \langle f, \nu \rangle| \leq \delta V(f) \tag{A.3.1}$$

for all functions  $f: \mathbb{R} \rightarrow \mathbb{R}$  with bounded variation. We shall do it step by step, knowing that (A.3.1) holds for functions of the form  $f(x) = \mathbf{1}(x \leq a)$ .

*Step 1:* assume  $f(x) = \mathbf{1}(x < a)$ . Letting  $f_n(x) = \mathbf{1}(x \leq a - \frac{1}{n})$ , observe that the sequence  $f_n$  is monotone increasing and converges pointwise to  $f$ . The desired inequality (A.3.1) now follows from the monotone convergence theorem.

*Step 2:* assume that  $f$  is a *basic jump function*, i. e.,

$$f(x) = \begin{cases} 0 & x < a \\ \theta & x = a \\ 1 & x > a \end{cases}$$

for some  $\theta \in [0, 1]$ . Noting that  $V(f) = 1$  and  $f(x) = 1 - (1 - \theta) \mathbf{1}(x \leq a) - \theta \mathbf{1}(x < a)$  and applying Step 1, we deduce that (A.3.1) is true for  $f$ .

*Step 3:* assume that  $f$  is a bounded monotone increasing *jump function*, i. e., of the form  $f = \sum_{n=1}^{\infty} c_n f_n$ , where  $\sum_{n=1}^{\infty} c_n < \infty$ . By part (3) of Proposition B.2.7, we have  $V(f) = \sum_{n=1}^{\infty} c_n$ . Clearly,  $\langle f, \mu \rangle = \sum_{n=1}^{\infty} \langle f_n, \mu \rangle$  and  $\langle f, \nu \rangle = \sum_{n=1}^{\infty} \langle f_n, \nu \rangle$ . Therefore, (A.3.1) is true for  $f$ .

*Step 4:* assume that  $f$  is absolutely continuous, bounded and monotone increasing. In this case,  $f(x) = \int_{-\infty}^x f'(a) da$  (see Theorem 6.4.2 of Heil [20]). By part (3) of Proposition B.2.7, we have  $V(f) = \int_{-\infty}^{\infty} f'(a) da$ . Rewriting the integral as  $f(x) = \int_{-\infty}^{\infty} f'(a) \mathbf{1}(x \geq a) da$  and applying Fubini's theorem, we find that  $\langle f, \mu \rangle = \int_{-\infty}^{\infty} \mu([a, \infty)) f'(a) da$  and  $\langle f, \nu \rangle = \int_{-\infty}^{\infty} \nu([a, \infty)) f'(a) da$ , so that (A.3.1) is true for  $f$ .

*Step 5:* assume that  $f$  is bounded and monotone increasing. By Lemma 1.6.31 (iii) of Tao [33],  $f$  can be decomposed as  $f = g + h$ , where  $g$  is a jump function  $g$  is absolutely continuous, and both functions are bounded and monotone increasing. Again, by part (3) of Proposition B.2.7, we have  $V(f) = V(g) + V(h)$ . Therefore, (A.3.1) follows from Steps 3 and 4.

*Step 6:* assume the general case where  $f$  has bounded variation. In this case, (A.3.1) follows from the previous step and the Jordan decomposition theorem (Theorem B.2.4).  $\square$

## A.4 Weak topologies

In Section A.1, we introduced a construction of a metric on the set of probability measures based on a given class of test functions  $\mathcal{F}$ . This metric induces a topology, which will be

referred to as the *metric topology* with respect to  $\mathcal{F}$ . However, the class  $\mathcal{F}$  introduces another very natural topology.

**Definition A.4.1.** Let the set  $\mathcal{M}$  of probability measures and the class  $\mathcal{F}$  of test functions be such that  $\mu(|f|) < \infty$  for all  $\mu \in \mathcal{M}$  and  $f \in \mathcal{F}$ . The *weak topology* on  $\mathcal{M}$  with respect to  $\mathcal{F}$  is the weakest topology, such that the functionals  $\mu \mapsto \langle f, \mu \rangle$  are continuous for all  $f \in \mathcal{F}$ . In other words, this is the topology with a subbasis consisting of all sets of the form

$$G(f, U) := \{\nu \in \mathcal{M} ; \langle f, \nu \rangle \in U\}, \quad (\text{A.4.1})$$

where  $f \in \mathcal{F}$  and  $U$  is an open set in  $\mathbb{R}$ .

**Remark A.4.2.** For each  $\mu \in \mathcal{M}$ , one can easily check that the sets

$$N(\mu; f_1, \dots, f_n; \varepsilon_1, \dots, \varepsilon_n) := \{\nu \in \mathcal{M} ; |\langle f_i, \nu \rangle - \langle f_i, \mu \rangle| < \varepsilon_i, i = 1, \dots, n\}, \quad (\text{A.4.2})$$

where  $f_1, \dots, f_n \in \mathcal{F}$  and  $\varepsilon_1, \dots, \varepsilon_n > 0$ , form a fundamental system of neighbourhoods of  $\mu$ . Therefore, in this topology, a sequence of probability measures  $\mu_n$  converges to a probability measure  $\mu$  if and only if the sequence  $\langle f, \mu_n \rangle$  converges to  $\langle f, \mu \rangle$  for each  $f \in \mathcal{F}$ . We say that the sequence  $\mu_n$  *weakly converges* to  $\mu$ .

Comparing the topology introduced here with the topology from Section A.1, we obtain the following result.

**Proposition A.4.3.** For given  $\mathcal{M}$  and  $\mathcal{F}$ , the underlying weak topology is weaker than the underlying metric topology.

PROOF. Let  $K(\mu, r)$  denote the open ball about  $\mu$  with radius  $r$ , that is,  $K(\mu, r) = \{\nu ; d(\mu, \nu) < r\}$ . It suffices to prove that each set  $N(\mu; f_1, \dots, f_n; \varepsilon_1, \dots, \varepsilon_n)$  contains an open ball of the form  $K(\mu, r)$ . However, this is fulfilled by choosing  $r := \min\{\varepsilon_1, \dots, \varepsilon_n\}$ .  $\square$

**Remark A.4.4.** Typically, the underlying metric topology is *strictly* stronger than the underlying weak topology, as illustrated in Example A.4.10 below.

Topologies have numerous important properties. The following assertion concerns two of them, which are satisfied by all metric topologies. The proof is left to the reader.

**Proposition A.4.5.**

- (1) The weak topology defined by the subbasis given in (A.4.1) is Hausdorff if and only if the space  $\mathcal{F}$  of test functions separates the probability measures in  $\mathcal{M}$ .
- (2) If  $\mathcal{F}$  is countable, then the weak topology is first countable.

$\square$

**Definition A.4.6.** Let  $S$  be endowed with a topology and let  $\mathcal{S}$  be the underlying Borel  $\sigma$ -algebra. The *usual weak topology* on  $\text{Pr}(S, \mathcal{S})$  is the one with respect to the class of all bounded continuous functions.

**Proposition A.4.7.** *If  $S$  is metrizable, then the class of all bounded continuous functions on  $S$  separates all probability measures.*

PROOF. We apply Proposition A.1.7: it suffices to prove that there exists a  $\pi$ -system  $\mathcal{P}$ , which generates  $\mathcal{S}$  and is such that for any set  $A \in \mathcal{P}$ , there exists a uniformly bounded sequence of continuous functions which converges pointwise to the indicator  $\mathbf{1}_A$  of the set  $A$ . However, since  $\mathcal{S}$  is the Borel  $\sigma$ -algebra, the family of all closed sets, which is a  $\pi$ -system, generates  $\mathcal{S}$ . Thus, for any closed set  $A$ , we need to construct a uniformly bounded sequence of continuous functions which converges pointwise to  $\mathbf{1}_A$ . If  $A$  is empty, we may take all of them to be zero. Otherwise, take a metric  $d$  on  $S$  and put

$$f_n(x) := (1 - n d(x, A))_+ .$$

This completes the proof. □

Combining Propositions A.4.5 and A.4.7 leads to

**Corollary A.4.8.** *If  $S$  is metrizable, then the usual weak topology on  $\text{Pr}(S, \mathcal{S})$  is Hausdorff.* □

The Hausdorff property is all that we shall need from the properties of the weak topology on a space of probability measures. However, much more can be proved.

**Theorem A.4.9.** *If  $S$  is separable and metrizable, then  $\text{Pr}(S, \mathcal{S})$  with the weak topology is also metrizable.*

We shall omit the proof. For the case where  $S$  is *complete*, see, for example, Rogers and Williams [27], pp. 205–209. This proof is very indirect (it refers to the Banach–Alaoglu theorem). However, the result can also be proved directly, by considering the *Lévy–Prokhorov metric* on  $\text{Pr}(S, \mathcal{S})$  based on a metric  $d$  in  $S$ :

$$\rho(\mu, \nu) := \inf\{\varepsilon \geq 0 ; \mu(A) \leq \nu(A^\varepsilon) + \varepsilon \text{ for all closed sets } A\} , \quad (\text{A.4.3})$$

where  $A^\varepsilon := \{x \in S ; d(x, A) < \varepsilon\}$ . Some rather involved calculation shows that this is really a metric which induces the weak topology. A substantial part of this is derived in Ethier and Kurtz [18] on pages 96–110, in fact all except for the fact that the weak topology is stronger than the topology induced by the Lévy–Prokhorov metric. This has to be verified directly – referring just to sequences does not suffice.

Now we turn to the example where the metric topology is strictly stronger than the weak topology with respect to the same class of test functions.

**Example A.4.10.** Let  $S$  be endowed with a topology, under which  $S$  is metrizable and is not discrete, and let  $\mathcal{S}$  be the Borel  $\sigma$ -algebra. Let  $\mathcal{F}$  be the class of all bounded continuous functions  $S \rightarrow [0, 1]$ . It is clear that the weak topology with respect to this class is just the usual weak topology. We shall show that it is strictly weaker than the underlying metric topology, noting that  $d_{\mathcal{F}}$  is a metric because  $\mathcal{F}$  separates all probability measures by Proposition A.4.7.

Since  $S$  is not discrete, there exists a point  $x$ , such that  $\{x\}$  is not an open set. Since  $S$  is metrizable, there exists a sequence of points  $x_n$  different from  $x$ , which converges to  $x$ . Then it is clear that the sequence of the Dirac measures  $\delta_{x_n}$  weakly converges to  $\delta_x$ . On the other hand, it is clear that for each  $n \in \mathbb{N}$ , there exists a function  $f \in \mathcal{F}$ , such that  $f(x) = 0$  and  $f(x_n) = 1$ , so that  $d_{\mathcal{F}}(\delta_x, \delta_{x_n}) = 1$ . As a result, the sequence does not converge in the metric topology.

On the real line, the usual weak topology is characterized by the following property.

**Proposition A.4.11.** *Let  $\mu_1, \mu_2, \dots$ , and  $\mu$  be Borel probability measures on the real line with the underlying cumulative distribution functions  $F_1, F_2, \dots$  and  $F$ , that is,  $F_n(a) = \mu_n((-\infty, a])$  and  $F(a) = \mu((-\infty, a])$ . Then the sequence  $\mu_1, \mu_2, \dots$  converges to  $\mu$  in the usual weak topology if and only if the values  $F_1(c), F_2(c), \dots$  converge to  $F(c)$  for each  $c$  where  $F$  is continuous, i. e., where  $\mu(\{c\}) = 0$ .*

PROOF. Suppose first that the measures weakly converge. Similarly as in the proof of Proposition A.2.9, for any  $a < b$ , define

$$f_{a,b}(x) := \begin{cases} 1 & ; x \leq a \\ \frac{b-x}{b-a} & ; a \leq x \leq b \\ 0 & ; x \geq b \end{cases} .$$

Observe that  $f_{a,b}(x) = \frac{1}{b-a} \int_a^b \mathbf{1}_{(-\infty, t]}(x) dt$ , so that

$$\langle f_{a,b}, \mu \rangle = \frac{1}{b-a} \int_a^b F(t) dt .$$

Now take a point  $c$  where  $F$  is continuous and take  $\varepsilon > 0$ . There exists  $\delta > 0$ , such that  $|F(x) - F(c)| < \varepsilon/2$  for all  $x$  with  $|x - c| < \delta$ . Consequently,

$$F(c) - \frac{\varepsilon}{2} < \langle f_{c-\delta, c}, \mu \rangle \leq F(c) \leq \langle f_{c, c+\delta}, \mu \rangle \leq F(c) + \frac{\varepsilon}{2} .$$

Since the measures weakly converge, there exists  $n_0$ , such that  $|\langle f_{c, c+\delta}, \mu_n \rangle - \langle f_{c, c+\delta}, \mu \rangle| < \varepsilon/2$  and  $|\langle f_{c-\delta, c}, \mu_n \rangle - \langle f_{c-\delta, c}, \mu \rangle| < \varepsilon/2$  for all  $n \geq n_0$ . Now estimate

$$F_n(c) - F(c) < \langle f_{c, c+\delta}, \mu_n \rangle - \langle f_{c, c+\delta}, \mu \rangle + \frac{\varepsilon}{2} < \varepsilon$$

and

$$F(c) - F_n(c) < \langle f_{c-\delta, c}, \mu \rangle + \frac{\varepsilon}{2} - \langle f_{c-\delta, c}, \mu_n \rangle < \varepsilon$$

This proves that the sequence  $F_1(c), F_2(c), \dots$  converges to  $F(c)$ .

Now we turn to the converse: suppose that the values  $F_1(c), F_2(c), \dots$  converge to  $F(c)$  for each  $c$  where  $F$  is continuous. Equivalently, the sequence  $\mu_1(\mathbf{1}_{(-\infty, c]}), \mu_2(\mathbf{1}_{(-\infty, c]}), \dots$  converges to  $\mu(\mathbf{1}_{(-\infty, c]})$ . Consequently, for any function of the form

$$g = a_0 + \sum_{k=1}^m a_k \mathbf{1}_{(-\infty, c_k]} ,$$

where  $F$  is continuous at all points  $c_1, \dots, c_m$  and where  $a_0, a_1, \dots, a_m$  are any real numbers, the sequence  $\langle g, \mu_1 \rangle, \langle g, \mu_2 \rangle, \dots$  converges to  $\langle g, \mu \rangle$ .

Now take a continuous function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with  $|f(x)| \leq B$  for all  $x \in \mathbb{R}$ . We can assume that  $B > 0$ . Take  $\varepsilon > 0$ . Recalling that continuity points of  $F$  are everywhere dense, observe that there exist points  $a < b$ , where  $F$  is continuous and are such that  $\mu(-\infty, a] < \varepsilon/(15B)$  and  $\mu([b, \infty)) < \varepsilon/(15B)$ . Next, there exists  $n_0$ , such that  $|F_n(a) - F(a)| < \varepsilon/(15B)$  and  $|F_n(b) - F(b)| < \varepsilon/(15B)$  for all  $n \geq n_0$ . As a result, we have  $\mu_n(-\infty, a] < 2\varepsilon/(15B)$  and  $\mu_n((b, \infty)) < 2\varepsilon/(15B)$ .

Since  $f$  is uniformly continuous on  $[a, b]$ , there exists  $\delta > 0$ , such that  $|f(x) - f(y)| < \varepsilon/15$  for all  $x, y$  with  $|x - y| < \delta$ . Next, there exist points  $a = c_0 < c_1 < \dots < c_m = b$ , such that  $c_k - c_{k-1} < \delta$  for all  $k = 1, 2, \dots, m$  and  $f$  is continuous at all  $c_k$ . Consider the function

$$g(x) = \begin{cases} f(c_1) & ; x \leq c_1 \\ f(c_2) & ; c_1 < x \leq c_2 \\ \dots & \\ f(c_{m-1}) & ; c_{m-2} < x \leq c_{m-1} \\ f(c_m) & ; x > c_{m-1}. \end{cases}$$

Observe that  $|g(x) - f(x)| < \varepsilon/15$  for all  $a \leq x \leq b$  and  $|g(x) - f(x)| < 2B$  for all other  $x$ . Therefore,

$$\begin{aligned} |\langle f, \mu \rangle - \langle g, \mu \rangle| &\leq \int_{(-\infty, a]} |g - f| d\mu + \int_{(a, b]} |g - f| d\mu + \int_{(b, \infty)} |g - f| d\mu \\ &< 2B \frac{\varepsilon}{15B} + \frac{\varepsilon}{15} + 2B \frac{\varepsilon}{15B} = \frac{5\varepsilon}{15} \end{aligned}$$

and similarly,

$$\begin{aligned} |\langle f, \mu_n \rangle - \langle g, \mu_n \rangle| &\leq \int_{(-\infty, a]} |g - f| d\mu_n + \int_{(a, b]} |g - f| d\mu_n + \int_{(b, \infty)} |g - f| d\mu_n \\ &< 2B \frac{2\varepsilon}{15B} + \frac{\varepsilon}{15} + 2B \frac{2\varepsilon}{15B} = \frac{9\varepsilon}{15} \end{aligned}$$

for all  $n \geq n_0$ . Finally, since

$$g = \sum_{k=1}^{m-1} (f(c_k) - f(c_{k+1})) \mathbf{1}_{(-\infty, c_k]} + f(c_m),$$

the sequence  $\langle g, \mu_1 \rangle, \langle g, \mu_2 \rangle, \dots$  converges to  $\langle g, \mu \rangle$ . Therefore, there exists  $n_1 \geq n_0$ , such that  $|\langle g, \mu_n \rangle - \langle g, \mu \rangle| < \varepsilon/15$  for all  $n \geq n_1$  and we have

$$|\langle f, \mu_n \rangle - \langle f, \mu \rangle| \leq |\langle f, \mu_n \rangle - \langle g, \mu_n \rangle| + |\langle g, \mu_n \rangle - \langle g, \mu \rangle| + |\langle g, \mu \rangle - \langle f, \mu \rangle| < \varepsilon,$$

so that the sequence  $\langle f, \mu_1 \rangle, \langle f, \mu_2 \rangle, \dots$  indeed converges to  $\langle f, \mu \rangle$ .  $\square$

## A.5 Change of class of test functions

We often consider different classes of test functions. In view of this, it is useful to compare the underlying metric and weak topologies. In this section, we shall consider three important relationships between the classes: inclusion, linear combinations and limits. The first case is immediate and is therefore stated without proof.

**Proposition A.5.1.** *Let  $\mathcal{F}$  and  $\mathcal{G}$  be classes of test functions with  $\mathcal{F} \subseteq \mathcal{G}$ . Assume that  $\langle |g|, \mu \rangle < \infty$  for all  $\mu \in \mathcal{M}$  and  $g \in \mathcal{G}$ .*

- (1) *Let  $d_{\mathcal{F}}$  and  $d_{\mathcal{G}}$  be defined as in (A.1.1). If  $d_{\mathcal{F}}$  is a metric on  $\mathcal{M}$ , so is  $d_{\mathcal{G}}$  and we have  $d_{\mathcal{F}} \leq d_{\mathcal{G}}$ . Thus, the metric  $d_{\mathcal{F}}$  is uniformly weaker than  $d_{\mathcal{G}}$ .*
- (2) *The weak topology with respect to  $\mathcal{F}$  is weaker than the weak topology with respect to  $\mathcal{G}$ .*

□

**Proposition A.5.2.** *Let  $\mathcal{F}$  be a class of test functions. Assume that  $\langle |f|, \mu \rangle < \infty$  for all  $\mu \in \mathcal{M}$  and  $f \in \mathcal{F}$ .*

- (1) *If*

$$\tilde{\mathcal{F}} = \left\{ \alpha_0 + \sum_{i=1}^n \alpha_i f_i ; f_i \in \mathcal{F}, \sum_{i=1}^n |\alpha_i| \leq 1, n \in \mathbb{N} \right\}, \quad (\text{A.5.1})$$

*then the metric with respect to  $\mathcal{F}$  agrees with the one with respect to  $\tilde{\mathcal{F}}$ .*

- (2) *The weak topology with respect to  $\mathcal{F}$  agrees with the weak topology with respect to  $\text{span}(\{1\} \cup \mathcal{F})$ .*

**Corollary A.5.3.** *If  $\mathcal{F} \subseteq \mathcal{G} \subseteq \tilde{\mathcal{F}}$ , the metric with respect to  $\mathcal{G}$  agrees with the metric with respect to  $\mathcal{F}$ . Similarly, if  $\mathcal{F} \subseteq \mathcal{G} \subseteq \text{span}(\{1\} \cup \mathcal{F})$ , the weak topology with respect to  $\mathcal{G}$  coincides with the weak topology with respect to  $\mathcal{F}$ . In particular, the weak topology with respect to  $\tilde{\mathcal{F}}$  coincides with the weak topology with respect to  $\mathcal{F}$ .* □

PROOF OF PROPOSITION A.5.2.

(1): Denote the underlying metrics by  $d_{\mathcal{F}}$  and  $d_{\tilde{\mathcal{F}}}$  and take any probability measures  $\mu, \nu \in \mathcal{M}$ . Clearly,  $d_{\mathcal{F}}(\mu, \nu) \leq d_{\tilde{\mathcal{F}}}(\mu, \nu)$ . However, we also have  $|\langle f, \nu \rangle - \langle f, \mu \rangle| \leq d_{\mathcal{F}}(\mu, \nu)$  for all  $f \in \tilde{\mathcal{F}}$ , so that  $d_{\tilde{\mathcal{F}}}(\mu, \nu) \leq d_{\mathcal{F}}(\mu, \nu)$ .

(2): It suffices to prove that the weak topology with respect to  $\text{span}(\{1\} \cup \mathcal{F})$  is weaker than the weak topology with respect to  $\mathcal{F}$ . In order to prove the latter, it suffices to prove that the set  $G(f, U)$  defined as in (A.4.1) is open in the weak topology with respect to  $\mathcal{F}$  for each  $f \in \text{span}(\{1\} \cup \mathcal{F})$  and each open set  $U \subseteq \mathbb{R}$ . Take arbitrary  $f = \alpha_0 + \sum_{i=1}^n \alpha_i f_i \in \text{span}(\{1\} \cup \mathcal{F})$ , where  $f_i \in \mathcal{F}$ , and arbitrary  $\mu \in G(f, U)$ , that is,  $\langle f, \mu \rangle = \alpha_0 + \sum_{i=1}^n \alpha_i \langle f_i, \mu \rangle \in U$ . Since the map  $(y_1, \dots, y_n) \mapsto \sum_{i=1}^n \alpha_i y_i$  is continuous,

there exist open sets  $U_i \subseteq \mathbb{R}$ , such that  $\langle f_i, \mu \rangle \in U_i$  and  $\alpha_0 + \sum_{i=1}^n \alpha_i y_i \in U$  for all  $y_i \in U_i$ . Therefore,  $\mu \in \bigcap_{i=1}^n G(f_i, U_i) \subseteq G(f, U)$ , implying that  $G(f, U)$  is indeed open in the weak topology with respect to  $\mathcal{F}$ . This completes the proof.  $\square$

**Proposition A.5.4.** *Let  $\mathcal{F}$  and  $\mathcal{G}$  be classes of measurable test functions on  $(S, \mathcal{S})$ . Assume that  $\langle |f|, \mu \rangle < \infty$  for all  $\mu \in \mathcal{M}$  and  $f \in \mathcal{F} \cup \mathcal{G}$ .*

- (1) *Let  $d_{\mathcal{F}}$  and  $d_{\mathcal{G}}$  denote the metrics with respect to  $\mathcal{F}$  and  $\mathcal{G}$ , respectively, by (A.1.1). If for each  $f \in \mathcal{F}$ , there exists a sequence of test functions  $f_n \in \mathcal{G}$ , such that  $\lim_{n \rightarrow \infty} \langle f_n, \mu \rangle = \langle f, \mu \rangle$  for all  $\mu \in \mathcal{M}$ , then  $d_{\mathcal{F}} \leq d_{\mathcal{G}}$ .*
- (2) *If for each  $f \in \mathcal{F}$  and each  $\mu \in \mathcal{M}$ , there exist sequences of functions  $f_n^+, f_n^- \in \mathcal{G}$ , such that  $f_n^- \leq f \leq f_n^+$  and  $\lim_{n \rightarrow \infty} \langle f_n^+, \mu \rangle = \lim_{n \rightarrow \infty} \langle f_n^-, \mu \rangle = \langle f, \mu \rangle$ , then the weak topology on  $\mathcal{M}$  with respect to  $\mathcal{G}$  is stronger than the one with respect to  $\mathcal{F}$ .*

**Corollary A.5.5.** *Under the assumptions of the preceding lemma along with the assumption that  $\mathcal{G} \subseteq \mathcal{F}$ , the metrics as well as the weak topologies based on  $\mathcal{F}$  and  $\mathcal{G}$  agree.*  $\square$

### Remarks.

- (1) The conditions of part (1) are fulfilled if the sequence  $f_n$  converges to  $f$  pointwise and is either monotone or uniformly bounded.
- (2) The conditions of part (2) are fulfilled if the sequences  $f_n^+$  and  $f_n^-$  pointwise converge to  $f$  and one of the following two conditions is fulfilled: either both sequences are uniformly bounded or  $f_n^+$  is monotonically decreasing, while  $f_n^-$  is monotonically increasing.

### PROOF OF PROPOSITION A.5.4.

(1): Immediate.

(2): For all  $\mu \in \mathcal{M}$ ,  $f \in \mathcal{F}$  and  $\varepsilon > 0$ , it suffices to construct a neighbourhood of  $\mu$  in the weak topology with respect to  $\mathcal{G}$ , which is contained in the set  $N(\mu; f; \varepsilon)$  defined as in (A.4.2). Taking  $\varepsilon > 0$  and appropriate sequences  $f_n^+$  and  $f_n^-$ , there exists  $n$ , such that  $\langle f, f_n^+ \rangle - \varepsilon/2 < \langle f, \mu \rangle < \langle f_n^-, \mu \rangle + \varepsilon/2$ . Define

$$U := \left\{ \nu \in \mathcal{M} ; \langle f_n^+, \nu \rangle - \langle f_n^+, \mu \rangle < \frac{\varepsilon}{2}, \langle f_n^-, \nu \rangle - \langle f_n^-, \mu \rangle < \frac{\varepsilon}{2} \right\}. \quad (\text{A.5.2})$$

Clearly,  $\mu \in U$  and  $U$  is  $\mathcal{G}$ -weakly open. Moreover,

$$\langle f, \nu \rangle \leq \langle f_n^+, \nu \rangle = \langle f_n^+, \nu \rangle - \langle f_n^+, \mu \rangle + \langle f_n^+, \mu \rangle < \langle f, \mu \rangle + \varepsilon \quad (\text{A.5.3})$$

and

$$\langle f, \nu \rangle \geq \langle f_n^-, \nu \rangle = \langle f_n^-, \nu \rangle - \langle f_n^-, \mu \rangle + \langle f_n^-, \mu \rangle > \langle f, \mu \rangle - \varepsilon \quad (\text{A.5.4})$$

for all  $\nu \in U$ . Therefore,  $U \subseteq N(\mu; f; \varepsilon)$  and the proof is complete.  $\square$

As a simple example, consider the following characterization of the total variation metric.

**Proposition A.5.6.** *We have*

$$d_{\text{TV}}(\mu, \nu) = \sup_{\substack{f \text{ is measurable} \\ 0 \leq f \leq 1}} |\langle f, \nu \rangle - \langle f, \mu \rangle|. \quad (\text{A.5.5})$$

**Remark A.5.7.** For the total variation metric, some authors take the following one instead:

$$d(\mu, \nu) = \sup_{\substack{f \text{ is measurable} \\ -1 \leq f \leq 1}} |\langle f, \nu \rangle - \langle f, \mu \rangle| = 2 d_{\text{TV}}(\mu, \nu). \quad (\text{A.5.6})$$

**PROOF OF PROPOSITION A.5.6.** By part (1) of Proposition A.5.2, the class of indicators can be replaced by the class of all measurable step functions from  $S$  to  $[0, 1]$ . This is because each measurable  $[0, 1]$ -valued step function can be represented as

$$f(x) = \begin{cases} y_0 & ; x \in H_0 \\ y_1 & ; x \in H_1 \\ \vdots & \vdots \\ y_n & ; x \in H_n, \end{cases}$$

where  $0 = y_0 \leq y_1 \leq \dots \leq y_n \leq 1$  and where  $H_0, H_1, \dots$  are disjoint measurable sets, and this can be further expressed as

$$f = y_1 \mathbf{1}_{H_1} + (y_2 - y_1) \mathbf{1}_{H_1 \cup H_2} + \dots + (y_n - y_{n-1}) \mathbf{1}_{H_1 \cup H_2 \cup \dots \cup H_n},$$

noting that  $\sum_{k=1}^n (y_k - y_{k-1}) = y_n \leq 1$ .

For each measurable  $[0, 1]$ -valued function, there exists a monotone sequence of measurable  $[0, 1]$ -valued step functions converging pointwise to it. Therefore, by part (1) of Proposition A.5.4, the class of all measurable  $[0, 1]$ -valued step functions can be further replaced by the class of all measurable functions from  $S$  to  $[0, 1]$ . This completes the proof.  $\square$

**Corollary A.5.8.** *The topology induced by the total variation metric is stronger than the usual weak topology.*

**Corollary A.5.9.** *If  $S$  is a bounded metric space, then the total variation metric on  $\text{Pr}(S, \mathcal{S})$  is stronger than the Wasserstein metric. More precisely, if  $D$  is the diameter of  $S$ , we have  $d_{\text{W}}(\mu, \nu) \leq D d_{\text{TV}}(\mu, \nu)$ .*  $\square$

**Proposition A.5.10.** *Let  $\mathcal{F}$  be a class of test functions on  $(S, \mathcal{S})$ , such that  $\langle |f|, \mu \rangle < \infty$  for all  $\mu \in \mathcal{M}$ . In addition, suppose that at least one of the following two conditions is fulfilled:*

- (1) *For each open set  $G \subseteq S$ , there exists a sequence of functions  $f_n \in \mathcal{F}$ ,  $f_n \leq \mathbf{1}_G$ , which is uniformly bounded from below and converges pointwise to  $\mathbf{1}_G$ .*
- (2) *For each closed set  $F \subseteq S$ , there exists a sequence of functions  $f_n \in \mathcal{F}$ ,  $f_n \geq \mathbf{1}_F$ , which is uniformly bounded from above and converges pointwise to  $\mathbf{1}_F$ .*



Then the weak topology with respect to  $\mathcal{F}$  is stronger than the usual weak topology on  $\mathcal{M}$ .

PROOF. First observe that, by part (2) of Proposition A.5.2, we may assume without loss of generality that  $\mathcal{F}$  is a vector space and that it contains all constants. Under this assumption,  $1 - f \in \mathcal{F}$  for all  $f \in \mathcal{F}$ , so that conditions (1) and (2) are equivalent. We may assume that both are fulfilled.

We shall use part (2) of Proposition A.5.4: for each bounded and continuous  $f: S \rightarrow \mathbb{R}$  and each  $\mu \in \mathcal{M}$ , we shall construct sequences of functions  $f_n^+, f_n^- \in \mathcal{F}$ , such that  $f_n^- \leq f \leq f_n^+$  and  $\lim_{n \rightarrow \infty} \langle f_n^+, \mu \rangle = \lim_{n \rightarrow \infty} \langle f_n^-, \mu \rangle = \langle f, \mu \rangle$ . Without loss of generality, we can assume that  $f$  is  $[0, 1]$ -valued. For any  $t \in \mathbb{R}$ , let:

$$F_t := \{x ; f(x) \geq t\}, \quad G_t := \{x ; f(x) > t\} \quad (\text{A.5.7})$$

One can easily check that

$$f - \frac{1}{n} \leq \frac{1}{n} \sum_{k=1}^{n-1} \mathbf{1}_{G_{k/n}} \leq f \leq \frac{1}{n} \sum_{k=0}^{n-1} \mathbf{1}_{F_{k/n}} \leq f + \frac{1}{n} \quad (\text{A.5.8})$$

for all  $n \in \mathbb{N}$ . From assumption (1) and the dominated convergence theorem, it follows that there exist functions  $f_{n,1}^-, \dots, f_{n,n-1}^- \in \mathcal{F}$ , such that  $f_{n,k}^- \leq \mathbf{1}_{G_{k/n}}$  and  $\langle f_{n,k}^-, \mu \rangle > \mu(G_{k/n}) - 1/n$  for all  $k$ . Similarly, from assumption (2), it follows that there exist functions  $f_{n,0}^+, \dots, f_{n,n-1}^+ \in \mathcal{F}$ , such that  $f_{n,k}^+ \geq \mathbf{1}_{F_{k/n}}$  and  $\langle f_{n,k}^+, \mu \rangle < \mu(F_{k/n}) + 1/n$  for all  $k$ . Put

$$f_n^+ := \sum_{k=0}^{n-1} f_{n,k}^+, \quad f_n^- := \sum_{k=1}^{n-1} f_{n,k}^-. \quad (\text{A.5.9})$$

Clearly,  $f_n^- \leq f \leq f_n^+$ . A short calculation shows that

$$\langle f, \mu \rangle - \frac{2}{n} \leq \langle f_n^-, \mu \rangle \leq \langle f, \mu \rangle \leq \langle f_n^+, \mu \rangle \leq \langle f, \mu \rangle + \frac{2}{n}, \quad (\text{A.5.10})$$

so that the sequences  $f_n^-$  and  $f_n^+$  satisfy the necessary conditions and the proof is complete.  $\square$

**Corollary A.5.11.** *If  $S$  is endowed with a metric  $d$ , then the topology induced by the Wasserstein metric is stronger than the usual weak topology.*

PROOF. First, by Proposition A.4.3, the topology induced by the Wasserstein metric is stronger than the weak topology with respect to the class of functions  $f$  with  $M_1(f) \leq 1$ . By part (2) of Proposition A.5.2, the latter coincides with the topology with respect to the class of all Lipschitz functions. Now observe that the latter class fulfills condition (2) of Proposition A.5.10: for a closed set  $F$ , the latter is satisfied with

$$f_n(x) := (1 - n d(x, F))_+ \quad (\text{A.5.11})$$

(under the convention that  $d(x, \emptyset) = \infty$ ). This completes the proof.  $\square$

**Corollary A.5.12.** *The topology on  $\text{Pr}(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  induced by the Kolmogorov metric is stronger than the usual weak topology.*

PROOF. Again, by Proposition A.4.3, the topology induced by the Kolmogorov metric is stronger than the weak topology with respect to the class of indicators of all closed half-lines. By part (2) of Proposition A.5.2, the latter coincides with the topology with respect to the class of all linear combinations of indicators of closed half-lines. Now we show that this class fulfills condition (1) of Proposition A.5.10. To see this, observe that each open set  $G$  on the real line is a countably infinite union of disjoint intervals of form  $(u, v]$ . Thus, we may write

$$\mathbf{1}_G = \sum_{k=1}^{\infty} \mathbf{1}_{(u_k, v_k]} .$$

As the underlying partial sums are linear combinations of indicators of closed half-lines and are smaller than  $\mathbf{1}_G$ , condition (1) is indeed satisfied. This completes the proof.  $\square$

# Appendix B

## Some real analysis

### B.1 Differentiation of absolutely continuous functions

Throughout this section, let  $I$  denote an interval on the real line.

The classical fundamental theorem of calculus says that a *continuously differentiable* function  $f: I \rightarrow \mathbb{R}$  satisfies

$$f(b) - f(a) = \int_a^b f'(x) \, dx \tag{B.1.1}$$

for all  $a, b \in I$ . This is also true for many functions which are only *almost everywhere* differentiable, i. e.,  $f(x) = |x|$ .

**Definition B.1.1.** A measurable function  $f': T \rightarrow \mathbb{R}$  is an *almost-everywhere* derivative of a function  $f: I \rightarrow \mathbb{R}$  if the set of points in  $I$  where  $f'$  is not the derivative of  $f$  has Lebesgue measure zero.

However, if  $f'$  is an almost-everywhere derivative of  $f$ , the fundamental theorem of calculus (B.1.1) need not be true. A prominent counter-example is the Cantor–Lebesgue function – see Theorem 5.1.2 of Heil [20].

The Cantor–Lebesgue function is continuous, but its continuity is not ‘nice’ enough. It turns out that the validity of (B.1.1) is equivalent to a stronger form of continuity.

**Definition B.1.2.** A function  $f: [a, b] \rightarrow \mathbb{R}$  is *absolutely continuous* if for each  $\varepsilon > 0$ , there exists  $\delta > 0$ , such that for any (finite or infinite) family of non-overlapping intervals  $[a_k, b_k] \subseteq [a, b]$  with  $\sum_k (b_k - a_k) < \delta$ , we have  $\sum_k |f(b_k) - f(a_k)| < \varepsilon$ .

**Remark B.1.3.** A restriction of an absolutely continuous function to a subinterval is also absolutely continuous.

**Definition B.1.4.** A function  $f: I \rightarrow \mathbb{R}$  is absolutely continuous if all restrictions to closed subintervals of  $I$  are absolutely continuous.

**Example B.1.5.** Every Lipschitz function  $f: I \rightarrow \mathbb{R}$  is absolutely continuous.

The following result states that absolute continuity is sufficient and necessary for the fundamental theorem of calculus to hold. For the proof, the reader is referred to Heil [20], Theorem 6.4.2.

**Definition B.1.6.** A function  $f: I \rightarrow \mathbb{R}$  is *locally Lebesgue integrable* if all restrictions of  $f$  to closed subintervals of  $I$  are Lebesgue integrable. The space of all locally Lebesgue integrable functions on  $I$  will be denoted by  $L^1_{\text{loc}}(I)$ .

**Theorem B.1.7.** A function  $f: I \rightarrow \mathbb{R}$  is *absolutely continuous* if and only if it is differentiable almost everywhere on  $I$  and some/any almost-everywhere derivative  $f'$  of  $f$  is in  $L^1_{\text{loc}}(I)$  (with respect to the Lebesgue measure) and satisfies (B.1.1) for all  $a, b \in I$ ; of course, the integral in the right hand side of (B.1.1) is interpreted as the Lebesgue integral over the underlying interval or its negative value if  $a > b$ .

**Convention B.1.8.** In the context of absolutely continuous functions  $f$ ,  $f'$  will denote an almost-everywhere derivative of  $f$  unless specified otherwise.

In particular, all bounded measurable functions on  $I$  are in  $L^1_{\text{loc}}(I)$ . For a function  $f: I \rightarrow \mathbb{R}$ , define

$$M_1(f) := \begin{cases} \text{ess sup}_{x \in I} |f'(x)| & ; f \text{ is absolutely continuous} \\ \infty & ; \text{otherwise} \end{cases} \quad (\text{B.1.2})$$

(clearly, the definition is independent of the version of  $f'$ ). Next, for  $r = 2, 3, 4, \dots$ , define

$$M_r(f) := \begin{cases} M_1(f^{(r-1)}) & ; f \text{ is } (r-1)\text{-times continuously differentiable on } I \\ \infty & ; \text{otherwise.} \end{cases} \quad (\text{B.1.3})$$

**Proposition B.1.9.** For any function  $f: I \rightarrow \mathbb{R}$ , we have

$$M_1(f) = \sup_{x \neq y} \frac{|f(x) - f(y)|}{|x - y|}.$$

SKETCH OF PROOF. Denote  $L(f) := \sup_{x \neq y} \frac{|f(x) - f(y)|}{|x - y|}$ . First, we show that  $M_1(f) \leq L(f)$ . Without loss of generality, we may assume that  $L(f) < \infty$ . In this case,  $f$  is Lipschitz and therefore absolutely continuous. By Theorem B.1.7, there exists a function  $f'$ , such that  $f'(x)$  is the classical derivative of  $f$  at  $x$  for Lebesgue-almost all  $x \in I$ . However, for such  $x$ , we have  $|f'(x)| \leq L(f)$ . Therefore,  $|f'(x)| \leq L(f)$  for Lebesgue-almost all  $x$ . As a result,  $M_1(f) \leq L(f)$ .

Now we prove that  $L(f) \leq M_1(f)$ . Similarly as before, we may assume that  $M_1(f) < \infty$ . Therefore,  $f$  is absolutely continuous. By Theorem B.1.7, we have  $|f(b) - f(a)| = \left| \int_a^b f'(x) dx \right| \leq |b - a| \text{ess sup}_{x \in I} |f'(x)| = |b - a| M_1(f)$ . Dividing by  $|b - a|$  and taking supremum over  $a$  and  $b$ , we find that  $L(f) \leq M_1(f)$ , completing the proof.  $\square$

## B.2 Functions with bounded variation

Like in the previous section,  $I$  will denote an interval on the real line throughout this section.

**Definition B.2.1.** The *total variation* of a function  $f: I \rightarrow \mathbb{R}$  is defined as

$$V(f) := \sup \sum_{i=1}^n |f(x_i) - f(x_{i-1})|,$$

where the supremum runs over all possible finite sequences  $x_0 \leq x_1 \leq x_2 \leq \dots \leq x_n$  of points in  $I$ . The function  $f$  can be initially defined on a larger set than  $I$ . In this case, denote by  $V(f; I)$  the total variation of the restriction of  $f$  to  $I$ , which will be also referred as the total variation of  $f$  on  $I$ .

A function  $f$  has *bounded variation* on  $I$  if its total variation on  $I$  is finite.

**Remark B.2.2.** Each function with bounded variation is bounded.

The proof of the following assertion is left to the reader as an exercise.

**Proposition B.2.3.** *If  $f: I \rightarrow \mathbb{R}$  is monotone increasing or monotone decreasing, then*

$$V(f) = \sup f - \inf f.$$

Moreover, let  $a_1, a_2, \dots$  be a sequence of points in  $I$  chosen as follows: let  $a$  be the lower and  $b$  the upper endpoint of  $I$ . If  $a \in I$ , let  $a_n = a$  for all  $n$ . Otherwise, let  $a_n$  converge to  $a$ . Similarly, choose a sequence  $b_1, b_2, \dots$ , replacing  $a$  with  $b$ . Then we have

$$V(f) = \lim_{n \rightarrow \infty} |f(b_n) - f(a_n)|.$$

□

**Theorem B.2.4** (Jordan decomposition theorem). *For each function  $f: I \rightarrow \mathbb{R}$  with bounded variation, there exist monotone increasing functions  $g$  and  $h$ , such that  $f = g - h$ . Moreover,  $g$  and  $h$  can be chosen so that  $V(f) = V(g) + V(h)$  and that  $g$  and  $h$  are left/right-continuous in all points where so is  $f$ .*

SKETCH OF PROOF. The existence of  $g$  and  $h$  is proved in Wheeden and Zygmund [34] for the case  $I = [a, b]$ , where  $-\infty < a \leq b < \infty$ : see Theorem 2.7 ibidem. From the construction of  $g$  and  $h$ , one can see that  $V(f) = V(g) + V(h)$  and that  $g$  and  $h$  are left/right-continuous in all points where so is  $f$  (exercise). The extension to the general case is also left to the reader as an exercise. □

**Corollary B.2.5.** *Each function with bounded variation is Borel measurable.*

SKETCH OF PROOF. Observe first that each monotone increasing function has at most countably many discontinuities (see Theorem 2.8 of Wheeden and Zygmund [34]) and is therefore Borel measurable (exercise). The result now follows from the Jordan decomposition theorem. □

The following result is proved in Wheeden and Zygmund [34] – see Theorem 7.31 ibidem.

**Proposition B.2.6.** *If  $f: I \rightarrow \mathbb{R}$  is absolutely continuous, then  $V(f; I) = \int_I |f'(x)| dx$ .*

□

**Proposition B.2.7.**

- (1) For any two functions  $g$  and  $h: I \rightarrow \mathbb{R}$ , we have  $V(g + h) \leq V(g) + V(h)$ .
- (2) Consider a measurable space  $(S, \mathcal{S})$ , a signed measure  $\kappa$  on  $S$  and a measurable function  $F: I \times S \rightarrow \mathbb{R}$ , such that the function  $y \mapsto F(x, y)$  is an element of  $L^1(|\kappa|)$  for all  $x \in I$ . Letting  $f_y(x) := F(x, y)$  and  $f(x) := \int_S F(x, y) \kappa(dy)$ , we have

$$V(f) \leq \int_S V(f_y) |\kappa|(dy). \quad (\text{B.2.1})$$

- (3) If all functions  $f_y$  are monotone increasing and  $\kappa$  is a positive measure, then (B.2.1) holds with equality.

PROOF. It suffices to prove the second part. Taking  $x_0 \leq x_1 \leq \dots \leq x_n$ , observe that

$$\begin{aligned} \sum_{i=1}^n |f(x_i) - f(x_{i-1})| &= \sum_{i=1}^n \left| \int_S (F(x_i, y) - F(x_{i-1}, y)) \kappa(dy) \right| \\ &\leq \int_S \sum_{i=1}^n |F(x_i, y) - F(x_{i-1}, y)| |\kappa|(dy) \\ &\leq \int_S V(f_y) |\kappa|(dy). \end{aligned}$$

Taking the supremum, we obtain (B.2.1). Finally, if the functions  $x \mapsto F(x, y)$  are monotone increasing and  $\kappa$  is a positive measure, take sequences  $a_n$  and  $b_n$  as in Proposition B.2.3, where, in addition,  $a_n$  is monotone decreasing and  $b_n$  is monotone increasing. Now observe that

$$\begin{aligned} V(f) &= \lim_{n \rightarrow \infty} (f(b_n) - f(a_n)) \\ &= \lim_{n \rightarrow \infty} \int_S (F(b_n, y) - F(a_n, y)) \kappa(dy) \\ &= \int_S \lim_{n \rightarrow \infty} (F(b_n, y) - F(a_n, y)) \kappa(dy) \\ &= \int_S V(f_y) \kappa(dy) \end{aligned}$$

by the monotone convergence theorem. This completes the proof.  $\square$

### B.3 The Riemann–Stieltjes integral

The results in this section are mostly just listed. For the proofs, the reader is referred to Chapter 2 and partly Chapter 7 of Wheeden and Zygmund [34], where the results may not be proved in the whole generality. However, the extension is not difficult and is left to the reader as an exercise.

**Definition B.3.1.** Let  $-\infty < a \leq b < \infty$ . Take functions  $f, g: [a, b] \rightarrow \mathbb{R}$ . The *Riemann–Stieltjes integral* of  $f$  with respect to  $g$  is the limit of the Riemann–Stieltjes sums

$$\sum_{k=1}^n f(\xi_k)(x_k - x_{k-1}),$$

where  $a = x_0 \leq \xi_1 \leq x_1 \leq \xi_2 \leq \dots \leq x_{n-1} \leq \xi_n \leq x_n = b$ . More precisely, a number  $J \in \mathbb{R}$  is the Riemann–Stieltjes integral of  $f$  with respect to  $g$  if for each  $\varepsilon > 0$ , there exists  $\delta > 0$ , such that

$$\left| \sum_{k=1}^n f(\xi_k)(x_k - x_{k-1}) - J \right| < \varepsilon$$

for all  $x_0, \dots, x_n$  and  $\xi_1, \dots, \xi_n$  being as above and such that  $x_k - x_{k-1} < \delta$  for all  $k = 1, 2, \dots, n$ . Write

$$J = \int_a^b f \, dg.$$

The function  $f$  is called *integrand*, while the function  $g$  is called *integrator*.

**Remark B.3.2.** The Riemann–Stieltjes integral does not always exist. However, if it exists, it is unique.

**Definition B.3.3.** An *improper Riemann–Stieltjes integral* of  $f$  with respect to  $g$  is a limit of Riemann–Stieltjes integrals in the sense of Definition B.3.1 as one or both endpoints approach their limits – the left endpoint from the right and the right endpoint from the left. Thus, the integral can be improper at one or both endpoints. When the integral is improper at both endpoints, the limit of the integral must exist and be the same regardless how the endpoints approach their limits.

Denoting the limits of the endpoints by  $a$  and  $b$  (which can also be infinite), we also denote the underlying improper integral by  $\int_a^b f \, dg$ , that is,

$$\int_a^b f \, dg = \lim_{a' \downarrow a} \int_{a'}^b f \, dg, \quad \lim_{b' \uparrow b} \int_a^{b'} f \, dg \quad \text{or} \quad \lim_{\substack{a' \downarrow a \\ b' \uparrow b}} \int_{a'}^{b'} f \, dg.$$

**Remark B.3.4.** When an endpoint is infinite, it is clear that the integral should be improper at that endpoint. However, when the endpoint is finite, we need to be careful as the proper and the underlying improper integral can both exist, but they can be different.

Unless specified otherwise, we the Riemann–Stieltjes integral will be consider proper at a certain endpoint if that endpoint is finite and both the integrand and the integrator are defined there.

**Proposition B.3.5** ([34], Theorem 2.16). *The following is true for proper as well as improper Lebesgue–Stieltjes integrals, provided that all improper integrals in the same formula are considered in the same sense (e. g., all improper at the left endpoint):*

- (1) If  $c \in \mathbb{R}$  and  $\int_a^b f \, dg$  exists, so do  $\int_a^b (cf) \, dg$  and  $\int_a^b f \, d(CG)$  and we have

$$\int_a^b (cf) \, dg = \int_a^b f \, d(CG) = c \int_a^b f \, dg.$$

(2) If  $\int_a^b f_1 dg$  and  $\int_a^b f_2 dg$  exist, so does  $\int_a^b (f_1 + f_2) dg$  and we have

$$\int_a^b (f_1 + f_2) dg = \int_a^b f_1 dg + \int_a^b f_2 dg.$$

(3) If  $\int_a^b f dg_1$  and  $\int_a^b f dg_2$  exist, so does  $\int_a^b f d(g_1 + g_2)$  and we have

$$\int_a^b f d(g_1 + g_2) = \int_a^b f dg_1 + \int_a^b f dg_2.$$

□

**Proposition B.3.6** ([34], Theorem 2.21). *If the proper integral  $\int_a^b f dg$  exists, then so does  $\int_a^b g df$  and we have*

$$\int_a^b f dg = f(b)g(b) - f(a)g(a) - \int_a^b g df.$$

*For improper integrals, the statement remains true with  $f(a)g(a)$  and/or  $f(b)g(b)$  replaced by the corresponding limits, which are assumed to exist. Thus, if the integral  $\int_a^b f dg$  exists in the improper sense at the right endpoint and  $\lim_{b' \uparrow b} f(b')g(b')$  exists, then the integral  $\int_a^b g df$  also exists in the improper sense at the right endpoint and we have*

$$\int_a^b f dg = \lim_{b' \uparrow b} f(b')g(b') - f(a)g(a) - \int_a^b g df.$$

□

**Proposition B.3.7** ([34], Chapter 2, Exercise 16). *The Riemann–Stieltjes integral  $\int_a^b f dg$  exists if  $f$  is continuous and  $g$  is of bounded variation on  $[a, b]$ .* □

**Proposition B.3.8.** *If  $f$  has bounded variation and  $g$  is absolutely continuous on  $[a, b]$ , then*

$$\int_a^b f dg = \int_a^b f(x) g'(x) dx. \quad (\text{B.3.1})$$

*Replacing the closed interval  $[a, b]$  with an open or half-open one (in this case, the endpoint which is not included can be infinite), the result remains true in the improper sense: if  $f$  has bounded total variation on all closed subintervals,  $g$  is absolutely continuous and the right hand side exists as an improper integral, so does the left hand side and they agree.* □

**PROOF.** We shall only prove the result for closed intervals; the extension to open and semi-open intervals is left to the reader as an exercise.

First, we show that both sides of (B.3.1) exist. The right hand side exists because  $f$  is bounded (Remark B.2.2) and  $g' \in L^1([a, b])$  (Theorem B.1.7). Next, since  $g$  is continuous and  $f$  has bounded variation,  $\int_a^b g df$  exists by Proposition B.3.7. The integral  $\int_a^b f dg$  then exists by Proposition B.3.6.



To prove that both sides in (B.3.8) agree, take an array of points  $a = x_0^{(n)} \leq x_1^{(n)} \leq \dots \leq x_n^{(n)} = b$ ,  $n \in \mathbb{N}$ , such that  $\lim_{n \rightarrow \infty} \max_{1 \leq k \leq n} (x_k - x_{k-1}) = 0$ , e. g.,  $x_k^{(n)} = a + \frac{k}{n}(b - a)$ . Write

$$\int_a^b f(x) g(x) dx = A_n + B_n,$$

where

$$A_n := \sum_{k=1}^n f(x_k) \int_{x_{k-1}}^{x_k} g'(x) dx = \sum_{k=1}^n f(x_k) (g(x_k) - g(x_{k-1})),$$

$$B_n := \sum_{k=1}^n \int_{x_{k-1}}^{x_k} (f(x) - f(x_k)) g'(x) dx = \int_a^b \delta_n(x) g'(x) dx$$

and  $\delta_n(x) := f(x) - f(x_k)$  for  $x_{k-1} < x \leq x_k$  and  $\delta_n(a) := 0$ . Clearly,  $\lim_{n \rightarrow \infty} A_n = \int_a^b f dg$ . Next, since  $f$  has bounded total variation, it has at most countably many discontinuities. However, if  $f$  is continuous at  $x$ , we have  $\lim_{n \rightarrow \infty} \delta_n(x) = 0$ . Thus, the latter holds for almost all  $x \in [a, b]$ . Since  $f$  is bounded and  $g' \in L^1([a, b])$ , we can apply the dominated convergence theorem to deduce that  $\lim_{n \rightarrow \infty} B_n = 0$ , completing the proof.  $\square$

## B.4 Signed measures

**Definition B.4.1.** A *signed measure* on a measurable space  $(S, \mathcal{S})$  is a function  $\mu: \mathcal{S} \rightarrow \mathbb{R}$ , such that for any sequence  $A_1, A_2, \dots$  of pairwise disjoint sets, we have

$$\mu\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{k=1}^{\infty} \mu(A_k). \quad (\text{B.4.1})$$

**Remark B.4.2.** Here, we only deal with finite signed measures. In general, it is also possible to consider signed measures with possible infinite values.

**Remark B.4.3.** As the sum in the right hand side of (B.4.1) is independent of the order of the summands, the series must converge absolutely.

**Definition B.4.4.** The *positive* and *negative part* of a signed measure  $\mu$  are set functions  $\mu^+$  and  $\mu^-$  defined by

$$\mu^+(A) := \sup\{\mu(B) ; B \subseteq A, B \in \mathcal{S}\}, \quad \mu^-(A) := \sup\{-\mu(B) ; B \subseteq A, B \in \mathcal{S}\}. \quad (\text{B.4.2})$$

The following result can be deduced from Corollary 4.1.6, Proposition 4.1.7 and Exercise 6 of Section 4.1 in Cohn [13]:

**Proposition B.4.5.** *The functions  $\mu^+$  and  $\mu^-$  defined as above are finite positive measures and we have  $\mu = \mu^+ - \mu^-$ .*  $\square$

The representation  $\mu = \mu^+ - \mu^-$  is called the *Jordan decomposition* of  $\mu$ .

**Definition B.4.6.** The *variation* of a signed measure  $\mu$  defined on  $(S, \mathcal{S})$  is the positive measure  $|\mu| := \mu^+ + \mu^-$ . The *total variation* of the  $\mu$  is defined as  $\|\mu\| := |\mu|(S)$ .

**Definition B.4.7.** For a function  $f \in L^1(|\mu|)$ , define the *Lebesgue integral* of  $f$  with respect to  $\mu$  as

$$\int f \, d\mu := \int f \, d\mu_+ - \int f \, d\mu^-,$$

where  $\mu = \mu^+ - \mu^-$  is called the Jordan decomposition of  $\mu$ .

**Remark B.4.8.** We have  $\left| \int f \, d\mu \right| \leq \int |f| \, d|\mu|$ .

**Definition B.4.9.** Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be a right-continuous function with bounded variation. The *Lebesgue–Stieltjes measure* associated to  $g$  is the signed measure  $\Lambda_g$  determined by

$$\Lambda_g((a, b]) = g(b) - g(a) \tag{B.4.3}$$

for all  $a \leq b$ .

**Remark B.4.10.** By Dynkin's  $\pi$ - $\lambda$  theorem (see, e. g., Theorem 1.6.2 of Cohn [13]), (B.4.3) determines  $\Lambda_g$  uniquely.

**Proposition B.4.11.** *Let  $g$  be a right-continuous function  $g$  with bounded variation.*

- (1) *The Lebesgue–Stieltjes measure  $\Lambda_g$  exists and we have  $\|\Lambda_g\| = V(g)$ .*
- (2) *For  $a \leq b$  and a continuous function  $f: [a, b] \rightarrow \mathbb{R}$ , the Riemann–Stieltjes integral  $\int_a^b f \, dg$  exists and we have*

$$\int_{(a,b]} f \, d\Lambda_g = \int_a^b f \, dg.$$

PROOF. By Theorem B.2.4, we may write  $g = h - k$ , where  $g$  and  $h$  are monotone increasing functions which are left/right-continuous at each point where so is  $f$ . Moreover,  $V(g) = V(h) + V(k)$ . Thus,  $h$  and  $k$  are right-continuous. By Theorem 11.10 of Wheeden and Zygmund [34], there exist Lebesgue–Stieltjes measures  $\Lambda_h$  and  $\Lambda_k$ , which are finite positive measures. From (B.4.3) and Dynkin's  $\pi$ - $\lambda$  theorem, it follows that  $\Lambda_g = \Lambda_h - \Lambda_k$ .

By Proposition B.2.3, we have  $V(h) = \lim_{n \rightarrow \infty} h(n) - h(-n) = \lim_{n \rightarrow \infty} \Lambda_h((-n, n]) = \Lambda_h(\mathbb{R})$  and similarly  $V(k) = \Lambda_k(\mathbb{R})$ . For each Borel set  $B \subseteq \mathbb{R}$ , we have  $\Lambda_g(B) \leq \Lambda_h(B)$  and  $-\Lambda_g(B) \leq \Lambda_k(B)$ . Taking the supremum over  $B$  and recalling (B.4.2), we obtain  $\Lambda_g^+(\mathbb{R}) \leq \Lambda_h(\mathbb{R}) = V(h)$  and  $\Lambda_g^-(\mathbb{R}) \leq \Lambda_k(\mathbb{R}) = V(k)$ . Therefore,  $\|\Lambda_g\| = \Lambda_g^+(\mathbb{R}) + \Lambda_g^-(\mathbb{R}) \leq V(h) + V(k) = V(g)$ .

For each  $\varepsilon > 0$ , there exists a sequence  $x_0 \leq x_1 \leq \dots \leq x_n$ , such that  $\sum_{k=1}^n |g(x_k) - g(x_{k-1})| > V(g) - \varepsilon$ . Now observe that  $\sum_{k=1}^n |\Lambda_g((x_{k-1}, x_k])| \leq \sum_{k=1}^n |\Lambda_g^+((x_{k-1}, x_k]) - \Lambda_g^-((x_{k-1}, x_k])| \leq \sum_{k=1}^n (\Lambda_g^+((x_{k-1}, x_k]) + \Lambda_g^-((x_{k-1}, x_k])) = \Lambda_g^+((x_0, x_n]) + \Lambda_g^-((x_0, x_n]) \leq \Lambda_g^+(\mathbb{R}) + \Lambda_g^-(\mathbb{R}) = \|\Lambda_g\|$ . Therefore,  $\|\Lambda_g\| > V(g) - \varepsilon$  for all  $\varepsilon > 0$ . As a result, we have  $\|\Lambda_g\| \geq V(g)$ . Combining with the preceding paragraph, we find that  $\|\Lambda_g\| = V(g)$ , proving part (1)

By Proposition B.3.7, the Riemann–Stieltjes integrals  $\int f dh$  and  $\int f dk$  both exist. By Proposition B.3.5, so does  $\int f dg$  and we have  $\int f dg = \int f dh - \int f dk$ . By Theorem 11.11 of Wheeden and Zygmund [34], we have  $\int_{(a,b]} f d\Lambda_h = \int_a^b f dh$  and  $\int_{(a,b]} f d\Lambda_k = \int_a^b f dk$ . Since  $\Lambda_g = \Lambda_h - \Lambda_k$ , we conclude that  $\int_{(a,b]} f d\Lambda_g = \int_a^b f dg$ , proving part (2).  $\square$

# Appendix C

## On the Mills ratio

### C.1 Basic properties

**Definition C.1.1.** Let  $\phi$  denote the standard Gaussian density on  $\mathbb{R}$ , that is,

$$\phi(w) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}w^2} \quad (\text{C.1.1})$$

and let  $\Phi$  denote its cumulative distribution function:

$$\Phi(w) := \int_{-\infty}^w \phi(x) \, dx. \quad (\text{C.1.2})$$

The *Mills ratio* is the ratio between these two functions as follows:

$$\psi(w) := \frac{\Phi(w)}{\phi(w)} = e^{\frac{1}{2}w^2} \int_{-\infty}^w e^{-\frac{1}{2}x^2} \, dx. \quad (\text{C.1.3})$$

**Remark C.1.2.** In the literature, the Mills ratio is often defined slightly differently, as the function  $w \mapsto \psi(-w)$  (in our notation). The benefit of that definition is that it gives a ‘tame’ function for positive  $w$ . However, Definition (C.1.3) also has its benefits.

The Mills ratio is important because it can serve to express solutions to the Stein equation. As one can easily check, the Mills ratio solves the equation

$$\psi'(w) = w\psi(w) + 1, \quad (\text{C.1.4})$$

which is a version of the Stein equation (3.2.1). Repeated differentiation gives further derivatives, some of which are listed below:

$$\psi''(w) = (w^2 + 1)\psi(w) + w, \quad (\text{C.1.5})$$

$$\psi'''(w) = (w^3 + 3w)\psi(w) + w^2 + 2, \quad (\text{C.1.6})$$

$$\psi^{(4)}(w) = (w^4 + 6w^2 + 3)\psi(w) + w^3 + 5w. \quad (\text{C.1.7})$$

The derivatives of the Mills ratio also satisfy a recurrent formula as stated below:

**Proposition C.1.3.** *For all  $r \in \mathbb{N}$  and all  $w \in \mathbb{R}$ , we have*

$$\psi^{(r+1)}(w) = w \psi^{(r)}(w) + r \psi^{(r-1)}(w). \quad (\text{C.1.8})$$

PROOF. For  $r = 1$ , the identity is immediate from (C.1.4) and (C.1.5). The induction step from  $r$  to  $r + 1$  can be performed by differentiating (C.1.8), leading to

$$\psi^{(r+2)}(w) = w \psi^{(r+1)}(w) + (r + 1) \psi^{(r)}(w), \quad (\text{C.1.9})$$

which is exactly (C.1.8) with  $r + 1$  in place of  $r$ .  $\square$

**Corollary C.1.4.** *The  $r$ -th derivative of the Mills ratio can be expressed as*

$$\psi^{(r)}(w) = P_r(w) \psi(w) + Q_r(w),$$

where  $P_r$  is a polynomial of degree  $r$  and  $Q_r$  is a polynomial of degree  $r - 1$ .  $\square$

**Proposition C.1.5.** *The Mills ratio and all its derivatives are strictly positive functions, i. e.,  $\psi^{(r)}(w) > 0$  for all  $r = 0, 1, 2, \dots$  and all  $w \in \mathbb{R}$ .*

PROOF. Introducing a new variable  $t := w - x$  into the integral in (C.1.3), we obtain

$$\psi(w) = e^{\frac{1}{2}w^2} \int_0^\infty e^{-\frac{1}{2}(w-t)^2} dt = \int_0^\infty e^{tw - \frac{1}{2}t^2} dt.$$

Repeated differentiation under the integral sign gives

$$\psi^{(r)}(w) = \int_0^\infty t^r e^{tw - \frac{1}{2}t^2} dt > 0 \quad (\text{C.1.10})$$

(one can easily check that this works because all these integrals converge uniformly and absolutely).  $\square$

Combining (C.1.4) and (C.1.5) with Proposition C.1.5 (which, among others, implies that  $\psi$  is strictly increasing), we obtain the following upper and lower bound on  $\psi$  for negative arguments: for all  $w > 0$ , we have

$$\frac{w}{w^2 + 1} < \psi(-w) < \min \left\{ \sqrt{\frac{\pi}{2}}, \frac{1}{w} \right\}.$$

Moreover, derivatives of  $\psi$  can be bounded similarly.

**Proposition C.1.6.** *For all  $r \in \mathbb{N}_0$  and all  $w > 0$ , we have*

$$\psi^{(r)}(-w) < \frac{r!}{w^{r+1}}. \quad (\text{C.1.11})$$

PROOF. From (C.1.10), we obtain

$$\psi^{(r)}(-w) = \int_0^\infty t^r e^{-tw - \frac{1}{2}t^2} dt < \int_0^\infty t^r e^{-tw} dt = \frac{r!}{w^{r+1}}, \quad (\text{C.1.12})$$

completing the proof.  $\square$

## C.2 Repeated integrals of the Gaussian density

There is a strong relationship between the derivatives of the Mills ratio and the repeated integrals of the standard Gaussian density  $\phi$ . Inductively, define

$$\Phi_0(w) := \phi(w), \quad \Phi_{r+1}(w) := \int_{-\infty}^w \Phi_r(x) dx. \quad (\text{C.2.1})$$

**Proposition C.2.1.** *The functions  $\Phi_r$  are all well defined (i. e., all integrals in (C.2.1) converge and we have*

$$\Phi_{r+1}(w) = \frac{1}{r!} \phi(w) \psi^{(r)}(w) \quad (\text{C.2.2})$$

for all  $r \in \mathbb{N}_0$ .

PROOF. For  $r = 0$ , (C.2.2) is immediate from the definition of the Mills ratio. Now make the induction step from  $r - 1$  to  $r$ . Thus, assume that

$$\Phi_r(w) = \frac{1}{(r-1)!} \phi(w) \psi^{(r-1)}(w). \quad (\text{C.2.3})$$

First, it follows from (C.1.12) that the function  $\Phi_r$  is integrable on all intervals  $(-\infty, w]$  and that (C.2.2) holds in the limit as  $w \rightarrow -\infty$ . Therefore, it suffices to prove the derivative of (C.2.2). By the induction hypothesis (C.2.3), we have

$$\begin{aligned} \frac{d}{dw} \left[ \Phi_{r+1}(w) - \frac{1}{r!} \phi(w) \psi^{(r)}(w) \right] &= \Phi_r(w) - \frac{1}{r!} \phi'(w) \psi^{(r)}(w) - \frac{1}{r!} \phi(w) \psi^{(r+1)}(w) \\ &= \frac{1}{r!} \phi(w) \left[ r \psi^{(r-1)}(w) + w \psi^{(r)}(w) - \psi^{(r+1)}(w) \right]. \end{aligned} \quad (\text{C.2.4})$$

However, by the recurrent formula (C.1.8), the preceding expression equals zero. This completes the proof.  $\square$

**Corollary C.2.2.** *For each  $r \in \mathbb{N}$ , there exists  $C_r$ , such that*

$$\Phi^{(r)}(-w) \leq C_r e^{-\frac{1}{2}w^2}$$

for all  $w \geq 0$ .

PROOF. By (C.2.2) and since  $\psi^{(r)}$  is increasing by Proposition C.1.5, we can estimate  $\Phi^{(r)} = \frac{1}{(r-1)!} \phi(w) \psi^{(r-1)}(-w) \leq \frac{1}{(r-1)!} \phi(w) \psi^{(r-1)}(0) = \frac{\psi^{(r-1)}(0)}{(r-1)! \sqrt{2\pi}} e^{-\frac{1}{2}w^2}$ .  $\square$

**Proposition C.2.3.** *For all  $r \in \mathbb{N}_0$  and all  $w \in \mathbb{R}$ , we have*

$$\Phi_r(-w) \psi^{(r)}(w) + \Phi_r(w) \psi^{(r)}(-w) = 1. \quad (\text{C.2.5})$$

PROOF. For  $r = 0$ , (C.2.5) reduces to the obvious identity  $\Phi(w) + \Phi(-w) = 1$ . For  $r \in \mathbb{N}$ , we apply Proposition C.2.1, which reduces (C.2.5) to

$$J_r(w) := \phi(w) \left[ \psi^{(r-1)}(-w) \psi^{(r)}(w) + \psi^{(r-1)}(w) \psi^{(r)}(-w) \right] = (r-1)!, \quad (\text{C.2.6})$$

and use induction over  $r$ . For  $r = 1$ , (C.1.4) implies

$$\begin{aligned}
 J_1(w) &= \phi(w) \left[ \psi(-w) \psi'(w) + \psi(w) \psi'(-w) \right] = \\
 &= \phi(w) \left[ \psi(-w) (w \psi(w) + 1) + \psi(w) (-w \psi(-w) + 1) \right] = \\
 &= \phi(w) (\psi(-w) + \psi(w)) = \\
 &= 1.
 \end{aligned} \tag{C.2.7}$$

Now perform the induction step from  $r$  to  $r + 1$ . This time, application of the recurrent formula (C.1.8) (instead of (C.1.4)) gives

$$\begin{aligned}
 J_{r+1}(w) &= \phi(w) \left[ \psi^{(r)}(-w) \psi^{(r+1)}(w) + \psi^{(r)}(w) \psi^{(r+1)}(-w) \right] = \\
 &= \phi(w) \left[ \psi^{(r)}(-w) (w \psi^{(r)}(w) + r \psi^{(r-1)}(w)) + \right. \\
 &\quad \left. + \psi^{(r)}(w) (-w \psi^{(r)}(-w) + r \psi^{(r-1)}(-w)) \right] = \\
 &= r J_r(w) = \\
 &= r!.
 \end{aligned} \tag{C.2.8}$$

This completes the proof. □

# Bibliography

- [1] R. ARRATIA, L. GOLDSTEIN, AND L. GORDON, *Poisson approximation and the Chen-Stein method*, *Statist. Sci.*, 5 (1990), pp. 403–434.
- [2] A. D. BARBOUR, *Stein's method*, in *Encyclopedia of Statistical Sciences*, vol. 1, Wiley, New York, 1997, pp. 513–521.
- [3] A. D. BARBOUR AND L. H. Y. CHEN, *An Introduction to Stein's Method*, vol. 4 of *Lecture Notes Series*, Institute for Mathematical Sciences, National University of Singapore, 2005.
- [4] —, *Stein's Method and Applications*, vol. 5 of *Lecture Notes Series*, Institute for Mathematical Sciences, National University of Singapore, 2005.
- [5] A. D. BARBOUR, L. HOLST, AND S. JANSON, *Poisson Approximation*, vol. 2 of *Oxford Studies in Probability*, Oxford University Press, New York, 1992.
- [6] A. D. BARBOUR, M. KAROŃSKI, AND A. RUCIŃSKI, *A central limit theorem for decomposable random variables with applications to random graphs*, *J. Combin. Theory Ser. B*, 47 (1989), pp. 125–145.
- [7] A. C. BERRY, *The accuracy of the Gaussian approximation to the sum of independent variates*, *Trans. Amer. Math. Soc.*, 49 (1941), pp. 122–136.
- [8] E. BOLTHAUSEN, *An estimate of the remainder in a combinatorial central limit theorem*, *Z. Wahrsch. verw. Gebiete*, 66 (1984), pp. 379–386.
- [9] L. H. Y. CHEN, *Poisson approximation for dependent trials*, *Ann. Probab.*, 3 (1975), pp. 534–545.
- [10] L. H. Y. CHEN, L. GOLDSTEIN, AND A. RÖLLIN, *Stein's method via induction*. ArXiv:1903.09319, 2019.
- [11] L. H. Y. CHEN AND Q.-M. SHAO, *A non-uniform Berry-Esseen bound via Stein's method*, *Probab. Theory Related Fields*, 120 (2001), pp. 236–254.
- [12] —, *Normal approximation under local dependence*, *Ann. Probab.*, 32 (2004), pp. 1985–2028.



- [13] D. L. COHN, *Measure theory*, Birkhäuser Advanced Texts: Basler Lehrbücher. [Birkhäuser Advanced Texts: Basel Textbooks], Birkhäuser/Springer, New York, second ed., 2013.
- [14] P. R. DE MONTMORT, *Essay de l'analyse sur les jeux de hazard*, Jacques Quillau, Paris, second ed., 1713.
- [15] A. DEMBO AND Y. RINOTT, *Some examples of normal approximations by Stein's method*, in Random discrete structures (Minneapolis, MN, 1993), vol. 76 of IMA Vol. Math. Appl., Springer, New York, 1996, pp. 25–44.
- [16] C.-G. ESSEEN, *On the Liapounoff limit of error in the theory of probability*, Ark. Mat. Astr. Fys., 28A (1942), p. 19.
- [17] C. G. ESSEEN, *A moment inequality with an application to the central limit theorem*, Skand. Aktuarietidskr., 39 (1956), pp. 160–170 (1957).
- [18] S. N. ETHIER AND T. G. KURTZ, *Markov Processes*, Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics, John Wiley & Sons Inc., New York, 1986.
- [19] L. GOLDSTEIN, *Berry–Esseen bounds for combinatorial central limit theorems and pattern occurrences, using zero and size biasing*, J. Appl. Probab., 42 (2005), pp. 661–683.
- [20] C. HEIL, *Introduction to Real Analysis*, vol. 280 of Graduate Texts in Mathematics, Springer, Cham, 2019.
- [21] P. L. HSU, *The approximate distributions of the mean and variance of a sample of independent variables*, Ann. Math. Statistics, 16 (1945), pp. 1–29.
- [22] V. Y. KOROLEV AND I. G. SHEVTSOVA, *An upper bound for the absolute constant in the Berry–Esseen inequality (Russian)*, Teor. Veroyatn. Primen., 54 (2009), pp. 671–695.
- [23] L. LE CAM, *An approximation theorem for the Poisson binomial distribution*, Pacific J. Math., 10 (1960), pp. 1181–1197.
- [24] M. MANETTI, *Topology*, vol. 91 of Unitext, Springer, Cham, 2015. Translated from the 2014 Italian edition by Simon G. Chiossi, La Matematica per il 3+2.
- [25] S. T. RACHEV, *The Monge–Kantorovich problem on mass transfer and its applications in stochastics*, Theory Probab. Appl., 29 (1984), pp. 647–676.
- [26] M. RAIČ, *CLT-related large deviation bounds based on Stein's method*, Adv. in Appl. Probab., 39 (2007), pp. 731–752.
- [27] L. C. G. ROGERS AND D. WILLIAMS, *Diffusions, Markov Processes, and Martingales. Vol. 1*, Cambridge Mathematical Library, Cambridge University Press, Cambridge, 2000.

- [28] I. G. SHEVTSOVA, *On the absolute constants in the Berry-Esseen type inequalities for identically distributed summands*. ArXiv:1111.6554, 2011.
- [29] —, *On the absolute constants in the Berry–Esseen inequality and its structural and nonuniform improvements (Russian)*, Inform. Primen., 7 (2013), pp. 124–125.
- [30] T. B. SINGH, *Introduction to Topology*, Springer, Singapore, 2019.
- [31] C. STEIN, *A bound for the error in the normal approximation to the distribution of a sum of dependent random variables*, in Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability (Univ. California, Berkeley, Calif., 1970/1971), Vol. II: Probability theory, Berkeley, Calif., 1972, Univ. California Press, pp. 583–602.
- [32] —, *Approximate Computation of Expectations*, Institute of Mathematical Statistics Lecture Notes—Monograph Series, 7, Institute of Mathematical Statistics, Hayward, CA, 1986.
- [33] T. TAO, *An Introduction to Measure Theory*, vol. 126 of Graduate Studies in Mathematics, American Mathematical Society, Providence, RI, 2011.
- [34] R. L. WHEEDEN AND A. ZYGMUND, *Measure and Integral*, Pure and Applied Mathematics (Boca Raton), CRC Press, Boca Raton, FL, second ed., 2015. An introduction to real analysis.