

Approximating the Conway–Maxwell–Poisson distribution normalization constant

Steven B. Gillispie * and Christopher G. Green
University of Washington, Seattle, USA

—
Department of Statistics
Technical Report no. 615

June 6, 2013

Abstract

By adding a second parameter, Conway and Maxwell created a new distribution for situations where data deviate from the standard Poisson distribution. This new distribution contains a normalization constant expressed as an infinite sum whose summation has no known closed-form expression. Shmueli *et al.* produced an approximation for this sum but proved only that it was valid for integer values of the second parameter, although they conjectured it was also valid for non-integers. Here we prove their conjecture to be true and discuss for what range of parameters the approximation can be accurately applied.

Keywords: Approximation; Conway–Maxwell–Poisson distribution; Normalization constant; Overdispersion; Underdispersion

**Address for correspondence:* Steven B. Gillispie, Department of Radiology, Box 357987, University of Washington, Seattle, Washington, 98195-7987, USA.

E-mail: gillisp@u.washington.edu

1 Introduction

Faced with a problem of queuing systems with state-dependent service rates, Conway and Maxwell (1962) introduced (by adding a second parameter ν) a discrete distribution that extended the Poisson distribution to better model situations where the data are either overdispersed or underdispersed relative to the Poisson distribution. This Conway–Maxwell–Poisson (CMP) distribution has apparently been little used in the last fifty years, but a recent attempt was made by Shmueli *et al.* (2005) to revive it. The CMP distribution can be described quite simply as (we use here the notation of Shmueli *et al.* instead of Conway and Maxwell):

$$P(X = x) = \frac{\lambda^x}{(x!)^\nu} \frac{1}{Z(\lambda, \nu)} \quad x = 0, 1, 2, \dots \quad (1)$$

where $Z(\lambda, \nu)$ is the normalization constant

$$Z(\lambda, \nu) = \sum_{j=0}^{\infty} \frac{\lambda^j}{(j!)^\nu}, \quad (2)$$

valid for $\nu \geq 0$ and most $\lambda > 0$. (The exception is discussed in Section 3.)

Shmueli *et al.* discuss approximating $Z(\lambda, \nu)$ and $1/Z(\lambda, \nu)$ via truncating this sum in their Appendix B, as well as the functional approximation

$$Z(\lambda, \nu) = \frac{\exp(\nu\lambda^{1/\nu})}{\lambda^{(\nu-1)/2\nu} (2\pi)^{(\nu-1)/2} \sqrt{\nu}} \left\{ 1 + O\left(\frac{1}{\lambda^{1/\nu}}\right) \right\}, \quad (3)$$

proven for all fixed *integer* $\nu \geq 1$ in the limit as $\lambda \rightarrow \infty$ (equivalently, $\lambda^{1/\nu} \rightarrow \infty$). They derive this approximation by converting the summation to a multiple integral, one integration per value of the integer ν , which ultimately can be reduced to a single complex contour integral. The approximation in Equation 3 can then be produced using Laplace’s method. Unfortunately this idea only works for integer ν , though the authors show using computed data that the approximation appears to be good for all ν . In our discussion here we prove that this conjecture is indeed true.

As it turns out, this exact summation was given as an example in asymptotic expansions of entire functions in a textbook on asymptotic approximation by Olver (1974), but only for $\nu \leq 4$. The technique shown there is to substitute the sum with a contour integral around the non-negative integers

where the proportional error to the approximation, represented by the ratio of the integral to the sum, goes to zero. But the proof is only valid for $\nu \leq 4$ because the contour integral is controlled by the integral of $1/\Gamma(z)$ along the imaginary axis where its growth rate is dominated by $\exp(\pi\nu y/2)$. This is multiplied by a term of approximately $\exp(-2\pi y)$, so that the integral

$$\int_{1/2}^{\infty} e^{\pi\nu y/2} e^{-2\pi y} dy = \int_{1/2}^{\infty} e^{\pi(\nu-4)y/2} dy \quad (4)$$

cannot converge unless $\nu \leq 4$.

For $\nu = 0$ the summation (2) only converges if $\lambda < 1$, but then the approximation is not needed as the sum is merely $1/(1 - \lambda)$. While Olver showed that the approximation is good for $0 < \nu \leq 4$, we prove here that it is valid for large λ for all $\nu > 0$. We do this in the following section, using a method very similar to that of Olver (which simultaneously gives more details on the techniques Olver used). The difference with our method is that we use a different contour to estimate the error integral. While the result of this error integral still tends to infinity for large λ given any fixed ν , the summation function $Z(\lambda, \nu)$ grows even larger; therefore the ratio of the error term to the approximation tends to zero. It is this result that confirms the conjecture by Shmueli *et al.* that their approximation was good for all $\nu > 0$, as supported by their computer-derived results.

2 Proof

Our method closely follows that of Olver, though with a different contour for the integral. We begin by defining the n th term approximation to $Z(\lambda, \nu)$:

$$Z_n(\lambda, \nu) = \sum_{j=0}^n \frac{\lambda^j}{(j!)^\nu} = \sum_{j=0}^n \frac{\lambda^j}{(\Gamma(j+1))^\nu} = \sum_{j=0}^n \frac{\lambda^j}{\exp(\nu \log \Gamma(j+1))}. \quad (5)$$

Since $\Gamma(z+1)$ is non-zero everywhere for $\text{Re } z > -1$, setting the branch point for the complex logarithm at $z = -1$ with its branch cut extending along the negative real axis to $-\infty$ means that the summand is analytic for any contour in the half-plane $\text{Re } z > -1$. Therefore, as long as we keep our contours in this half-plane, we can simplify the notation and just refer to $(\Gamma(j+1))^\nu$ without need of the logarithm. By the Cauchy residue theorem,

$$\oint_{z=j} f(z) \cot(\pi z) dz = \oint_{z=j} \frac{f(z)}{\pi} \cos(\pi z) \frac{\pi(z-j)}{\sin(\pi z)} \frac{1}{z-j} dz = 2i f(j). \quad (6)$$

Letting $f(z)$ be the analytic summand from (5), it is clear that we can combine the separate contour integrals in (6) around each of the non-negative integers j into a single contour. Let $\zeta = \lambda^{1/\nu}$ for notational simplicity. Then combining (5) and (6) produces

$$Z_n(\lambda, \nu) = \frac{1}{2i} \oint_C \left(\frac{\zeta^z}{\Gamma(z+1)} \right)^\nu \cot(\pi z) dz. \quad (7)$$

Our contour involves an angle β defined below but the description of the contour is simplified by first creating a dependent value $\delta \equiv (n/2+1/4) \tan \beta$. Then our contour C (see Figure 1) is described as starting from the positive real axis at $z = n + 1/2$ and first proceeding vertically to the point $z = n+1/2+i\delta$. This section of the contour is labeled C_a . From there the contour runs parallel to the real axis back to the point $z = (n+1/2)/2 + i\delta = x_p + i\delta$ where $x_p \equiv (n+1/2)/2$. This section is labeled C_b . Next, the contour proceeds toward the origin at an angle β to a point $z = x_d(1+i \tan \beta)$ that is at a distance of $1/2$ from the origin along this ray; this section is labeled C_c . (So $x_d = (\cos \beta)/2$.) The section C_d is a semi-circular arc of radius $1/2$ around the origin to the point $z = x_d(1-i \tan \beta)$. From there sections C_e , C_f , and C_g are mirror reflections of C_c , C_b and C_a , respectively, moving down, across, and back up to the starting point at $z = n + 1/2$ and bending at the points $z = x_p - i\delta$ and $z = n + 1/2 - i\delta$. We also refer to the half of the contour C above the real axis as C_1 and C_2 as the half below it. Our goal will be to eventually let $n \rightarrow \infty$ so that $Z_n(\lambda, \nu) \rightarrow Z(\lambda, \nu)$.

Starting with the Euler relation $e^{i\pi z} = \cos(\pi z) + i \sin(\pi z)$, then dividing by $2i \sin(\pi z) = e^{i\pi z} - e^{-i\pi z}$, we derive

$$\frac{1}{2i} \cot(\pi z) = -\frac{1}{2} - \frac{1}{e^{-i2\pi z} - 1}. \quad (8)$$

Substituting this into (7) for C_1 , and its alternative form after applying $\cot(-\pi z) = -\cot(\pi z)$ for C_2 , equates Z_n to a sum of four integrals. The contours of the two integrals without the exponential in the denominator can be completed by adding and subtracting the same integrals but now solely along the real axis between the two points where C crosses the real line. But since the integrands of these two integrals no longer have any singularities their now-closed contour integrals evaluate to zero. This leaves the two newly added integrals along the real line, which sum together, and

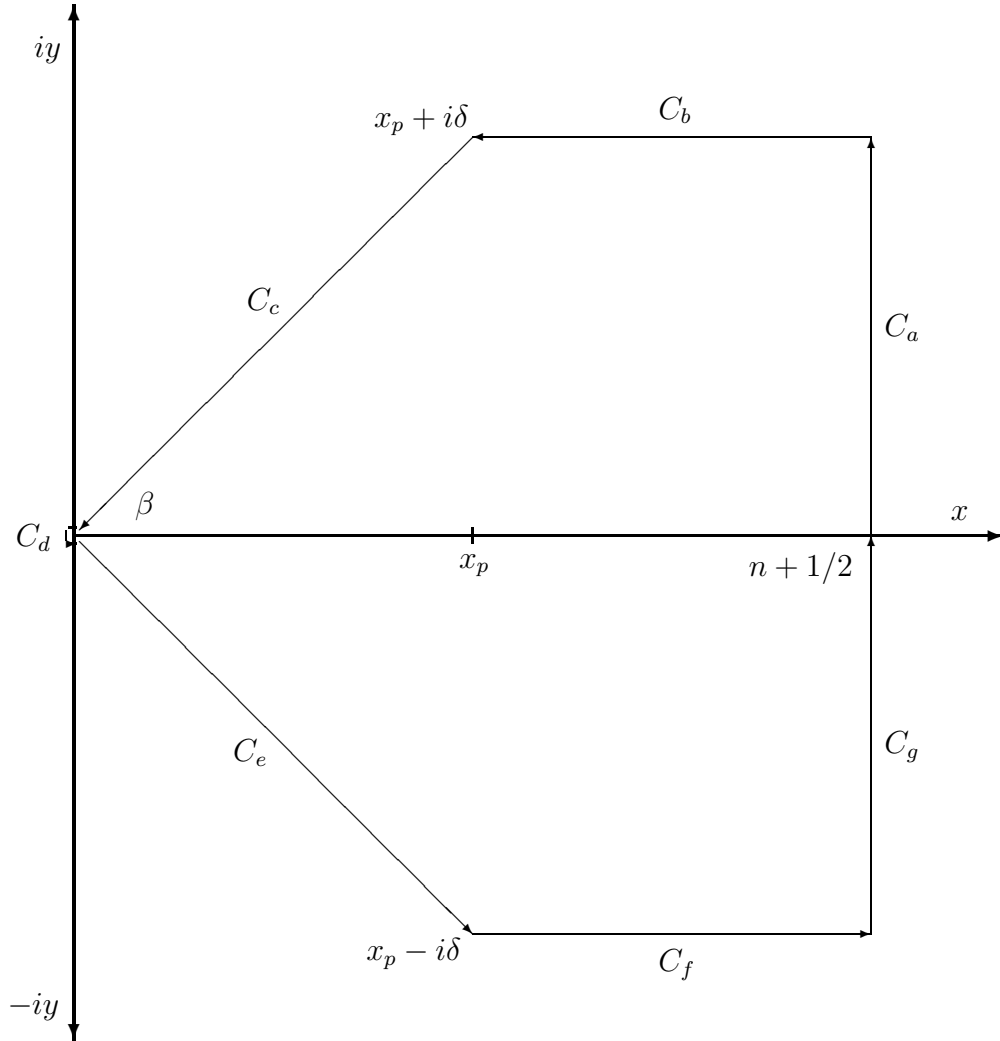


Figure 1: The contour used in (7) and subsequent integrals. Here $\beta = \pi/4$.

the two remaining contour integrals. The result is

$$Z_n(\lambda, \nu) = \int_{-1/2}^{n+1/2} \left(\frac{\zeta^x}{\Gamma(x+1)} \right)^\nu dx \quad (9)$$

$$- \oint_{C_1} \left(\frac{\zeta^z}{\Gamma(z+1)} \right)^\nu \frac{1}{e^{-i2\pi z} - 1} dz + \oint_{C_2} \left(\frac{\zeta^z}{\Gamma(z+1)} \right)^\nu \frac{1}{e^{i2\pi z} - 1} dz \quad (10)$$

where line (9) represents the approximating function and line (10) represents the error.

We now evaluate the integral in (9). We wish to use Laplace's method, which states that $\int_a^b \exp\{f(x)\}g(x) dx \approx \exp\{f(x_m)\} \sqrt{-2\pi/f''(x_m)} g(x_m)$ where x_m maximizes $f(x)$, but the integrand is not in the proper form. Nevertheless, because the gamma function has a minimum it must be that the integrand has a maximum. The only part of the integrand's derivative that could be zero is

$$\log \zeta - \psi(z+1) = \log \zeta - \psi(z) - \frac{1}{z} \approx \log \zeta - \log z - \frac{1}{2z} \quad (11)$$

where $\psi(z) = \Gamma'(z)/\Gamma(z)$ is the digamma function and its approximation is from Abramowitz and Stegun (1964). Thus the maximum of the integrand must occur at a point μ slightly less than ζ . We use Stirling's approximation of $\Gamma(z+1) = (z/e)^z \sqrt{2\pi z} \{1 + O(1/z)\}$ and also make the simplifying change of variable $t = x/\zeta - 1$. This gives the integrand the proper form for Laplace's method, where the function in the exponential is maximized at $t = \mu/\zeta - 1$. That is, if $a = -(1 + 1/2\zeta)$ and $b = (n + 1/2)/\zeta - 1$,

$$\frac{\zeta}{(2\pi\zeta)^{\nu/2}} \int_a^b \frac{e^{\zeta\nu(1+t)(1-\log(1+t))}}{(1+t)^{\nu/2}} dt = \frac{e^{\nu\mu}}{(2\pi\mu)^{(\nu-1)/2} \sqrt{\nu} \mu^2} \frac{\zeta^2}{\mu^2} \left\{ 1 + O\left(\frac{1}{\zeta}\right) \right\}. \quad (12)$$

For large ζ , the ratio $\zeta/\mu \rightarrow 1$ and the desired functional approximation (3) is the result.

The next step is to evaluate the contour integrals on line (10). We will be considering the magnitudes of the integrals so their signs and directions of integration no longer matter. The contour C_d has length less than 2π and the integrand on it is bounded; therefore the integral must also be bounded. We can see by observation that the result will be dominated by $A^\nu \zeta^{c\nu}$ for some constants A and c , but the approximation function in (3) is dominated

by e^{ζ^ν} so the ratio of the contour integral around C_d to (3) must go to zero for large enough ζ .

Along C_a (the vertical contour portion of C_1), $z = n + 1/2 + iy$ so that

$$\left| \frac{1}{e^{-i2\pi z} - 1} \right| = \left| \frac{1}{e^{-i2\pi(n+1/2+iy)} - 1} \right| = \frac{1}{e^{2\pi y} + 1} \leq e^{-2\pi y}. \quad (13)$$

Now let $z = |z|e^{i\theta}$ where $\tan \theta = y/(n + 1/2)$ and y varies from 0 to $\delta = x_p \tan \beta$. Then, remembering that $x_p = (n + 1/2)/2$,

$$\theta = \arctan \frac{y}{n + 1/2} \leq \arctan \frac{x_p \tan \beta}{n + 1/2} = \arctan \frac{\tan \beta}{2} \equiv \Theta \quad (14)$$

where Θ is a constant.

The integrand will also contain (after applying Stirling's formula) a term

$$\left| \left(\frac{1}{(n + 1/2 + iy)^{n+1+iy}} \right)^\nu \right| \leq \left(\frac{e^{\theta y}}{|n + 1/2 + iy|^{n+1}} \right)^\nu \leq \frac{e^{\Theta \nu y}}{n^{\nu(n+1)}} \quad (15)$$

using the result of (14). Then, applying the Stirling approximation and (13) to the contour integral $|I_a|$ along C_a where $|dz| = |dy|$,

$$\begin{aligned} |I_a| &\leq \int_0^\delta \left| \left(\frac{\zeta^{n+1/2+iy}}{\Gamma(n + 1/2 + iy + 1)} \right)^\nu \right| \frac{1}{e^{2\pi y} + 1} dy \\ &\leq \frac{(\zeta e)^{\nu(n+1/2)}}{(2\pi)^{\nu/2} n^{\nu(n+1)}} \int_0^\delta e^{(\Theta\nu - 2\pi)y} dy = \frac{(\zeta e)^{\nu(n+1/2)}}{(2\pi)^{\nu/2} n^{\nu(n+1)}} \frac{e^{(\Theta\nu - 2\pi)\delta} - 1}{\Theta\nu - 2\pi} \end{aligned} \quad (16)$$

and as $n \rightarrow \infty$ this goes to zero. The same result is produced along the C_g portion of C_2 .

Next, for the two integrals C_b and C_c in C_1 of line (10), we first consider the magnitude of the exponential portion of the integrands:

$$\sqrt{\frac{1}{(e^{-i2\pi z} - 1)(e^{i2\pi \bar{z}} - 1)}} = \sqrt{\frac{1}{e^{4\pi y} - 2e^{2\pi y} \cos(2\pi x) + 1}} \leq \frac{1}{|e^{2\pi y} - 1|}. \quad (17)$$

This result will also be true along C_2 after substituting $|y|$ for y .

Along the C_b portion of the C_1 contour, $z = x + i\delta$ and a similar reasoning as above along C_a , again with $\theta = \arg z$, shows that

$$\left| \frac{1}{(x + i\delta)^{x+1/2+i\delta}} \right| \leq \frac{e^{\theta\delta}}{|x + i\delta|^{x+1/2}} \leq \frac{e^{\beta\delta}}{x^{x+1/2}}, \quad (18)$$

where the final inequality will be true for all $x \in [x_p, \infty)$. Then, after applying (17), (18), and Stirling's formula,

$$\begin{aligned}
|I_b| &\leq \frac{e^{\beta\nu\delta}}{|e^{2\pi\delta} - 1|} \int_{x_p}^{n+1/2} \left(\frac{\zeta^x}{\sqrt{2\pi} e^{-x} x^{x+1/2}} \right)^\nu dx \\
&\approx (2\pi)^{-\nu/2} e^{(\beta\nu-2\pi)\delta} \int_{x_p}^{2x_p} \frac{\exp\{\nu(1 + \log \zeta - \log x)x\}}{x^{\nu/2}} dx \\
&\leq (2\pi)^{-\nu/2} e^{(\beta\nu-2\pi)\delta} \int_{x_p}^{\infty} e^{\nu(1+\log \zeta - \log x)x} dx. \tag{19}
\end{aligned}$$

If x_p is large enough then $1 + \log \zeta - \log x < 0$ for all $x \geq x_p$ so that

$$\nu(1 + \log \zeta - \log x) \leq \nu(1 + \log \zeta - \log x_p) \equiv -k < 0 \tag{20}$$

and then

$$\begin{aligned}
|I_b| &\leq (2\pi)^{-\nu/2} e^{(\beta\nu-2\pi)\delta} \int_{x_p}^{\infty} e^{-kx} dx = (2\pi)^{-\nu/2} e^{\tan \beta(\beta\nu-2\pi)x_p} \frac{e^{-kx_p}}{k} \\
&= (2\pi)^{-\nu/2} \frac{\exp([\{\tan \beta(\beta\nu - 2\pi) + \nu(1 + \log \zeta)\} - \nu \log x_p] x_p)}{\nu(\log x_p - 1 - \log \zeta)}. \tag{21}
\end{aligned}$$

No matter what the positive values of β , ν , or ζ are, if $x_p = (n+1/2)/2$ is large enough then the $\nu \log x_p$ term will make the exponent negative. Therefore as $n \rightarrow \infty$ the magnitude of the contour integral I_b goes to zero. As can be easily shown, this same result is produced for the C_f portion of the C_2 contour.

We next move on to the C_c portion of C_1 , where all of the points lie on a ray from the origin at angle β . Let $z = re^{i\beta} = x \sec \beta e^{i\beta} = x + ix \tan \beta$. First note that the exponential term will be bounded similarly as in (17) but now with $y = x \tan \beta$. Then for the contour integral I_c along C_c ,

$$\begin{aligned}
|I_c| &\leq \oint_{C_c} \left| \left(\frac{\zeta^z}{\Gamma(z+1)} \right)^\nu \frac{1}{e^{-i2\pi z} - 1} \right| |dz| \\
&\leq \int_{x_d}^{x_p} \left| \left(\frac{\zeta^{x+ix \tan \beta}}{\Gamma(x \sec \beta e^{i\beta} + 1)} \right)^\nu \frac{1}{e^{-i2\pi(x+ix \tan \beta)} - 1} \right| \sec \beta dx \tag{22}
\end{aligned}$$

$$\begin{aligned}
&\approx \sec \beta \int_{x_d}^{x_p} \left(\frac{\zeta^x e^{\beta x \tan \beta} e^x}{\sqrt{2\pi} x^{x+1/2} (\sec \beta)^{x+1/2}} \right)^\nu \frac{1}{|e^{2\pi x \tan \beta} - 1|} dx \\
&= \sec \beta \left(\frac{\cos \beta}{2\pi} \right)^{\nu/2} \int_{x_d}^{x_p} \frac{e^{\nu x(\kappa - \log x)}}{x^{\nu/2}} \frac{1}{|e^{2\pi x \tan \beta} - 1|} dx \tag{23}
\end{aligned}$$

for large x_p , where $\kappa = 1 + \beta \tan \beta + \log \cos \beta + \log \zeta$, a constant. Thus no matter what the value of κ is, if x_p is big enough then the $\log x$ term will eventually become larger so the integral will converge.

If we let $f(x) = \nu x(\kappa - \log x)$ then $f'(x) = \nu(\kappa - 1 - \log x)$ and $f''(x) = -\nu/x$. Thus $f(x)$ has a single maximum at

$$x_M = e^{\kappa-1} = \cos \beta e^{\beta \tan \beta} \zeta = \rho \zeta \quad (24)$$

where $\rho \equiv \cos \beta \exp(\beta \tan \beta)$. (As an example, if $\beta = \pi/4$ then $\rho \approx 1.55$.) Therefore if $x_p \gg \rho \zeta$ we can again use Laplace's method to evaluate (23).

Recalling that $x_p = (n+1/2)/2$, if n is chosen large enough then $x_p \gg \rho \zeta$ can be guaranteed. Therefore, using x_M as above with Laplace's method,

$$|I_c| \leq \frac{e^{\nu \rho \zeta}}{(2\pi \rho \zeta)^{(\nu-1)/2} \sqrt{\nu}} \frac{e^{-2\pi \tan \beta \rho \zeta}}{(\sec \beta)^{(\nu-2)/2}} \left\{ 1 + O\left(\frac{1}{\nu}\right) \right\}. \quad (25)$$

As $n \rightarrow \infty$ the upper endpoint of the integral, x_p , increases but the integral from x_p and beyond was already insignificant so the result above continues to hold. But unlike the results along C_a and C_b , the value of this integral cannot be guaranteed to go to zero with large ζ for all values of ν . What can be seen, though, is that the first parts of the integrands in (10) have the property that $\overline{f(z)} = f(\bar{z})$. This observation, combined with the result of (17) and the C_2 contour being a reflection of the C_1 contour, shows that the magnitude of the integral for C_2 will be the same as that for C_1 . Thus, if we can show that twice the ratio of the contour integral I_c to the approximation function (3) approaches zero then we will have achieved our goal of validating the approximation for all ν .

Therefore we divide (25) by (12) (the $Z(\lambda, \nu)$ approximation) and multiply by 2 to produce the error as a proportion of the approximation:

$$\eta = \frac{2}{\sqrt{\rho}} \left(\frac{\cos \beta}{\rho} \right)^{(\nu-2)/2} \exp \left\{ \left(\nu - \frac{2\pi \tan \beta \rho}{\rho - 1} \right) (\rho - 1) \zeta \right\}. \quad (26)$$

As $\beta \rightarrow 0$ both $\cos \beta$ and ρ approach 1 so the leading coefficient approaches 2. Inside the first parentheses in the $\exp\{\}$ brackets, both $\tan \beta$ and $\rho - 1$ approach zero from above but an application of l'Hôpital's Rule shows that the $\rho - 1$ approaches faster. The consequence is that, given any ν , it is possible to find a β such that the value inside the $\exp\{\}$ term is negative as $\zeta \rightarrow \infty$ (equivalently, $\lambda \rightarrow \infty$), so that the error ratio $\eta \rightarrow 0$. Thus the approximation (12) to $Z(\lambda, \nu)$ is valid for all $\nu > 0$.

However, note that the $\rho - 1$ in the second parentheses also becomes very small as $\beta \rightarrow 0$. This means that, even though the error eventually becomes small for small β , it does so increasingly slower. Why this occurs will be discussed further in Section 4 after first looking at the qualitative shapes of both $Z(\lambda, \nu)$ and its approximation, for various ν and λ , in the following section.

3 Function behaviors

We want to examine how $Z(\lambda, \nu)$ and its approximation behave for various values of ν and λ . For simplicity we designate the $Z(\lambda, \nu)$ approximation function (3) to be $\tilde{Z}(\lambda, \nu)$. In this section it usually will be helpful to think of both Z and \tilde{Z} as univariate functions of λ with ν simply as a parameter.

3.1 The behavior of $Z(\lambda, \nu)$

The simplest observation is to note that when $\lambda = 0$ then

$$Z(0, \nu) = \frac{0^0}{(0!)^\nu} + \sum_{j=1}^{\infty} \frac{0^j}{(j!)^\nu} = 1 \quad (27)$$

for any ν . Next,

$$\frac{\partial Z}{\partial \lambda} = \sum_{j=1}^{\infty} \frac{j \lambda^{j-1}}{(j!)^\nu} > 0 \quad (28)$$

also for all ν . In particular,

$$\left. \frac{\partial Z}{\partial \lambda} \right|_{\lambda=0} = \sum_{j=1}^{\infty} \left. \frac{j \lambda^{j-1}}{(j!)^\nu} \right|_{\lambda=0} = \frac{0^0}{1^\nu} + \sum_{j=2}^{\infty} \left. \frac{j \lambda^{j-1}}{(j!)^\nu} \right|_{\lambda=0} = 1 \quad (29)$$

for all ν . Thus, at $\lambda = 0$, $Z = 1$ with a slope of 1. It is obvious that the second derivative $\partial^2 Z / \partial \lambda^2 > 0$ for all finite ν as well, so that Z is increasing at a continually faster rate.

When $\nu = 1$ then Z reverts to a particularly simple form:

$$Z(\lambda, 1) = \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} = e^\lambda. \quad (30)$$

This corresponds to the Poisson distribution form of the CMP distribution.

Considering the change in Z with ν ,

$$\frac{\partial Z}{\partial \nu} = - \sum_{j=0}^{\infty} \log(j!) \frac{\lambda^j}{(j!)^\nu} < 0 \quad (31)$$

because $\log(j!) \lambda^j / (j!)^\nu \geq 0$, as long as ν is not zero or infinity. That is, as ν increases Z decreases at every fixed λ and, similarly, as ν decreases Z increases. The result is that the graph for Z flattens as ν increases from 1, and sharpens upward as ν decreases from 1.

For the first of the two endpoints of the range for ν ,

$$\lim_{\nu \rightarrow 0} Z(\lambda, \nu) = \sum_{j=0}^{\infty} \lambda^j. \quad (32)$$

As long as $0 \leq \lambda < 1$ then $Z(\lambda, 0) = 1/(1 - \lambda)$, which leads to the geometric distribution. But as $\lambda \rightarrow 1$ the sum approaches infinity. That is, for $\nu = 0$ the graph of Z curves up so sharply that λ cannot even get to 1 (or beyond). This is the restriction on λ referred to in Section 1.

The only remaining ν value is its second endpoint, where

$$\lim_{\nu \rightarrow \infty} Z(\lambda, \nu) = 1 + \lambda + \lim_{\nu \rightarrow \infty} \sum_{j=2}^{\infty} \frac{\lambda^j}{(j!)^\nu} = 1 + \lambda. \quad (33)$$

That is, the curve has flattened out completely and is now just a straight line with slope 1 and y -intercept of 1. The same simplification occurs to the distribution, though:

$$P(X = x) = \frac{\lambda^x}{(x!)^\nu Z(\lambda, \nu)} \quad (34)$$

so that $P(X = 0) = 1/(1 + \lambda)$ and $P(X = 1) = \lambda/(1 + \lambda)$ are the only survivors, producing a Bernoulli distribution.

The essential points of all of this are that $Z(\lambda, \nu)$ starts out at 1 and increases upward: sharper for smaller values of ν and flatter for larger values. In particular, as $\nu \rightarrow \infty$ the function becomes completely linear: $Z(\lambda, \infty) = 1 + \lambda$.

3.2 The behavior of $\tilde{Z}(\lambda, \nu)$

Just as in Section 2, it will be more convenient to use $\zeta = \lambda^{1/\nu}$ instead of λ when considering \tilde{Z} . As above, this results in the asymptotic approximation

$$\tilde{Z}(\lambda, \nu) = \frac{\exp(\nu\lambda^{\frac{1}{\nu}})}{\lambda^{\frac{\nu-1}{2\nu}}(2\pi)^{\frac{\nu-1}{2}}\sqrt{\nu}} = \frac{1}{(2\pi)^{\frac{\nu-1}{2}}\sqrt{\nu}} \zeta^{\frac{1-\nu}{2}} e^{\nu\zeta}. \quad (35)$$

Its behavior is more complicated than that for Z , in part because it does not behave similarly for all ν . Since $Z(\lambda, 0) = 1/(1-\lambda)$ and $Z(\lambda, \infty) = 1 + \lambda$, neither of which need approximation, we assume here for simplicity that $\nu \in (0, \infty)$ only. In this range an immediate observation is that $\lim_{\lambda \rightarrow \infty} \tilde{Z} \rightarrow \infty$ for all ν , which is just as we expect.

We first consider the single point when $\nu = 1$. In that case,

$$\tilde{Z}(\lambda, 1) = e^\lambda = Z(\lambda, 1), \quad (36)$$

after comparison with (30), so the approximation is exact. This is not the case when $\nu \neq 1$ so here it will be helpful to know the partial derivative $\partial\tilde{Z}/\partial\lambda$, again expressed in terms of ζ :

$$\frac{\partial\tilde{Z}}{\partial\lambda} = \frac{\partial\tilde{Z}}{\partial\zeta} \frac{\partial\zeta}{\partial\lambda} = \frac{\nu^{-3/2}}{2(2\pi)^{(\nu-1)/2}} (2\nu\zeta + 1 - \nu)\zeta^{(1-3\nu)/2} e^{\nu\zeta}. \quad (37)$$

Consider next the case when $0 < \nu < 1$. From (35) we see that

$$\lim_{\lambda \rightarrow 0} \tilde{Z} = \lim_{\zeta \rightarrow 0} \tilde{Z} = \frac{1}{(2\pi)^{\frac{\nu-1}{2}}\sqrt{\nu}} \lim_{\zeta \rightarrow 0} \zeta^{\frac{1-\nu}{2}} = 0 \quad (38)$$

and since \tilde{Z} increases indefinitely for large λ we should expect the slope at zero also to be increasing. But examination of (37) shows a more complicated picture. For $1/3 < \nu < 1$ the slope at zero is infinite; for $\nu = 1/3$ the slope is a positive constant; and for $0 < \nu < 1/3$ the slope is zero. However, looking at the higher derivatives shows that eventually $\partial\tilde{Z}/\partial\lambda$ does become positive for any ν , and so \tilde{Z} increases just as Z does.

Lastly, consider $\nu > 1$. Then

$$\lim_{\lambda \rightarrow 0} \tilde{Z} = \lim_{\zeta \rightarrow 0} \tilde{Z} = \frac{1}{(2\pi)^{\frac{\nu-1}{2}}\sqrt{\nu}} \lim_{\zeta \rightarrow 0} \frac{e^{\nu\zeta}}{\zeta^{(\nu-1)/2}} \rightarrow \infty \quad (39)$$

and, similarly,

$$\lim_{\lambda \rightarrow 0} \frac{\partial \tilde{Z}}{\partial \lambda} = \lim_{\zeta \rightarrow 0} \frac{\partial \tilde{Z}}{\partial \lambda} = \frac{\nu^{-3/2}}{2(2\pi)^{(\nu-1)/2}} (1 - \nu) \lim_{\zeta \rightarrow 0} \frac{1}{\zeta^{(3\nu-1)/2}} \rightarrow -\infty, \quad (40)$$

though the slope becomes finite (but still large and negative) as soon as $\lambda > 0$.

Thus \tilde{Z} starts out at infinity when $\lambda = 0$, but falls very rapidly. However, for large λ it increases indefinitely so the function must reverse direction. As can be seen from (37), a critical point can occur in \tilde{Z} when $2\nu\zeta + 1 - \nu = 0$, or

$$\lambda_m = \left(\frac{\nu - 1}{2\nu} \right)^\nu = \frac{1}{2^\nu} \left(1 - \frac{1}{\nu} \right)^\nu \leq \frac{1}{2^\nu} \leq \frac{1}{2}. \quad (41)$$

Examination of the second derivative shows that it is indeed positive there, as expected, so λ_m is a minimum. Then, combining (35) and (41) and recalling that $e = \lim_{n \rightarrow \infty} (1 + 1/n)^n$,

$$\tilde{Z}(\lambda_m, \nu) = \frac{\left(\frac{\nu-1}{2\nu} \right)^{\frac{1-\nu}{2}} e^{\frac{\nu-1}{2}}}{(2\pi)^{\frac{\nu-1}{2}} \sqrt{\nu}} = \sqrt{\frac{1}{\nu} \left(\frac{e\nu}{\pi(\nu-1)} \right)^{\nu-1}} \leq \sqrt{\frac{e}{\nu} \left(\frac{e}{\pi} \right)^{\nu-1}}, \quad (42)$$

and this is less than one certainly for $\nu \geq 3$. But (27) and (28) from Subsection 3.1 show that $Z(\lambda, \nu) \geq 1$ for all λ and ν . Therefore, for almost all ν , the function \tilde{Z} starts out at infinity with a negatively infinite slope, falls rapidly below 1 before $\lambda = 1/2^\nu$, crossing the curve of Z , and then turns upward to approach Z from below.

4 Discussion

The first question that might arise is that since our contour integrals representing the error in (10) are the same as Olver's and are analytic over the half-plane $\text{Re } z > -1$ where all of these contours lie, then any modifications to the contours should produce a net result of zero; thus we should find the same result as Olver. The explanation is that Olver only estimated the integrals via an upper bound. Of course, we did this as well, but the difference is that the contours we used produced better estimates and therefore our results are not limited to $\nu \leq 4$ as was the case with Olver's result summarized in (4). As can be seen on examining the definitions for both Z and \tilde{Z} there

seems to be no reasonable expectation why $\nu = 4$ should actually be a cutoff point for the accuracy of the approximation, and both our results and those of Shmueli *et al.* indicate that this cutoff point was merely an artifact of Olver's choice of contour.

Having said this the next question is why the error ratio η in (26) is so poor for large ν , requiring ever larger values of ζ before approaching a small proportion of the approximation. But after reviewing the behaviors of the two functions Z and \tilde{Z} in the previous section this should now be clear. For large ν the function \tilde{Z} dives from infinity at $\lambda = 0$ (as shown by (39) and (40)) to nearly zero (as shown by (42)) at $\lambda \leq 1/2^\nu$ (as shown by (41)). At this point \tilde{Z} reverses direction and rises up to meet an increasingly linear $Z(\lambda, \nu) \rightarrow 1 + \lambda$. For the exponential-like \tilde{Z} to approach the almost-linear function Z from below like this means that it must avoid it for as long as possible. The result is that the error proportion η remains large until λ is very large, diminishing only very slowly. And this is exactly what (26) demonstrates.

At the same time, though, this expression for η provides an estimate of how good \tilde{Z} should be as an approximation of Z for any given ν and λ . For η to approach zero the first parentheses in (26) must be negative, though only just barely, and the error will diminish most rapidly if the exponential is as large as possible. In that case (26) provides a minimum bound on the rate the error ratio η can be expected to diminish as a function of λ given any fixed value of ν : the error might go to zero faster than what this estimate provides, but it will not go any slower.

To better understand the effect that β has, first see that when $\beta = 0$ then $\rho = 1$. Then, for $0 < \beta < \pi/2$,

$$\frac{d\rho}{d\beta} = \cos \beta e^{\beta \tan \beta} (\beta \sec^2 \beta + \tan \beta) - \sin \beta e^{\beta \tan \beta} = \frac{\beta}{\cos \beta} e^{\beta \tan \beta} > 0 \quad (43)$$

so that ρ is an increasing function of β . If we describe ρ as $\exp(\beta \tan \beta) / \sec \beta$ then an application of l'Hôpital's Rule shows that $\lim_{\beta \rightarrow \pi/2} \rho \rightarrow \infty$, so ρ increases indefinitely as $\beta \rightarrow \pi/2$. This suggests that to maximize $\rho - 1$ one should maximize β , subject to keeping the first parentheses of (26) negative. But it is also possible to see that the ratio in the first parentheses of (26), $2\pi \tan \beta \rho / (\rho - 1)$, grows to infinity as $\beta \rightarrow \pi/2$ as well as when $\beta \rightarrow 0$, and has a minimum when $\rho - 1 - \beta \tan \beta = 0$. Therefore one might be tempted to let $\beta \rightarrow \pi/2$, causing $\rho - 1 \rightarrow \infty$ and thus getting a very large error reduction rate for any ν . But this would be a mistake: not only does this seem

impossible given the discussions on Z and \tilde{Z} in the previous section but it violates the conditions required to make the error contour integrals converge as shown by (24). Remember that ρ is not just an arbitrary parameter but instead describes the relationship between x_M and ζ . That is, $x_M = \rho\zeta$ and it is assumed that x_M is included in the C_c contour along the ray at angle β . If ρ is very large then it could happen that $x_M > x_p$ and thus no longer be in C_c , thus invalidating (26) as a measure of the error ratio. The result is that, for large ν , one must bring the contour close to the real line by making β small and therefore ρ close to 1 with the error only slowly decreasing, as matches what was seen in Section 3.

Computer calculations show that $\rho - 1 - \beta \tan \beta = 0$ when $\beta \approx 0.998$, making $\rho \approx 2.55$. Then the ratio in the first parentheses of (26) is about 16.04. This means that, for any $\nu \leq 16$, it does no good to try to increase β to increase $\rho - 1$ to make the error ratio decrease faster because then one risks raising ρ so much as to enter the region where $x_M > x_p$. Alternatively, for $\nu > 16$, one does need to start giving up \tilde{Z} approximation strength by decreasing β and, consequently, $\rho - 1$. In the data set provided by Shmueli *et al.* as an underdispersion ($\nu > 1$) example, the estimated ν was about 2.15. If this is an example of the usual deviation from $\nu = 1$ (the Poisson case) for non-Poisson data then one could expect most ν arising in practice to fall into this situation of being less than 16. For this same data set, $\lambda \approx 7.74$, $\zeta \approx 2.59$, and $\eta \approx 10^{-24}$, an entirely acceptable approximation!

Of more interest, then, would be estimates of when the approximation was not at the level of 10^{-24} . Since $Z(0, \nu) = 1$ and $\tilde{Z}(0, \nu > 1) = 0$, clearly \tilde{Z} should be a poor approximation for small λ . Using the minimum $\beta \approx 0.998$ above so that

$$\eta \approx 1.25 (0.213)^{(\nu-2)/2} e^{(\nu-16)1.55\zeta} \quad (44)$$

produces $\zeta \approx 0.541$ when $\nu = 2$ and $\eta = 0.00001$. This corresponds to $\lambda \approx 0.293$, meaning that for any $\lambda > 0.293$, \tilde{Z} could approximate Z with an error of less than 1/1000th of a percent. If instead $\nu = 15$ then the same accuracy is achieved for $\lambda > 3.47$. These are excellent approximations, but they do not seem to match what was shown in the previous section.

The explanation is that η only estimates the error of approximating Z with \tilde{Z} using the method above of contour integrals. So while it shows that the approximation is indeed very good for almost all ν and λ , it overlooks the error in the approximation method not from the contour integrals on line (10) but from using Laplace's method to evaluate the original approxi-

mation integral on line (9). This error is $O(1/\zeta)$ and for the Shmueli *et al.* example it is about 39%. (Though the computer comparisons by Shmueli *et al.* indicate that the actual errors are not nearly that large.) Thus, for those instances of real data where ζ is not large, this is the main source of error in calculating $\tilde{Z}(\lambda, \nu)$. As a consequence, for these values of ζ , a better estimate of $Z(\lambda, \nu)$ would be to use a truncated summation of the original equation (2). Note that if the summands of (2) are denoted as s_j , then

$$s_j = \frac{\lambda^j}{(j!)^\nu} = \frac{\lambda^{j-1}}{((j-1)!)^\nu} \frac{\lambda}{j^\nu} = s_{j-1} \frac{\lambda}{j^\nu} \quad (45)$$

and this type of formulation is especially amenable to calculations in a computer spreadsheet. In fact, setting up such a spreadsheet shows that for a large range of ν and λ the summation converges with a very high precision using less than twenty terms. On the other hand, the approximation function \tilde{Z} is also easy to calculate so one can simply do both and compare the two numbers. Thus, when the value of $\zeta = \lambda^{1/\nu}$ is small then it is likely that the truncated summation will be the most accurate method of calculating $Z(\lambda, \nu)$. But when ζ is large, so that many terms in the sum are required, then the summation method risks being incorrect due to the aggregation of multiple rounding errors; in that case the use of \tilde{Z} may be the better method.

A third method is to use numerical integration to calculate the value of the integral in (9) since the error, as described above, occurs primarily due to the use of Laplace's method to derive the $\tilde{Z}(\lambda, \nu)$ approximation. But the very reason Laplace's method works is because the integrand of this integral is very sharply peaked: the numerator is rising quickly but the denominator falls off even faster once it has recovered from its maximum. The result is a very narrow range over x for which the integrand has significant magnitude, with very steep approaches on both sides to the peak value. This requires very high quality software to accurately evaluate such an integral, and that depends on what a particular researcher may have at hand. Nevertheless, this is another option to consider.

In conclusion, we have shown here that the conjecture of Shmueli *et al.* that the function $\tilde{Z}(\lambda, \nu)$ is a valid approximation of $Z(\lambda, \nu)$ for almost all ν and λ is true. At the same time, however, we have also shown that simply using the truncated original summation may be the most accurate calculation for a large number of data sets encountered in actual practice.

References

- Abramowitz, M. and I. A. Stegun (1964). *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables* (ninth Dover printing, tenth GPO printing ed.). New York: Dover.
- Conway, R. W. and W. L. Maxwell (1962). A queuing model with state dependent service rates. *J. Industrial Engineering* 12, 132–136.
- Olver, F. W. J. (1974). *Asymptotics and Special Functions*. New York: Academic Press.
- Shmueli, G., T. P. Minka, J. B. Kadane, S. Borle, and P. Boatwright (2005). A useful distribution for fitting discrete data: revival of the Conway–Maxwell–Poisson distribution. *J. R. Statist. Soc. C* 54(1), 127–142.