

# Using the **PDQutils** package

Steven E. Pav \*

January 4, 2016

## Abstract

Example computations via the **PDQutils** package are illustrated.

The **PDQutils** package provides tools for approximating the density, distribution, and quantile functions, and for generation of random variates of distributions whose cumulants and moments can be computed. The PDF and CDF are computed approximately via the Gram Charlier A series, while the quantile is computed via the Cornish Fisher approximation. [3, 7] The random generation function uses the quantile function and draws from the uniform distribution.

## 1 Gram Charlier Expansion

Given the raw moments of a probability distribution, we can approximate the probability density function, or the cumulative distribution function, via a Gram-Charlier A expansion. This is typically developed as an approximation to the normal distribution using Hermite polynomials, but here we follow a more general derivation, which allows us to approximate distributions which are more like a gamma or beta.

Let  $w(x)$  be some non-negative ‘weighting function’, typically the PDF of a known probability distribution. Let  $p_n(x)$  be polynomials which are orthogonal with respect to this weighting function. That is

$$\int_{-\infty}^{\infty} w(x) p_n(x) p_m(x) dx = \delta_{n,m} h_n, \quad (1)$$

where  $\delta_{m,n}$  is the Kronecker delta, equal to one only when  $m = n$ , otherwise equal to zero. We furthermore suppose that the polynomials  $p_n(x)$  are complete: any reasonably smooth function can be represented as a linear combination of these polynomials.

Then we can expand the probability density of some random variable,  $f(x)$  in terms of this basis. Let

$$f(x) = \sum_{n=0}^{\infty} c_n p_n(x) w(x). \quad (2)$$

By the orthogonality relationship, we can find the constants  $c_n$  by multiplying both sides by  $p_m(x)$  and integrating:

$$\int_{-\infty}^{\infty} p_m(x) f(x) dx = \sum_{n=0}^{\infty} c_n \int_{-\infty}^{\infty} p_m(x) p_n(x) w(x) dx = c_m h_m. \quad (3)$$

---

\*shabbychef@gmail.com

Thus

$$c_n = \frac{1}{h_n} \int_{-\infty}^{\infty} p_n(x) f(x) dx$$

When the coefficients of the polynomial  $p_n(x)$  and the uncentered moments of the probability distribution are known, the constant  $c_n$  can easily be computed.

Thus the density  $f(x)$  can be approximated by truncating the infinite sum as

$$f(x) \approx \sum_{n=0}^m p_n(x) w(x) \left[ \frac{1}{h_n} \int_{-\infty}^{\infty} p_n(z) f(z) dz \right]. \quad (4)$$

To approximately compute the cumulative distribution function, one can compute the integral of the approximate density. The approximation is

$$F(x) \approx \sum_{n=0}^m \int_{-\infty}^{\infty} p_n(y) w(y) dy \left[ \frac{1}{h_n} \int_{-\infty}^{\infty} p_n(z) f(z) dz \right]. \quad (5)$$

In summary, to approximate the PDF or CDF of a distribution via the Gram Charlier series, one must know the moments of the distribution, and be able to compute  $w(x)$ ,  $p_n(x)$ ,  $h_n$ , and  $\int p_n(y) w(y) dy$ . These are collected in Table 1 for a few different families of probability distributions. [1, 22.2] The traditional Gram Charlier ‘A’ series corresponds to case where  $w(x)$  is the PDF of the standard normal distribution and  $p_n(x)$  is the (probabilist’s) Hermite polynomial. Also of interest are the cases where  $w(x)$  is PDF of the gamma distribution (including Chi-squares), in which case  $p_n(x)$  are the generalized Laguerre polynomials; the case where  $w(x)$  is the PDF of the (shifted) Beta distribution, and  $p_n(x)$  are the Jacobi polynomials. As special cases of the Beta distribution, one also has the Arcsine distribution (with Chebyshev polynomials of the first kind), the Wigner distribution (Chebyshev of the second kind), and the uniform distribution (Legendre polynomials). [1]

## 2 Edgeworth Expansion

Another approximation of the probability density and cumulative distribution functions is the Edgeworth Expansions. These are expressed in terms of the cumulants of the distribution, and also include the Hermite polynomials. However, the derivation of the Edgeworth expansion is rather more complicated than of the Gram Charlier expansion. [3] The Edgeworth series for a zero-mean unit distribution is

$$f(x) = \frac{1}{\sigma} \phi\left(\frac{x}{\sigma}\right) \left[ 1 + \sum_{1 \leq s} \sigma^s \sum_{\{k_m\}} He_{s+2r}(x/\sigma) \prod_{1 \leq m \leq s} \frac{1}{k_m!} \left( \frac{S_{m+2}}{(m+2)!} \right)^{k_m} \right],$$

where the second sum is over partitions  $\{k_m\}$  such that  $k_1 + 2k_2 + \dots + sk_s = s$ , where  $r = k_1 + k_2 + \dots + k_s$ , and where  $S_n = \frac{\kappa_n}{\sigma^{2n-2}}$  is a semi-normalized cumulant.

class	support	$w(x)$
Normal/Hermite	$(-\infty, \infty)$	$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$
Gamma/generalized Laguerre	$[0, \infty)$	$g_{a+1}(x) = \frac{1}{\Gamma(a+1)} x^a \exp(-x)$
Beta/Jacobi	$(-1, 1)$	$f_{a+1,b+1}(x) = \frac{(1-x)^a (1+x)^b}{B(a+1,b+1) 2^{a+b+1}}$

class	$p_n(x)$	$h_n$
Normal	$He_n(x)$	$n!$
Gamma	$L_n^{(a)}(x)$	$\frac{\Gamma(n+a+1)}{\Gamma(a+1)n!}$
Beta	$P_n^{(a,b)}(x)$	$\frac{1}{2n+a+b+1} \frac{1}{B(a+1,b+1)} \frac{\Gamma(n+a+1)\Gamma(n+b+1)}{n!\Gamma(n+a+b+1)}$

class	$\int p_n(y) w(y) dy$
Normal	$-\phi(y) He_{n-1}(y)$
Gamma	$\left(\frac{a+1}{n}\right) g_{a+2}(y) L_{n-1}^{(a+1)}(y)$
Beta	$\left(\frac{-2}{n}\right) \left(\frac{B(a+2,b+2)}{B(a+1,b+1)}\right) f_{a+2,b+2}(y) P_{n-1}^{(a+1,b+1)}(y)$

Table 1: Different classes of orthogonal polynomials are presented. In each case the weight function,  $w(x)$  is the PDF of a common distribution, while the orthogonal polynomials come from a well known family. The constant  $h_n$  is the normalizing constant. The last table gives the integral of the polynomial times the weighting function, a value which is needed for approximating the CDF. Values are given for: the normal PDF, with probabilist's Hermite polynomials; the Gamma PDF, with generalized Laguerre polynomials; the Beta PDF with Jacobi polynomials. As special cases of the latter, one has the Arcsine, Wigner, and Uniform distributions, with Chebyshev and Legendre polynomials.

### 3 Cornish Fisher Approximation

The Cornish Fisher approximation is the Legendre inversion of the Edgeworth expansion of a distribution, but ordered in a way that is convenient when used on the mean of a number of independent draws of a random variable.

Suppose  $x_1, x_2, \dots, x_n$  are  $n$  independent draws from some probability distribution. Letting

$$X = \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} x_i,$$

the Central Limit Theorem assures us that, assuming finite variance,

$$X \rightarrow \mathcal{N}(\sqrt{n}\mu, \sigma),$$

with convergence in  $n$

The Cornish Fisher approximation gives a more detailed picture of the quantiles of  $X$ , one that is arranged in decreasing powers of  $\sqrt{n}$ . The quantile function is the function  $q(p)$  such that  $P(X \leq q(p)) = q(p)$ . The Cornish Fisher expansion is

$$q(p) = \sqrt{n}\mu + \sigma \left( z + \sum_{3 \leq j} c_j f_j(z) \right),$$

where  $z = \Phi^{-1}(p)$  is the normal  $p$ -quantile, and  $c_j$  involves standardized cumulants of the distribution of  $x_i$  of order up to  $j$ . Moreover, the  $c_j$  include

decreasing powers of  $\sqrt{n}$ , giving some justification for truncation. When  $n = 1$ , however, the ordering is somewhat arbitrary.

## 4 An Example: Sum of Nakagamis

The Gram Charlier and Cornish Fisher approximations are most convenient when the random variable can be decomposed as the sum of a small number of independent random variables whose cumulants can be computed. For example, suppose  $Y = \sum_{1 \leq i \leq k} \sqrt{X_i/\nu_i}$  where the  $X_i$  are independent central chi-square random variables with degrees of freedom  $\nu_1, \nu_2, \dots, \nu_k$ . I will call this a ‘snak’ distribution, since each summand follows a Nakagami distribution. We can easily write code that generates variates from this distribution given a vector of the degrees of freedom:

```
rsnak <- function(n, dfs) {
  samples <- Reduce("+", lapply(dfs, function(k) {
    sqrt(rchisq(n, df = k)/k)
  }))
}
```

Let’s take one hundred thousand draws from this distribution. A q-q plot of this sample against normality is shown in Figure 1. The normal model is fairly decent, although possibly unacceptable in the tails. Using a Cornish Fisher approximation, we can do better.

```
n.samp <- 1e+05
dfs <- c(8, 15, 4000, 10000)
set.seed(18181)
# now draw from the distribution
rvs <- rsnak(n.samp, dfs)
data <- data.frame(draws = rvs)

library(ggplot2)
mu <- mean(rvs)
sigma <- sd(rvs)
ph <- ggplot(data, aes(sample = draws)) + stat_qq(dist = function(p) {
  qnorm(p, mean = mu, sd = sigma)
}) + geom_abline(slope = 1, intercept = 0, colour = "red") +
  theme(text = element_text(size = 8)) + labs(title = "Q-Q plot (against normality)")

print(ph)
```

Using the additivity property of cumulants, we can compute the cumulants of  $Y$  easily if we have the cumulants of the  $X_i$ . These in turn can be computed from the raw moments. The  $j$ th moment of a chi distribution with  $\nu$  degrees of freedom has form

$$2^{j/2} \frac{\Gamma((\nu + j)/2)}{\Gamma(\nu/2)}.$$

The following function computes the cumulants of the ‘snak’ distribution:

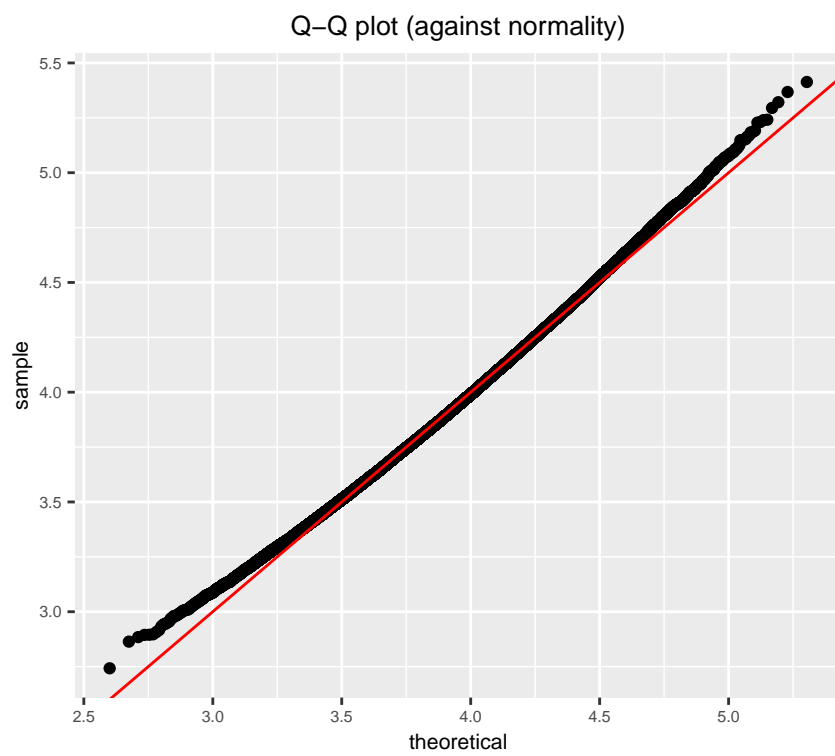


Figure 1: A q-q plot of  $1e+05$  draws from a sum of Nakagamis distribution is shown against normality.

```

# for the moment2cumulant function:
library(PDQutils)

# compute the first ord.max raw cumulants of the
# sum of Nakagami variates
snak_cumulants <- function(dfs, ord.max = 10) {
  # first compute the raw moments
  moms <- lapply(dfs, function(nu) {
    ords <- 1:ord.max
    moms <- 2^(ords/2) * exp(lgamma((nu + ords)/2) -
      lgamma(nu/2))
    # we are dividing the chi by sqrt the d.f.
    moms <- moms/(nu^(ords/2))
    moms
  })
  # turn moments into cumulants
  cumuls <- lapply(moms, moment2cumulant)
  # sum the cumulants
  tot.cumul <- Reduce("+", cumuls)
  return(tot.cumul)
}

```

We can now trivially implement the ‘dpq’ functions trivially using the Gram-Charlier and Cornish-Fisher approximations, via [PDQutils](#), as follows:

```

library(PDQutils)

dsnak <- function(x, dfs, ord.max = 10, ...) {
  raw.moment <- cumulant2moment(snak_cumulants(dfs,
    ord.max))
  retval <- dapx_gca(x, raw.moment, support = c(0,
    Inf), ...)
  return(retval)
}

psnak <- function(q, dfs, ord.max = 10, ...) {
  raw.moment <- cumulant2moment(snak_cumulants(dfs,
    ord.max))
  retval <- papx_gca(q, raw.moment, support = c(0,
    Inf), ...)
  return(retval)
}

qsnak <- function(p, dfs, ord.max = 10, ...) {
  raw.cumul <- snak_cumulants(dfs, ord.max)
  retval <- qapx_cf(p, raw.cumul, support = c(0,
    Inf), ...)
  return(retval)
}

```

An alternative version of the PDF and CDF functions using the Edgeworth expansion would look as follows:

```

dsnak_2 <- function(x, dfs, ord.max = 10, ...) {
  raw.cumul <- snak_cumulants(dfs, ord.max)

```

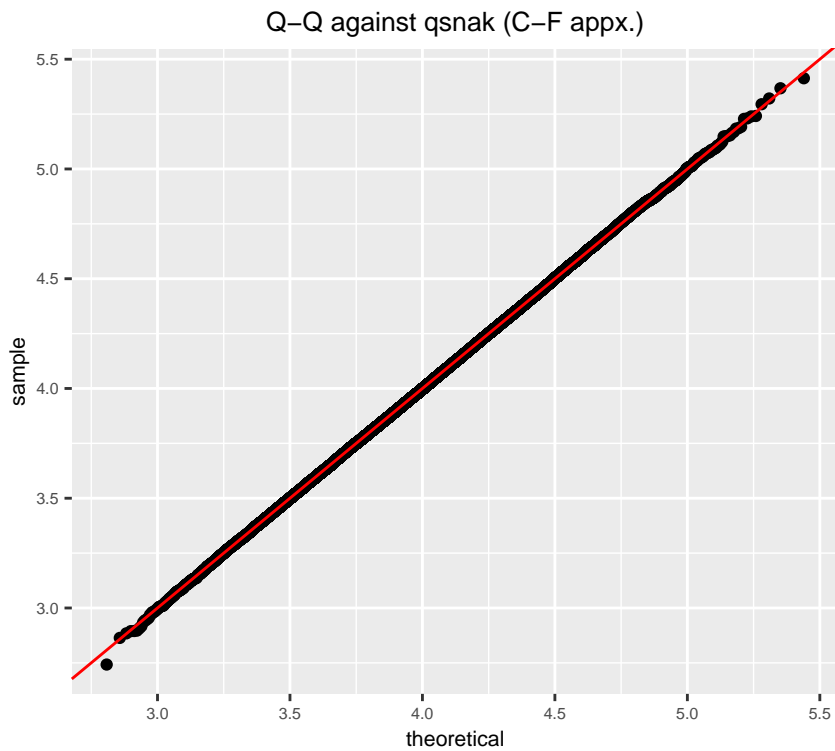


Figure 2: A q-q plot of  $1e+05$  draws from a sum of Nakagamis distribution is shown against quantiles from the ‘qsnak’ function.

```

    retval <- dapx_edgeworth(x, raw.cumul, support = c(0,
      Inf), ...)
    return(retval)
  }
  psnak_2 <- function(q, dfs, ord.max = 10, ...) {
    raw.cumul <- snak_cumulants(dfs, ord.max)
    retval <- papx_edgeworth(q, raw.cumul, support = c(0,
      Inf), ...)
    return(retval)
  }

```

Using this approximate quantile function, the q-q plot looks straighter, as shown in Figure 2.

```

data <- data.frame(draws = rvs)
library(ggplot2)
ph <- ggplot(data, aes(sample = draws)) + stat_qq(dist = function(p) {
  qsnak(p, dfs = dfs)
}) + geom_abline(slope = 1, intercept = 0, colour = "red") +
  theme(text = element_text(size = 8)) + labs(title = "Q-Q against qsnak (C-F appx.)")
print(ph)

```

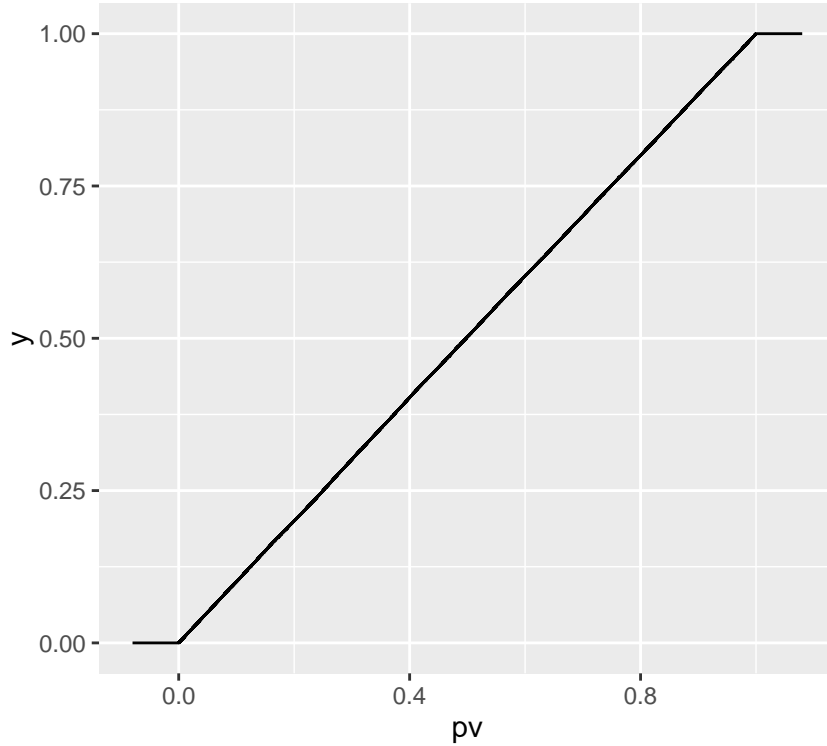


Figure 3: The empirical CDF of the approximate CDF of a sum of Nakagamis distribution on  $1e+05$  draws is shown.

Note that the q-q plot uses the approximate quantile function, computed via the Cornish-Fisher expansion. We can test the Gram Charlier expansion by computing the approximate CDF of the random draws and checking that it is nearly uniform, as shown in Figure 3.

```
apx.p <- psnak(rvs, dfs = dfs)
require(ggplot2)
ph <- ggplot(data.frame(pv = apx.p), aes(x = pv)) +
  stat_ecdf(geom = "step")
print(ph)
```

## 5 A warning on convergence

Blinnikov and Moessner note that the the Gram Charlier expansion will actually diverge for some distributions when more terms in the expansion are considered, behaviour which is not seen for the Edgeworth expansion. [3] Here, we will replicate their example of the chi-square distribution with 5 degrees of freedom. Blinnikov and Moessner actually transform the chi-square to have zero mean and unit variance. They plot the true PDF of this normalized distribution, along with the 2- and 6-term Gram Charlier approximations, as shown in Figure 4.



```

# compute moments and cumulants:
df <- 5
max.ord <- 20
subords <- 0:(max.ord - 1)
raw.cumulants <- df * (2^subords) * factorial(subords)
raw.moments <- cumulant2moment(raw.cumulants)

# compute the PDF of the rescaled variable:
xvals <- seq(-sqrt(df/2) * 0.99, 6, length.out = 1001)
chivals <- sqrt(2 * df) * xvals + df
pdf.true <- sqrt(2 * df) * dchisq(chivals, df = df)

pdf.gca2 <- sqrt(2 * df) * dapx_gca(chivals, raw.moments = raw.moments[1:2],
  support = c(0, Inf))
pdf.gca6 <- sqrt(2 * df) * dapx_gca(chivals, raw.moments = raw.moments[1:6],
  support = c(0, Inf))

all.pdf <- data.frame(x = xvals, true = pdf.true, gca2 = pdf.gca2,
  gca6 = pdf.gca6)

# plot it by reshaping and ggplot'ing
require(reshape2)
arr.data <- melt(all.pdf, id.vars = "x", variable.name = "pdf",
  value.name = "density")

require(ggplot2)
ph <- ggplot(arr.data, aes(x = x, y = density, group = pdf,
  colour = pdf)) + geom_line()
print(ph)

```

Compare this with the 2 and 4 term Edgeworth expansions, shown in Figure 5.

```

# compute the PDF of the rescaled variable:
xvals <- seq(-sqrt(df/2) * 0.99, 6, length.out = 1001)
chivals <- sqrt(2 * df) * xvals + df
pdf.true <- sqrt(2 * df) * dchisq(chivals, df = df)

pdf.edgeworth2 <- sqrt(2 * df) * dapx_edgeworth(chivals,
  raw.cumulants = raw.cumulants[1:4], support = c(0,
  Inf))
pdf.edgeworth4 <- sqrt(2 * df) * dapx_edgeworth(chivals,
  raw.cumulants = raw.cumulants[1:6], support = c(0,
  Inf))

all.pdf <- data.frame(x = xvals, true = pdf.true, edgeworth2 = pdf.edgeworth2,
  edgeworth4 = pdf.edgeworth4)

# plot it by reshaping and ggplot'ing
require(reshape2)
arr.data <- melt(all.pdf, id.vars = "x", variable.name = "pdf",
  value.name = "density")

require(ggplot2)

```

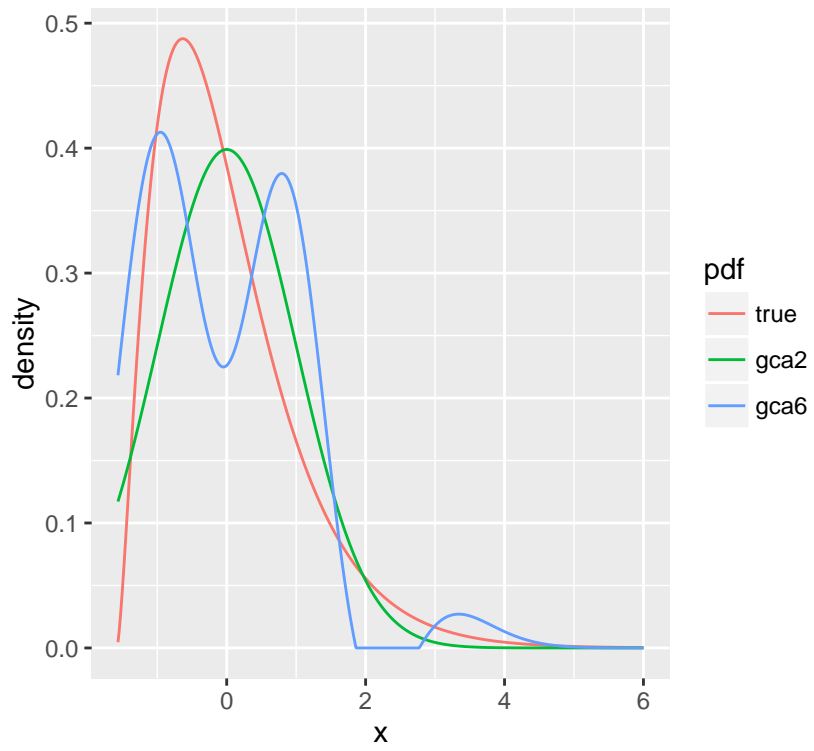


Figure 4: The true PDF of a normalized  $\chi_5^2$  distribution is shown, along with the 2- and 6-term Gram Charlier approximations. This replicates Figure 1 of Blinnikov and Moessner. [3]

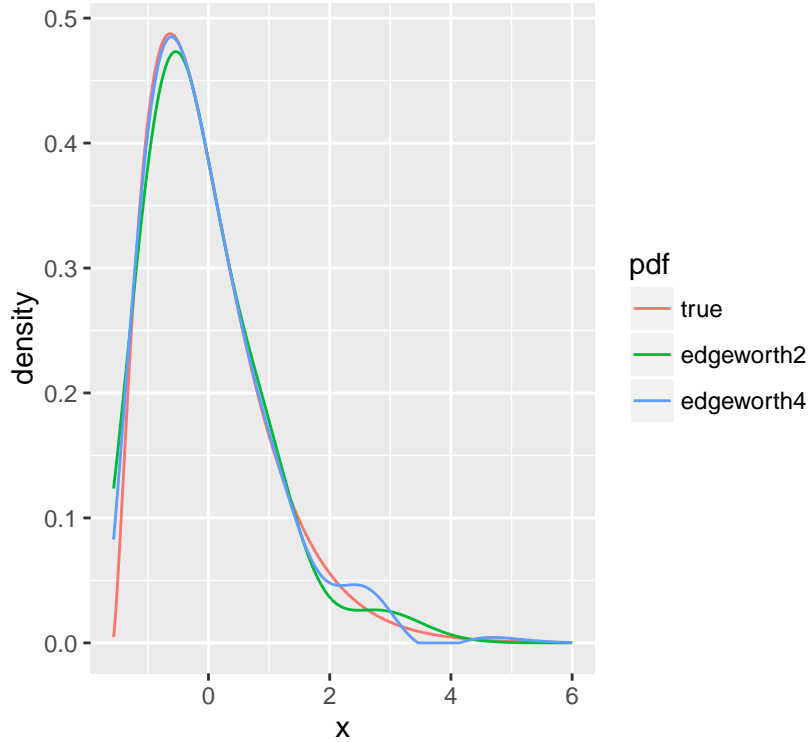


Figure 5: The true PDF of a normalized  $\chi_5^2$  distribution is shown, along with the 2- and 4-term Edgeworth expansions. This replicates Figure 6 of Blinnikov and Moessner. [3]

```
ph <- ggplot(arr.data, aes(x = x, y = density, group = pdf,
  colour = pdf)) + geom_line()
print(ph)
```

## References

- [1] Milton Abramowitz and Irene A. Stegun. *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Dover, New York, ninth dover printing, tenth gpo printing edition, 1964. URL <http://people.math.sfu.ca/~cbm/aands/toc.htm>.
- [2] Mário N. Berberan-Santos. Expressing a probability density function in terms of another PDF: A generalized Gram-Charlier expansion. *Journal of Mathematical Chemistry*, 42(3):585–594, 2007. ISSN 0259-9791. doi: 10.1007/s10910-006-9134-5. URL <http://web.ist.utl.pt/ist12219/data/115.pdf>.
- [3] S. Blinnikov and R. Moessner. Expansions for nearly Gaussian distributions. *Astronomy and Astrophysics Supplement*, 130:193–205, May

1998. doi: 10.1051/aas:1998221. URL <http://arxiv.org/abs/astro-ph/9711239>.
- [4] P. K. Cheah, D. A. S. Fraser, and N. Reid. Some alternatives to Edgeworth. *Canadian Journal of Statistics*, 21(2):131–138, 1993. URL <http://fisher.utstat.toronto.edu/dfraser/documents/174.pdf>.
  - [5] Victor Chernozhukov, Iván Fernández-Val, and Alfred Galichon. Rearranging Edgeworth-Cornish-Fisher expansions. Privately Published, 2007. URL <http://arxiv.org/abs/0708.1627>.
  - [6] Leon Cohen. On the generalization of the Edgeworth/Gram-Charlier series. *Journal of Mathematical Chemistry*, 49(3):625–628, 2011. ISSN 0259-9791. doi: 10.1007/s10910-010-9787-y. URL <http://dx.doi.org/10.1007/s10910-010-9787-y>.
  - [7] Stefan R. Jaschke. The Cornish-Fisher-expansion in the context of Delta - Gamma - Normal approximations. Technical Report 2001,54, Humboldt University of Berlin, Interdisciplinary Research Project 373: Quantification and Simulation of Economic Processes, 2001. URL <http://www.jaschke-net.de/papers/CoFi.pdf>.
  - [8] Yoong-Sin Lee and Ting-Kwong Lin. Algorithm AS 269: High order Cornish-Fisher expansion. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 41(1):pp. 233–240, 1992. ISSN 00359254. URL <http://www.jstor.org/stable/2347649>.
  - [9] Paul A. Samuelson. Fitting general Gram-Charlier series. *The Annals of Mathematical Statistics*, 14(2):pp. 179–187, 1943. ISSN 00034851. URL <http://projecteuclid.org/euclid.aoms/1177731459>.
  - [10] Vladimir V. Ulyanov. Cornish Fisher expansions. In Miodrag Lovric, editor, *International Encyclopedia of Statistical Science*, pages 312–315. Springer Berlin Heidelberg, 2014. ISBN 978-3-642-04897-5. doi: 10.1007/978-3-642-04898-2\_193. URL [http://dx.doi.org/10.1007/978-3-642-04898-2\\_193](http://dx.doi.org/10.1007/978-3-642-04898-2_193).