

Fast evaluation of the inverse Poisson CDF

Mike Giles

University of Oxford
Mathematical Institute

Ninth IMACS Seminar on Monte Carlo Methods

July 16, 2013

Outline

- problem specification
- incomplete Gamma function
- CPUs versus GPUs
- asymptotic Normal approximation
- asymptotic Temme approximation
- Temme asymptotic evaluation
- putting it all together

Poisson CDF and inverse

The CDF for Poisson rate λ is

$$\bar{C}(n) \equiv \mathbb{P}(N \leq n) = e^{-\lambda} \sum_{m=0}^n \frac{\lambda^m}{m!}.$$

The inverse CDF is defined as $\bar{C}^{-1}(u) = n$ where n is the smallest integer such that

$$u \leq e^{-\lambda} \sum_{m=0}^n \frac{\lambda^m}{m!} \quad (\text{bottom-up})$$

or

$$1 - u \geq e^{-\lambda} \sum_{m=n+1}^{\infty} \frac{\lambda^m}{m!} \quad (\text{top-down})$$

Poisson CDF and inverse

When λ is fixed and not too large ($\lambda < 10^4$?) can pre-compute $\overline{C}(n)$ and perform a table lookup.

When λ is variable but small ($\lambda < 10$?) can use bottom-up/top-down summation.

When λ is variable and large, then rejection methods can be used to generate Poisson r.v.'s, but the inverse CDF is sometimes helpful:

- stratified sampling
- Latin hypercube
- QMC

This is the problem I am concerned with — approximating $\overline{C}^{-1}(u)$ at a cost similar to the inverse Normal CDF, or inverse error function.

Incomplete Gamma function

If X is a positive random variable with CDF

$$C(x) \equiv \mathbb{P}(X < x) = \frac{1}{\Gamma(x)} \int_{\lambda}^{\infty} e^{-t} t^{x-1} dt.$$

then integration by parts gives

$$\mathbb{P}(\lfloor X \rfloor \leq n) = \frac{1}{n!} \int_{\lambda}^{\infty} e^{-t} t^n dt = e^{-\lambda} \sum_{m=0}^n \frac{\lambda^m}{m!}$$

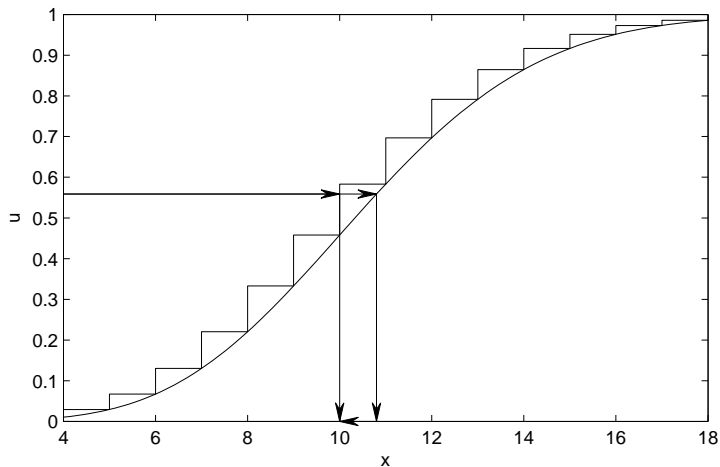
$$\implies \bar{C}^{-1}(u) = \lfloor C^{-1}(u) \rfloor$$

We will approximate $Q(u) \equiv C^{-1}(u)$ so that $|\tilde{Q}(u) - Q(u)| < \delta \ll 1$

This will round down correctly except when $Q(u)$ is within δ of an integer – then we need to check some $\bar{C}(m)$

Incomplete Gamma function

Illustration of the rounding down of $Q(u) \equiv C^{-1}(u)$ to give $\bar{C}^{-1}(u)$



CPUs and GPUs

On a CPU, if the costs of $\tilde{Q}(u)$ and $\overline{C}(m)$ are C_Q and C_C , the average cost is approximately

$$C_Q + 2\delta C_C.$$

However, on a GPU with a vector length of 32, the C_C penalty is incurred if any element needs it, so the average cost is

$$C_Q + (1 - (1 - 2\delta)^{32}) C_C \approx C_Q + 64\delta C_C \quad \text{if } \delta \ll 1.$$

This pushes us to more accurate approximations for GPUs.

Normal approximation

It's well known that

$$C(x) \approx \Phi\left(\frac{x-\lambda}{\sqrt{\lambda}}\right)$$

which motivates the following change of variables

$$x = \lambda + \sqrt{\lambda} y, \quad t = \lambda + \sqrt{\lambda} (y-z)$$

giving

$$C(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^y I(y, z) dz$$

where

$$\log I = \frac{1}{2} \log(2\pi) - \log \Gamma(x) - t + (x-1) \log t - \frac{1}{2} \log \lambda$$

Normal approximation

An asymptotic expansion in powers of $\varepsilon \equiv \lambda^{-1/2}$ yields

$$l(y, z) = \exp\left(-\frac{1}{2}z^2\right) \left(1 + \sum_{n=1}^{\infty} \varepsilon^n p_n(y, z)\right)$$

where $p_n(y, z)$ are polynomial in y and z . Integrating by parts gives

$$C(x) \approx \Phi(y) + \phi(y) \left(\varepsilon \left(-\frac{1}{3} - \frac{1}{6} y^2\right) + \varepsilon^2 \left(\frac{1}{12} y + \frac{1}{72} y^3 - \frac{1}{72} y^5\right) + \varepsilon^3 \left(-\frac{1}{540} - \frac{23}{540} y^2 + \frac{7}{2160} y^4 + \frac{5}{648} y^6 - \frac{1}{1296} y^8\right) \right)$$

and inverting this gives the asymptotic expansion

$$Q(u) = \lambda + \sqrt{\lambda} w + \left(\frac{1}{3} + \frac{1}{6} w^2\right) + \lambda^{-1/2} \left(-\frac{1}{36} w - \frac{1}{72} w^3\right) + \lambda^{-1} \left(-\frac{8}{405} + \frac{7}{810} w^2 + \frac{1}{270} w^4\right) + O(\lambda^{-3/2})$$

where $w = \Phi^{-1}(u)$.

Normal approximation

All asymptotic expansions were performed using MATLAB's Symbolic Toolbox.

This gives three approximations:

$$\tilde{Q}_{N1}(u) = \lambda + \sqrt{\lambda} w + \left(\frac{1}{3} + \frac{1}{6} w^2\right)$$

$$\tilde{Q}_{N2}(u) = \tilde{Q}_{N1}(u) + \lambda^{-1/2} \left(-\frac{1}{36} w - \frac{1}{72} w^3\right)$$

$$\tilde{Q}_{N3}(u) = \tilde{Q}_{N2}(u) + \lambda^{-1} \left(-\frac{8}{405} + \frac{7}{810} w^2 + \frac{1}{270} w^4\right)$$

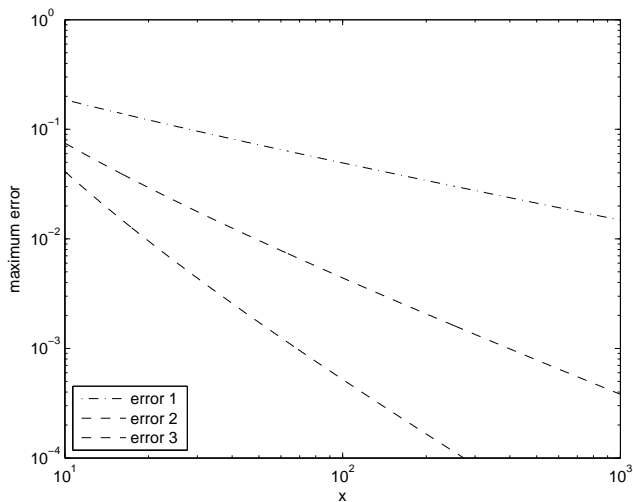
and suggests an error bound for \tilde{Q}_{N2} :

$$\delta = \lambda^{-1} \left(\frac{1}{40} + \frac{1}{80} w^2 + \frac{1}{160} w^4\right)$$

with $\mathbb{E}[\delta] = \frac{9}{160} \lambda^{-1}$.

Normal approximation

Maximum error over range $|w| \leq 3$:



Temme approximation

The Normal approximation is not good when w is large, and also the asymptotic convergence is poor in powers of $\lambda^{-1/2}$.

Temme (1979) derived a uniformly convergent asymptotic expansion for $C(x)$ of the form

$$C(x) = \Phi\left(\lambda^{\frac{1}{2}}f(r)\right) + \lambda^{-\frac{1}{2}}\phi\left(\lambda^{\frac{1}{2}}f(r)\right) \sum_{n=0}^{\infty} \lambda^{-n} a_n(r)$$

where $r = x/\lambda$ and

$$f(r) \equiv \sqrt{2(1-r+r\log r)},$$

with the sign of the square root matching the sign of $r-1$.

Temme approximation

To leading order, the quantile function is

$$Q(u) \approx \lambda r + c_0(r)$$

where

$$r = f^{-1}(w/\sqrt{\lambda}), \quad w = \Phi^{-1}(u)$$

and

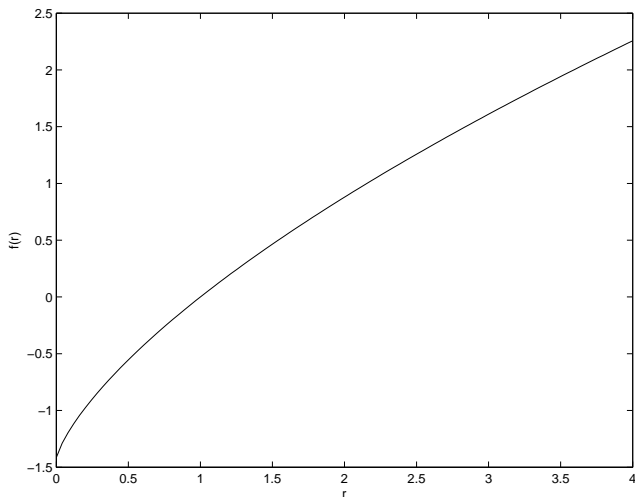
$$c_0(r) = \frac{\log(f(r)\sqrt{r}/(r-1))}{\log r}$$

The key is that both $f^{-1}(s)$ and $c_0(r)$ can be approximated very accurately by polynomials, and an additional *ad hoc* correction gives

$$\tilde{Q}_{T3}(u) = \lambda r + p_2(r) + p_3(r)/\lambda$$

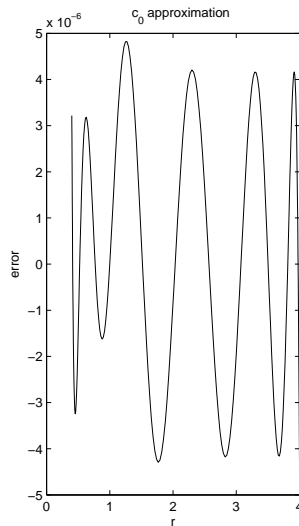
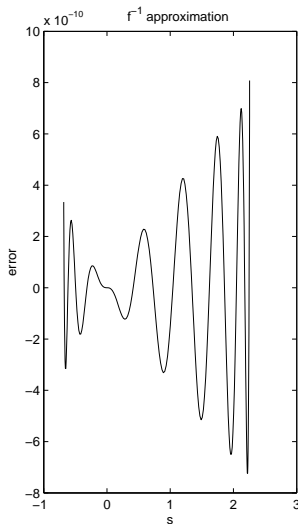
Temme approximation

The function $f(r)$



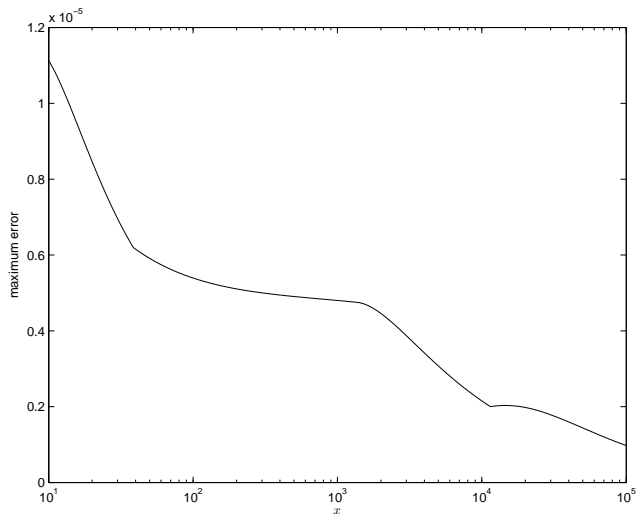
Temme approximation

Errors in $f^{-1}(s)$ and $c_0(r)$ approximations:



Temme approximation

Maximum error in \tilde{Q}_{T3} approximation:



$C(m)$ evaluation

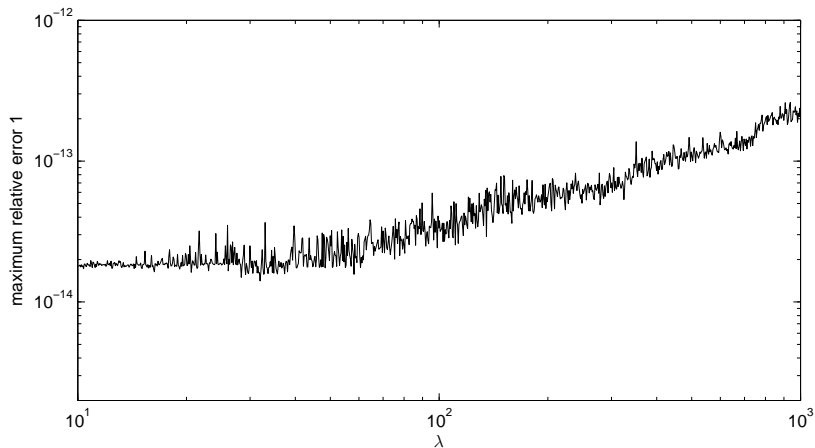
When $\tilde{Q}(u)$ is too close to an integer, we need to evaluate $C(m)$ for integer m to choose between m and $m+1$.

When $\frac{1}{2}\lambda \leq m \leq 2\lambda$, this can be done very accurately using another approximation due to Temme (1987).

Outside this range, a modified version of bottom-up / top-down summation can be used, because successive terms decrease by factor 2 or more.

Temme approximation

Maximum relative error in Temme approximation for $C(m)$



The CPU algorithm

given inputs: λ, u

if $\lambda > 2.5$

$$w := \Phi^{-1}(u)$$

if $|w| < 3$

$$x := \tilde{Q}_{N2}(w)$$

$$\delta := \lambda^{-1} \left(\frac{1}{40} + \frac{1}{80} w^2 + \frac{1}{160} w^4 \right)$$

else

$$r := f^{-1}(w/\sqrt{\lambda})$$

$$x := \lambda r + c_0(r)$$

$$x := x - (4.1/805)/(x + 0.025\lambda)$$

$$\delta := 0.01/\lambda$$

end

$$n := \lfloor x + \delta \rfloor$$

The CPU algorithm

```
if  $x > 10$ 
  if  $x - n > \delta$ 
    return  $n$ 
  else if  $C(n) < u$ 
    return  $n$ 
  else
    return  $n - 1$ 
  end
end
end

if  $u \leq 0.5$ 
  use bottom-up summation to determine  $n$ 
else
  use top-down summation to determine  $n$ 
end
```

The GPU algorithm

given inputs: λ , u

if $\lambda > 2.5$

$$w := \Phi^{-1}(u)$$

$$s := w/\sqrt{\lambda}$$

$$\delta := 5 \times 10^{-7} \sqrt{\lambda} |w|$$

if $s_{min} < s < s_{max}$

$$r := p_1(s)$$

$$x := \lambda r + p_2(r) + p_3(r)/\lambda$$

$$\delta := \delta + 1.2 \times 10^{-5}$$

else

$$r := f^{-1}(w/\sqrt{\lambda})$$

$$x := \lambda r + c_0(r)$$

$$x := x - (4.1/805)/(x + 0.025\lambda)$$

$$\delta := \delta + 0.01/\lambda$$

end

$$n := \lfloor x + \delta \rfloor$$

The GPU algorithm

```
if  $x > 10$ 
  if  $x - n > \delta$ 
    return  $n$ 
  else if  $C(n) < u$ 
    return  $n$ 
  else
    return  $n - 1$ 
  end
end
end
```

use bottom-up summation to determine n

if not accurate enough

use top-down summation to determine n

end

Conclusions

- By approximating the inverse incomplete Gamma function, have developed an approach for inverting the Poisson CDF for $\lambda > 2.5$
- Computational cost is roughly cost of inverse error function plus three polynomials of degree 8–12
- Slower than using rejection method for generating Poisson r.v.'s (at least on CPUs – may be competitive on GPUs) but can be used for Latin Hypercube sampling and QMC
- Open source implementation should be finished soon
- Student is starting work on the extension to the Binomial CDF using the incomplete Beta function

References

“The asymptotic expansion of the incomplete gamma functions”,
NM Temme, *SIAM Journal of Mathematical Analysis*, 10(4):757-766, 1979

“On the computation of the incomplete gamma functions for large values
of the parameters”, NM Temme, in *Algorithms for Approximation*,
Clarendon Press, New York, 1987

“Fast evaluation of the inverse Poisson cumulative distribution function”,
MBG, in preparation, 2013.

Software will be available from my homepage when it's ready:

<http://people.maths.ox.ac.uk/gilesm/>