# Computer exercise 2

## Simon Sigurdhsson, 900322–0291

This is a report on exercise 8.43 in Rice (2007), which concerns a data set of interarrival times of photons. The exercise consists of 6 parts labeled (a) through (f). The exercise has been solved using the R statistical software and the code used is included throughout the report. The preamble of this code simply loads a couple of packages along with the data set:

```
1  library(ggplot2)
2  library(tikzDevice)
3  data <- read.csv("gamma-arrivals.txt", header=FALSE,
4                   col.names=c("arrivals"))
```

Here, the `ggplot2` library is used to generate figures and the `tikzDevice` library is used to export them as TikZ code.

**a.** *Histogram of the interarrival times.*

```
7  tikz("tikz/histogram.tex", width=5, height=3)
8  p <- ggplot(data, aes(x=arrivals)) +
9    geom_histogram(binwidth=5) + xlim(0,750)
```

As shown by figure 1 on the following page, the short interarrival times are most common, while the number of occurrences of any interarrival time decreases in an exponential manner when the interrarival time gets large. This seems to fit well with a $\Gamma$ distribution with a shape parameter $k$ close to 1.

**b.** *Estimated parameters for $\Gamma(k,\theta)$.*

```
12  m1 <- mean(data$arrivals)
13  m2 <- mean(data$arrivals^2)
14  moments.k <- m1^2 / (m2-m1^2)
15  moments.theta <- (m2 - m1^2) / m1
```

The method of moments uses the statistics $\mathrm{E}[X]$ and $\mathrm{E}\left[X^2\right]$ (with point estimates $\overline{X}$ and $\overline{X^2}$) to estimate the parameters of the distribution. For the $\Gamma$
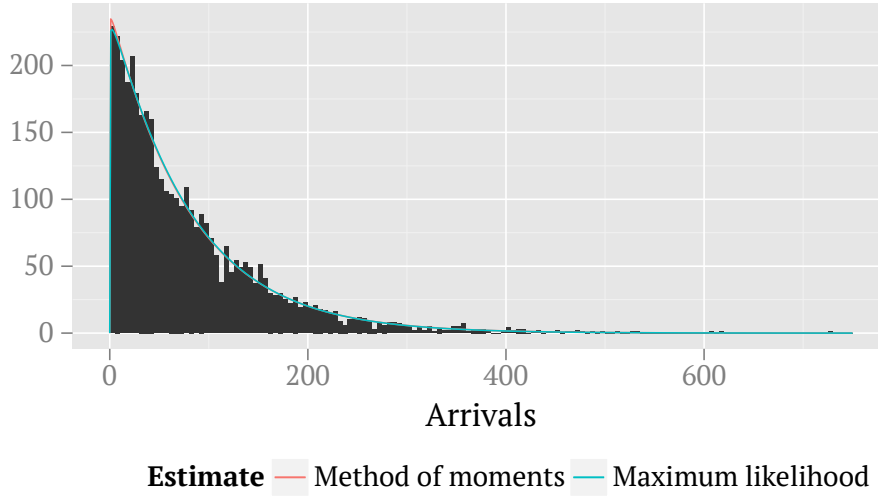
Figure 1: *Histogram of interarrival times.*

distribution, expressions for the parameters are easily derived:

$$\begin{cases} \mathrm{E}\,[X] = k\theta, \\[2mm] \mathrm{E}\,\big[X^2\big] = \theta^2 k(k+1) \end{cases} \implies \begin{cases} \hat{k} = \dfrac{\overline{X}^2}{\overline{X^2} - \overline{X}^2}, \\[4mm] \hat{\theta} = \dfrac{\overline{X^2} - \overline{X}^2}{\overline{X}} \end{cases}.$$

Using these expressions with the data set yields estimated parameter values $\hat{k}_{\mathrm{MM}} = 1.012$ and $\hat{\theta}_{\mathrm{MM}} = 78.96$, which seems to correspond to the theory that the data is $\Gamma$-distributed with $k$ close to 1.

```
18  f <- function(k) abs((log(k)-digamma(k)) -
19    (log(mean(data$arrivals))-mean(log(data$arrivals))))
20  r <- nlm(f, moments.k)
21  likelihood.k <- r$estimate
22  likelihood.theta <- mean(data$arrivals)/likelihood.k
```

The likelihood function given a $\Gamma$ distribution is $L(k,\theta) = \prod_{i=1}^{N}(\Gamma(k)\theta^k)^{-1}x_i^{k-1}\exp\left(-{^{x_i}\!/_\theta}\right)$. By instead considering the log-likelihood funcion,

$$l(k,\theta) = (k-1)\sum_{i=1}^{N}\log(x_i) - \sum_{i=1}^{N}\frac{x_i}{\theta} - Nk\log(\theta) - N\log(\Gamma(k)),$$

the resulting maximum likelihood problem is simplified. Taking the derivative with respect to each of the parameters $k$ and $\theta$ and setting it to zero (thus finding

2

the maximum) yields an estimate for $\theta$,

$$\hat{\theta} = \frac{1}{\hat{k}N} \sum_{i=1}^{N} x_i,$$

and a formula for $k$,

$$\log\left(\hat{k}\right) - \frac{\Gamma'(\hat{k})}{\Gamma(\hat{k})} = \log\left(\frac{1}{N} \sum_{i=1}^{N} x_i\right) - \frac{1}{N} \sum_{i=1}^{N} \log\left(x_i\right)$$

which has no closed-form solution but from which $\hat{k}$ can be calculated numerically.

Doing so results in estimates $\hat{k}_{\text{ML}} = 1.026$ and $\hat{\theta}_{\text{ML}} = 77.88$,

**c.** *Comparing estimated distributions.*

```
25  estimated.pdfs <- data.frame(x=0:750,
26    moments=dgamma(0:750, moments.k, scale=moments.theta),
27    likelihood=dgamma(0:750, likelihood.k, scale=likelihood.theta))
28  scalefactor = 5*length(data$arrivals) # 5 = binwidth
29  p <- p + geom_line(data=estimated.pdfs,
30    aes(x=x, y=moments*scalefactor, colour='#0072B2'))
31  p <- p + geom_line(data=estimated.pdfs,
32    aes(x=x, y=likelihood*scalefactor, colour='#D55E00'))
33  p + scale_colour_hue("Estimate",
34    labels=c('Method of moments','Maximum likelihood')) +
35    opts(legend.direction = "horizontal", legend.position = "bottom")
36  dev.off()
```

As shown by figure 1 on the previous page, the estimated distributions correspond well to the histogram.

**d.** *Estimate errors with bootstrapping.*

```
39  sample.moments <- function(){
40    samples <- rgamma(length(data$arrivals), moments.k, scale=moments.theta)
41    m1 <- mean(samples)
42    m2 <- mean(samples^2)
43    k <- m1^2 / (m2-m1^2)
44    theta <- (m2 - m1^2) / m1
45    c(k, theta)
46  }
47  simulated.moments <- t(replicate(1000, sample.moments()))
48  sample.likelihood <- function(){
49    samples <- rgamma(length(data$arrivals), likelihood.k, scale=likelihood.theta)
50    f <- function(k) abs((log(k)-digamma(k)) -
51    (log(mean(samples))-mean(log(samples))))
52    r <- nlm(f, likelihood.k)
53    k <- r$estimate
```

```
54    theta <- mean(samples)/k
55    c(k, theta)
56  }
57  simulated.likelihood <- t(replicate(1000, sample.likelihood()))
58  moments.k.error <- sd(simulated.moments[,1])
59  moments.theta.error <- sd(simulated.moments[,2])
60  likelihood.k.error <- sd(simulated.likelihood[,1])
61  likelihood.theta.error <- sd(simulated.likelihood[,2])
```

Bootstrapping the parameters by generating 1000 samples of the same size as
the initial dataset reveals the error of the four estimated parameters to be

$$s_{\hat{k}_{\mathrm{MM}}} = 0.032, \qquad\qquad s_{\hat{\theta}_{\mathrm{MM}}} = 2.83,$$
$$s_{\hat{k}_{\mathrm{ML}}} = 0.020, \qquad\qquad s_{\hat{\theta}_{\mathrm{ML}}} = 1.98.$$

As expected, the maximum likelihood estimate has a slightly smaller error.

**e.** *Confidence intervals with bootstrapping.*

```
64  cs <- unname(quantile(simulated.moments[,1], c(0.025, 0.975)))
65  moments.k.ci <- c(2*moments.k - cs[2], 2*moments.k - cs[1])
66  cs <- unname(quantile(simulated.moments[,2], c(0.025, 0.975)))
67  moments.theta.ci <- c(2*moments.theta - cs[2], 2*moments.theta - cs[1])
68  cs <- unname(quantile(simulated.likelihood[,1], c(0.025, 0.975)))
69  likelihood.k.ci <- c(2*likelihood.k - cs[2], 2*likelihood.k - cs[1])
70  cs <- unname(quantile(simulated.likelihood[,2], c(0.025, 0.975)))
71  likelihood.theta.ci <- c(2*likelihood.theta - cs[2], 2*likelihood.theta - cs[1])
```

Using the bootstrapping from (d), one can calculate confidence intervals for the
parameters. What one finds is that

$$k_{\mathrm{MM}} \in [0.948, 1.076], \qquad\qquad \theta_{\mathrm{MM}} \in [73.30, 84.39],$$
$$k_{\mathrm{ML}} \in [0.983, 1.062], \qquad\qquad \theta_{\mathrm{ML}} \in [74.22, 81.67].$$

Again, quite expectedly, the maximum likelihood estimate has a tighter confid-
ence interval, which of course is due to the smaller error.

**f.** *Are arrival times a Poisson process?* Since all operations above indicate that $k$
is close to 1, let us assume that we have $k = 1$. Since it is known (and easily
provable) that $\Gamma(1,\theta) \sim \mathrm{Exp}\left(1/\theta\right)$, and since the exponential distribution in fact
describes the time between events in a Poisson process (this can be derived from
the definition of a Poisson process), one can conclude that the interarrival time
distribution is fairly consistent with a Poisson process model for arrival times.

# References

Rice, John (2007). *Mathematical statistics and data analysis*. Australia Belmont, CA:
    Thompson/Brooks/Cole. ISBN: 9780495118688.