

# Bayesian Tests for Goodness of Fit

## 1. Introduction

In hydrology one often wants to know whether or not a series of observations corresponds to a series of samples from a given probability distribution. Examples are high flows, precipitation sums and so on. Traditionally such questions are phrased in terms of  $p$ -values. Given the popularity of the Bayesian approach elsewhere in hydrology, it seems natural to examine how it could be used in this context.

## 2. Conventions

Upper case  $S, T, X, Y, Z$ : random variables

$\vec{x}$ : vector

$\vec{X}$ : random vector

$\theta, \vec{\theta}$ : parameter or parameter vector of probability distribution

## 3. Problem statement

We have a random variable  $X$  distributed according to an unknown probability measure  $P$ . Our hypothesis to be tested is:  $H_0 = "P$  is a member of a family of probability measures  $\mathcal{P}"$ . We assume that to each element  $P \in \mathcal{P}$  we can assign parameters  $\vec{\theta} \in L \subset \mathbb{R}^m$  and that for each  $P$  the random variable  $X$  induces a probability density function (pdf)  $f(x | \vec{\theta})$  and a cumulative distribution function (cdf)  $F(x | \vec{\theta})$  on  $\mathbb{R}$ .

Now suppose  $\vec{Y}$  is a random vector of length  $n$  with independent identically distributed (iid) random variables, each with the same distribution as  $X$ . We wish to make a statement about our hypothesis based on a sample  $\vec{y}$  of  $\vec{Y}$ . To do this we assume we have a function  $s_n$  from  $\mathbb{R}^n \times \mathbb{R}^m$  to  $\mathbb{R}$  such that

$$\Pr(s_n(\vec{Y}, \vec{\theta}_0) > t | P \in \mathcal{P} \text{ and } \vec{\theta} = \vec{\theta}_0)$$

is decreasing for increasing  $t$ .

## 4. A practical example

Let  $H_0 = "P$  is a member of the Gumbel distribution family with location and scale parameters", so

$$F(x | \vec{\theta}) = F(x | \langle \xi, \zeta \rangle) = \exp\left(-\exp\left(-\frac{x - \xi}{\zeta}\right)\right)$$

For  $s_n$  we take the Kolmogorov-Smirnov distance  $D_{KS}$  between the empirical cdf (ecdf) and the Gumbel cdf.

$$D_{KS}(\vec{y}, \langle \xi, \zeta \rangle) = \sup_x |F(x | \langle \xi, \zeta \rangle) - F_n(x, \vec{y})|$$

where

$$F_n(x, \vec{y}) = \frac{\text{number of } y_i \text{ smaller than } x}{n}$$

We used R 2.15.0 (R Development Core Team, 2012) and the additional package `evd` (Stephenson, 2002) to perform the experiments.

### 4.1 Classical $p$ -value

Classical  $p$ -value for known parameters

$$\Pr(s_n(\vec{Y}, \langle 0, 1 \rangle) > s_n(\vec{y}, \langle 0, 1 \rangle) | H_0 \text{ and } \xi = 0, \zeta = 1)$$

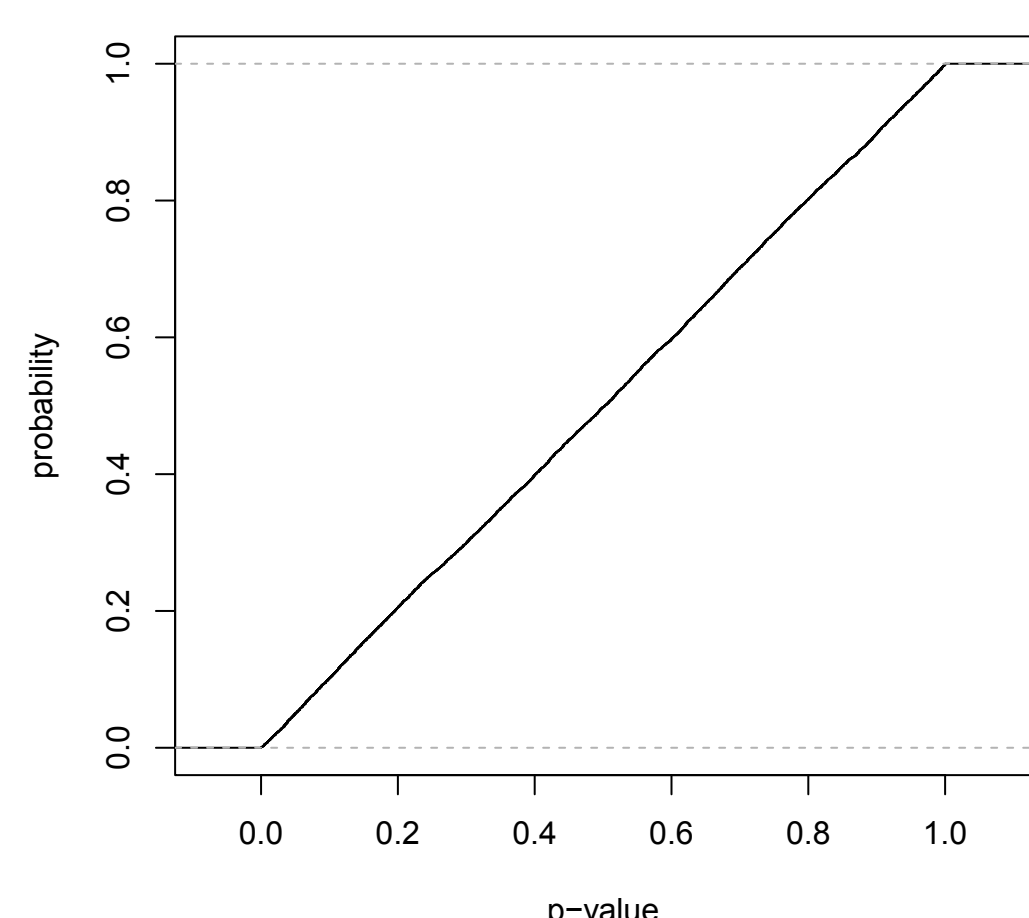


Figure 1: ecdf of  $p$ -values Gumbel with  $\langle \xi, \zeta \rangle = \langle 0, 1 \rangle$  (sample of 10000 vectors with  $n = 32$ )

Plug-in  $p$ -value. Let  $\langle \hat{\xi}, \hat{\zeta} \rangle$  be the maximum likelihood estimate of the parameters based on the sample  $\vec{y}$ .

$$\Pr(s_n(\vec{Y}, \langle \hat{\xi}, \hat{\zeta} \rangle) > s_n(\vec{y}, \langle \hat{\xi}, \hat{\zeta} \rangle) | H_0, \xi = \hat{\xi}, \zeta = \hat{\zeta})$$

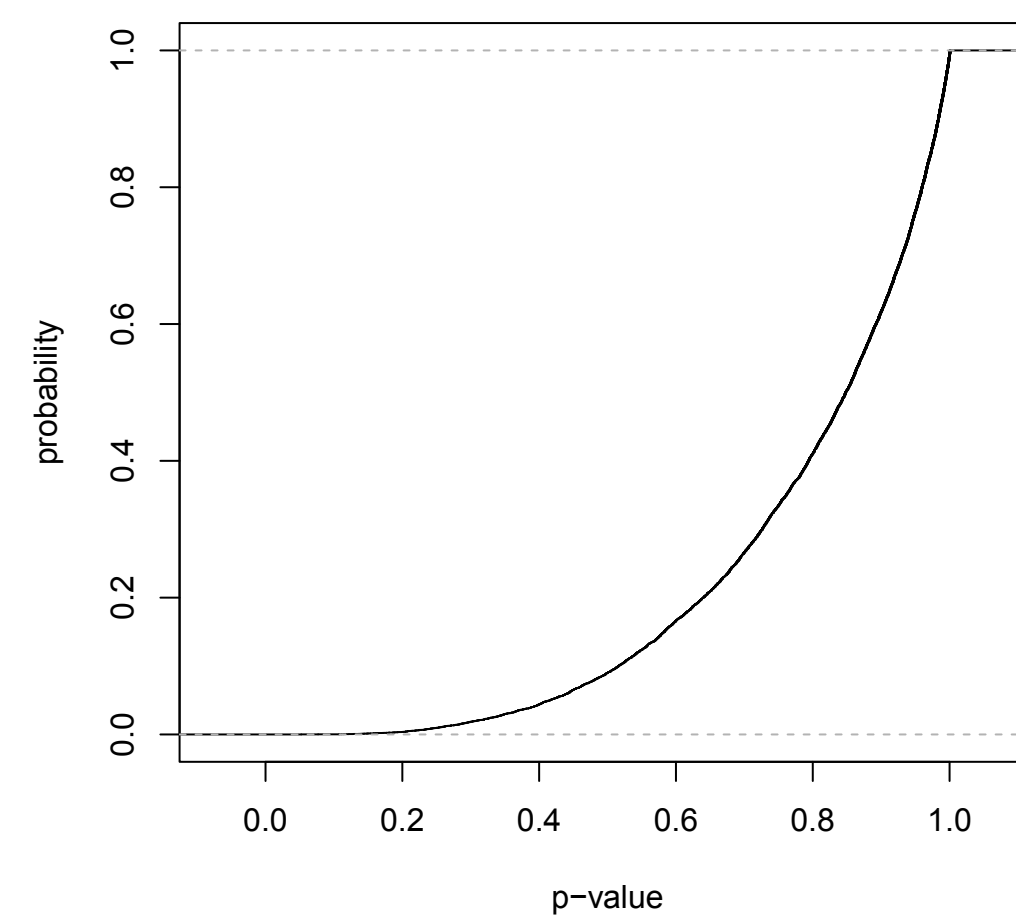


Figure 2: ecdf of plug-in  $p$ -values for samples from Gumbel with  $\langle \xi, \zeta \rangle = \langle 0, 1 \rangle$  (sample of 10000 vectors with  $n = 32$ )

The problem (the  $p$ -value is no longer equal to the probability of its occurrence) is caused by the double use of the sample. Instead of using the distribution of

$$s_n(\vec{Y}, \langle \xi(\vec{y}), \zeta(\vec{y}) \rangle)$$

we need find the distribution of

$$s_n(\vec{Y}, \langle \hat{\xi}(\vec{Y}), \hat{\zeta}(\vec{Y}) \rangle)$$

One way to recalibrate: Let our original sample be  $\vec{y}_0$ . Determine  $\hat{\xi}_0 = \hat{\xi}(\vec{y}_0)$ ,  $\hat{\zeta}_0 = \hat{\zeta}(\vec{y}_0)$ . Calculate  $d_0 = s_n(\vec{y}_0, \langle \hat{\xi}_0, \hat{\zeta}_0 \rangle)$ . Next generate samples  $\vec{y}_k$  with  $k = 1, 2, \dots, K$  from Gumbel with  $\xi = \hat{\xi}_0, \zeta = \hat{\zeta}_0$  and calculate  $d_k = s_n(\vec{y}_k, \langle \hat{\xi}(\vec{y}_k), \hat{\zeta}(\vec{y}_k) \rangle)$ . The recalibrated  $p$ -value is

$$\frac{\text{number of } d_k \text{ smaller than } d_0}{K}$$

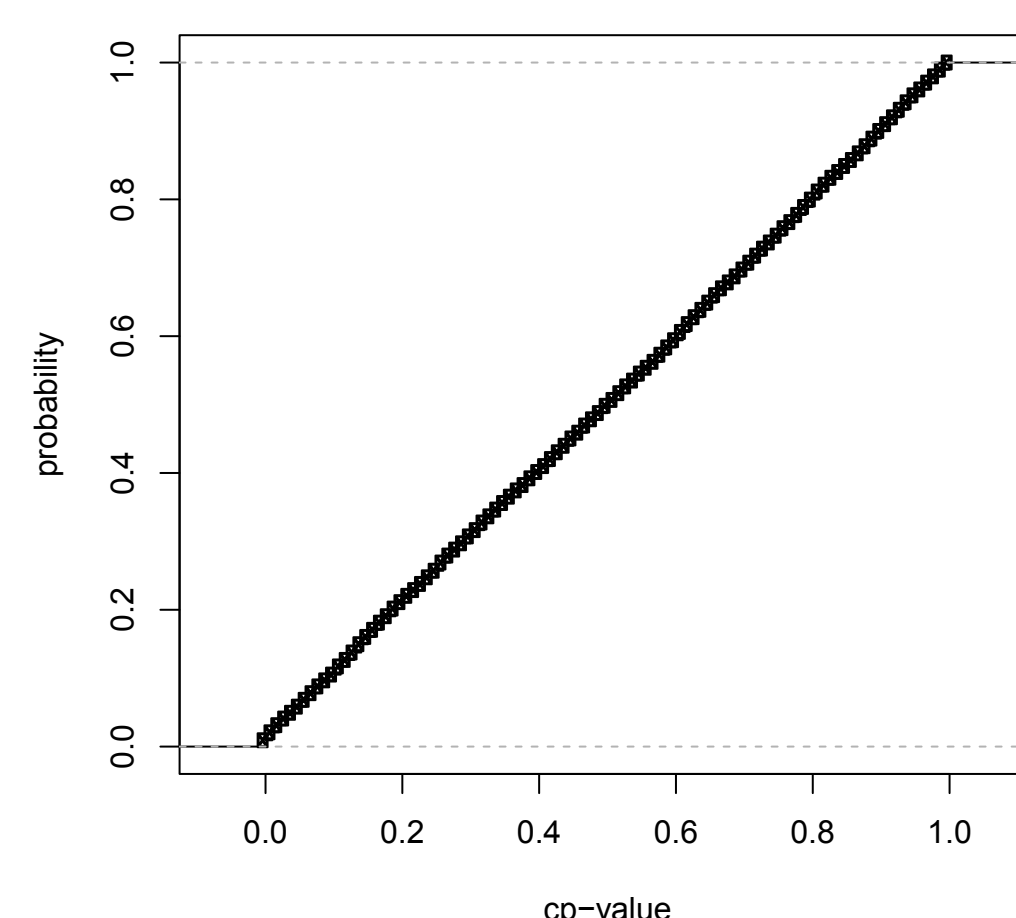


Figure 3: ecdf of recalibrated plug-in  $p$ -values for samples from Gumbel with  $\langle \xi, \zeta \rangle = \langle 0, 1 \rangle$  (sample of 10000 vectors with  $n = 32, K = 100$ )

### 4.2 Posterior predictive $p$ -value

For a sample  $\vec{y}_0$  the posterior predictive  $p$  value is defined as

$$\text{ppp}(\vec{y}_0) = \Pr(s_n(\vec{Y}, \vec{\Theta}) \geq d(\vec{y}_0, \vec{\Theta}) | \vec{y}_0)$$

where  $\langle \vec{Y}, \vec{\Theta} \rangle$  has pdf

$$f(\langle \vec{y}, \vec{\theta} \rangle | \vec{y}_0, H_0) \propto f(\vec{y} | \vec{\theta}) \pi(\vec{\theta} | \vec{y}_0)$$

and

$$\pi(\vec{\theta} | \vec{y}_0) \propto f(\vec{y}_0 | \vec{\theta}) \pi(\vec{\theta})$$

Again the sample is used twice, first in the calculation  $s_n$  and then again in the determination of the posterior distribution of  $\langle \vec{Y}, \vec{\Theta} \rangle$ .

Again we take a sample from the Gumbel distribution with parameters  $\langle 0, 1 \rangle$ . For the purpose of this example we use a proper prior,

$$\pi(\langle \xi, \zeta \rangle) = \begin{cases} 0 & \xi < -4 \text{ or } \xi > 4 \text{ or } \zeta < 0.25 \text{ or } \zeta > 4 \\ \frac{1}{8 \times 3.75 \zeta} & -4 \leq \xi \leq 4 \text{ and } 0.25 < \zeta < 4 \end{cases}$$

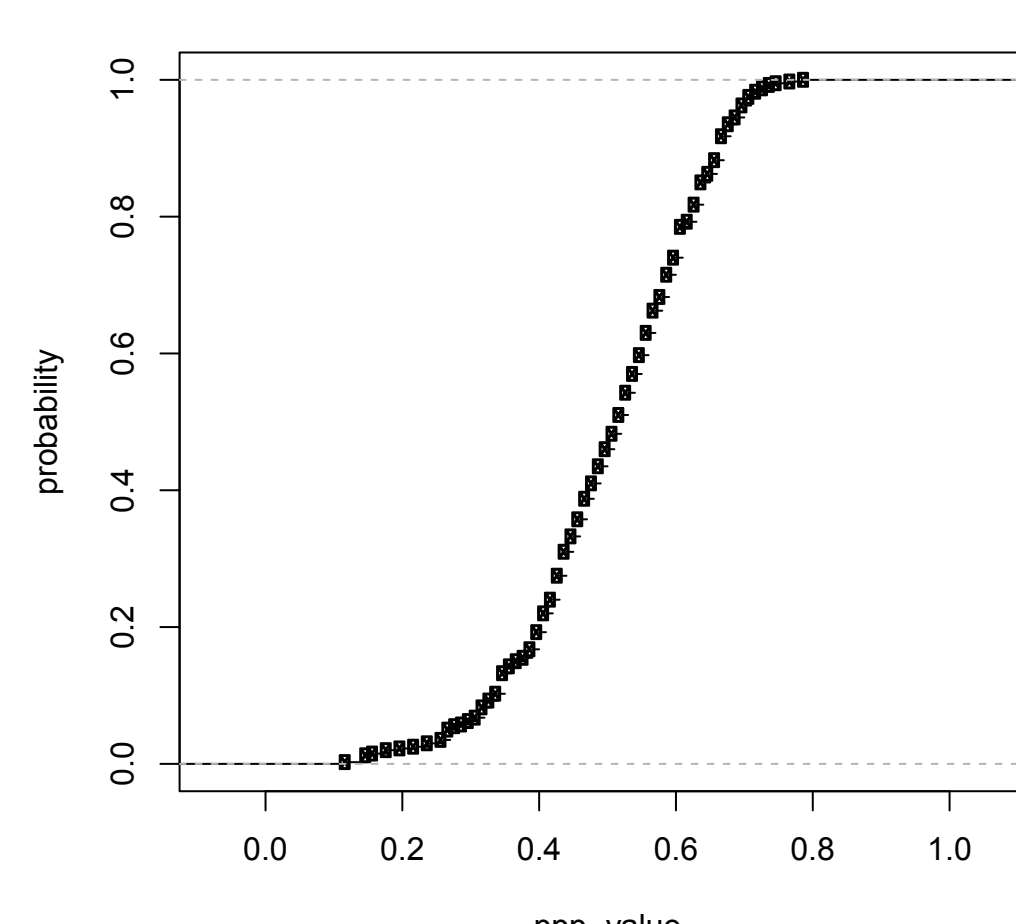


Figure 4: ecdf of ppp-values for samples from Gumbel with  $\langle \xi, \zeta \rangle = \langle 0, 1 \rangle$  (400 vectors with  $n = 32$ )

Again we recalibrate, we use the procedure outlined in (Hjort et al., 2006),

For given  $\vec{y}_0$  take  $J$  samples  $\{\vec{\theta}_j\}_{j=1}^J$  from  $\pi(\vec{\theta} | \vec{y}_0)$ .

Next for each  $j$  sample  $\vec{y}_j$  from  $f_{\vec{Y}}(\vec{y} | \vec{\theta}_j)$ .

This results in a sample  $\{\langle \vec{y}_j, \vec{\theta}_j \rangle\}_{j=1}^J$  from a distribution with density  $\propto f_{\vec{Y}}(\vec{y} | \vec{\theta}, H_0) \pi(\vec{\theta} | \vec{y}_0)$ . Define

$$p_{\text{ppp}}^{\text{MC}}(\vec{y}_0) = \frac{\text{number of } s_n(\vec{y}_j, \vec{\theta}_j) \geq s_n(\vec{y}_0, \vec{\theta}_j)}{J}$$

Algorithm 1: Posterior predictive  $p$ -value by Monte Carlo

For given  $\vec{y}_0$  determine  $p_{\text{ppp}}^{\text{MC}}(\vec{y}_0)$  (Algorithm 1).

Take  $K$  samples from  $\pi(\vec{\theta})$ , resulting in  $\{\vec{\theta}_k\}_{k=1}^K$ .

For each  $k$ , sample  $\vec{y}_k$  from  $f_{\vec{Y}}(\vec{y} | \vec{\theta}_k)$ .

Result: a sample  $\{\langle \vec{y}_k, \vec{\theta}_k \rangle\}_{k=1}^K$  from  $\propto f_{\vec{Y}}(\vec{y} | \vec{\theta}) \pi(\vec{\theta})$ .

For each  $\vec{y}_k$  determine  $p_{\text{ppp}}^{\text{MC}}(\vec{y}_k)$  (Algorithm 1) and define

$$p_{\text{cPPP}}^{\text{MC}}(\vec{y}_0) = \frac{\text{number of } p_{\text{ppp}}^{\text{MC}}(\vec{y}_k) \geq p_{\text{ppp}}^{\text{MC}}(\vec{y}_0)}{K}$$

Algorithm 2: Calibrated posterior predictive  $p$ -value by Monte Carlo

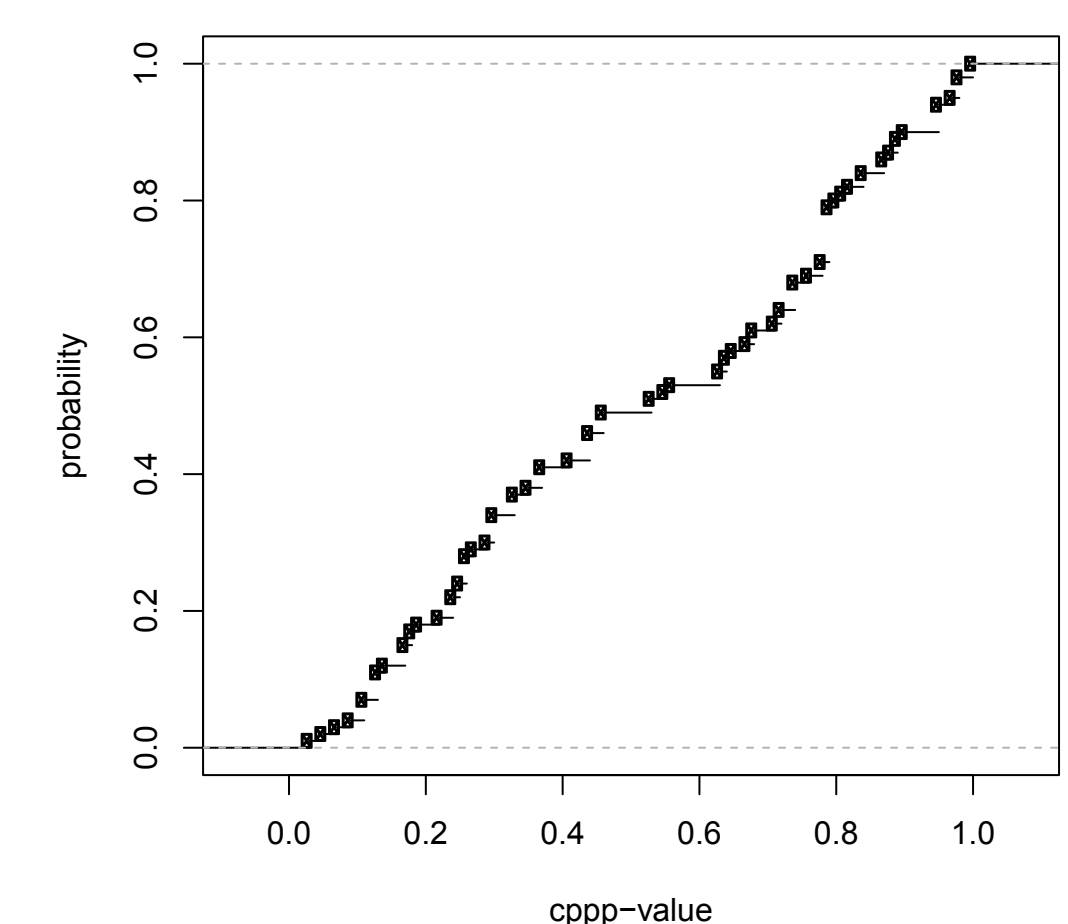


Figure 5: ecdf of cPPP-values for samples from Gumbel with  $\langle \xi, \zeta \rangle = \langle 0, 1 \rangle$  (400 vectors of length 32)

## 5. Conclusions

The literature is clear on the fact that frequentist  $p$ -values for composite hypothesis need to be interpreted carefully, see for example D'Agostino and Stephens (1986). Bayesian  $p$ -values are not as well known, but here too the literature advises caution, see also Bayarri and Berger (2000), Hjort et al. (2006). We presented an example where both the frequentist and Bayesian  $p$ -values needed an additional calibration.

## References

- M. J. Bayarri and James O. Berger.  $P$  values for composite null models. *Journal of the American Statistical Association*, 95(452):pp. 1127–1142, 2000.
- Ralph B. D'Agostino and Michael A. Stephens, editors. *Goodness-of-fit techniques*. Statistics: Textbooks and Monographs. Marcel Dekker Inc., New York, 1986.
- Nils Lid Hjort, Fredrik A. Dahl, and Gunnhildur Högnadóttir Steinbakk. Post-processing posterior predictive  $p$  values. *Journal of the American Statistical Association*, 101(475):pp. 1157–1174, 2006.
- R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2012. ISBN 3-900051-07-0.
- A. G. Stephenson. `evd`: Extreme value distributions. *R News*, 2(2), June 2002.