

# Hypothesis testing based on nonparametric entropies

October 2, 2023

Consider the problem of planning hypothesis tests based on nonparametric entropies from the functional (??) and (??): Let  $f(\cdot)$  be a density,

$$T_{\mathcal{D}}(f) = \int_{-\infty}^{\infty} f(x)\Phi(f(x))dx = \int_{-\infty}^{\infty} f(x)h\left(\frac{1}{f(x)}\right)dx, \quad (1)$$

where  $\Phi(x)$  and  $h(x)$  are real-valued functions on  $[0, \infty)$  satisfying certain regularity conditions. As an example of (??), the definition of  $\Phi(x) = h(1/x) = -\log(x)$  becomes  $T_{\mathcal{D}}(f)$  the Shannon entropy formula of a distribution on  $(-\infty, \infty)$ . From now on, we use  $T_{\mathcal{D}} := T_{\mathcal{D}}(f)$  as a notation for a possible or well-defined entropy (such as Shannon, Rényi, and Tsallis).

A consistent estimator for  $T_{\mathcal{D}}$  is given for the next class of statistics: Let  $X_{1:n} \leq X_{2:n} \leq \dots \leq X_{n:n}$  be the ordered sample from an  $n$ -point random sample and  $m$  be an integer such that  $1 \leq m \leq n$ ,

$$T_{n,m} = \frac{1}{2} \frac{1}{n-m} \sum_{j=1}^{n-m} \Phi\left(\frac{m}{n+1} [X_{j+m:n} - X_{j:n}]\right)^{-1} = \frac{1}{n} \sum_{j=1}^{n-m} h\left(\frac{n}{m} [X_{j+m:n} - X_{j:n}]\right).$$

?? have derived general asymptotic results for functions of spacings; while ?? have developed a correction for the case of Shannon entropy.]After ?? and ??, the next result applies.

**Lemma 0.1** *Suppose that  $f(\cdot)$  is a bounded density bounded away from zero and satisfies a Lipschitz condition on its support. Then the next asymptotic results follow: If  $m, n \rightarrow \infty$  and  $m = o(n^{1/2})$ , then:*

(i)

$$\sqrt{n} \left( V_{m,n}^* + \int_{-\infty}^{\infty} f(x) \log f(x) dx \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \text{Var}(\log f(X))),$$

where  $o(\cdot)$  represents the little-o and

$$V_{m,n}^* = \frac{1}{n-m} \sum_{j=1}^{n-m} \log \left( \frac{n+1}{m} [X_{j+m:n} - X_{j:n}] \right) + \sum_{k=m}^n \frac{1}{k} + \log m - \log(n+1).$$

(ii)

$$\sigma^{-1} \sqrt{n} \left( T_{m,n}^* - \int_{-\infty}^{\infty} h(1/f(x)) f(x) dx \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1),$$

where  $h(1/f(x)) = \Phi(f(x))$  provides a unification between the results of **?** and **?**,  $F(x)$  means the cumulative distribution function (cdf) with respect to  $f(x)$ ,

$$T_{n,m}^* = \frac{1}{n} \sum_{j=1}^{n-m} \left[ h \left( \frac{n}{m} [X_{j+m:n} - X_{j:n}] \right) + \frac{[(n/m)(X_{j+m:n} - X_{j:n})]^2}{2m} h'' \left( \frac{n}{m} [X_{j+m:n} - X_{j:n}] \right) \right]$$

and

$$\sigma^2 = \int_{-\infty}^{\infty} h' \left( \frac{1}{f(x)} \right) \frac{1}{f(x)} dx + \int_{-\infty}^{\infty} \left[ \int_{-\infty}^y \left( \frac{1}{f(x)} - \frac{F(x)f'(x)}{f^3(x)} \right) h' \left( \frac{1}{f(x)} \right) f(x) dx \right]^2 \frac{f(y)}{F^2(y)} dy$$

for  $h'(x)$  and  $h''(x)$  as first and second order derivatives, respectively.

Let us now consider, starting from the previous lemma, the test of the null hypothesis  $\mathcal{H}_0 : T_{\mathcal{D}} = D_0$  as opposed to one of the other three:

$$(i) \mathcal{H}_1 : T_{\mathcal{D}} \neq D_0, \quad (ii) \mathcal{H}_1 : T_{\mathcal{D}} > D_0, \quad \text{or} \quad (iii) \mathcal{H}_1 : T_{\mathcal{D}} < D_0.$$

For this purpose, we can use the test statistics:

$$Z_{m,n} = \begin{cases} \frac{\sqrt{n}(V_{m,n}^* - D_0)}{\sqrt{\text{Var}(\log f(X))}}, & \text{for } \Phi(x) = \log(x), \\ \frac{\sqrt{n}(T_{m,n}^* - D_0)}{\sigma}, & \text{for any } h(x) \text{ possible.} \end{cases}$$

So the null hypothesis should be rejected if (i)  $Z_{n,m} > z_{\alpha/2}$  or  $Z_{n,m} < -z_{\alpha/2}$  for  $\Phi_{\mathcal{N}}(z_{\alpha/2}) = 1 - \alpha/2$  and  $\Phi_{\mathcal{N}}$  being the standard normal cdf, (ii)  $Z_{n,m} > z_{\alpha/2}$  or (iii)  $Z_{n,m} < -z_{\alpha/2}$ .

The power function for case (i) (two-sided test) at  $t \neq D_0$  is given by

$$\pi_{m,n}(t) = 1 - \Phi_{m,n} \left( z_{\alpha/2} - \frac{\sqrt{n}(Z_{m,n} - D_0)}{\sigma} \right) + \Phi_{m,n} \left( -z_{\alpha/2} - \frac{\sqrt{n}(Z_{m,n} - D_0)}{\sigma} \right),$$

for a sequence of cdfs  $\Phi_{m,n}(x)$  which tends uniformly to  $\Phi_{\mathcal{N}}(x)$ .