

# Statistical Inference

Robert L. Wolpert  
Institute of Statistics and Decision Sciences  
Duke University, Durham, NC, USA

Week 12. Testing and Kullback-Liebler Divergence

## 1. Likelihood Ratios

Let  $X_1, X_2, X_2, \dots$  be independent, identically distributed random variables taking values in some space  $\mathcal{X}$ , whose distribution has a probability density function  $f(x | \theta)$ . If we believe that one of the two hypotheses

$$H_0 : [\theta = \theta_0] \qquad H_1 : [\theta = \theta_1]$$

is true but we are unsure which, we can compute the *likelihood ratio* against the null hypothesis  $H_0$ ,

$$\lambda_n = \frac{f_1(x_1)f_1(x_2) \cdots f_1(x_n)}{f_0(x_1)f_0(x_2) \cdots f_0(x_n)}$$

where we have simplified the notation by writing  $f_0(x)$  for  $f(x | \theta_0)$  and  $f_1(x)$  for  $f(x | \theta_1)$ .

A Likelihoodist Statistician would find the likelihood ratio  $\lambda_n$  to be the best direct measure of the relative support of the data for these two hypotheses; a Bayesian statistician with prior probabilities  $\pi_0 = \mathbf{P}[H_0]$  and  $\pi_1 = \mathbf{P}[H_1] = (1 - \pi_0)$  would find the best measure to be the posterior probabilities

$$\begin{aligned} \mathbf{P}[H_0 | X_n] &= \frac{\pi_0 f_0(x_1)f_0(x_2) \cdots f_0(x_n)}{\pi_0 f_0(x_1)f_0(x_2) \cdots f_0(x_n) + \pi_1 f_1(x_1)f_1(x_2) \cdots f_1(x_n)} \\ &= \frac{\pi_0}{\pi_0 + \pi_1 \lambda_n} \end{aligned}$$

or, equivalently, the posterior odds

$$\frac{P[H_1 | X_n]}{P[H_0 | \mathbf{x}_n]} = \frac{\pi_1}{\pi_0} \lambda_n$$

Finally, the Classical (or Frequentist) statistician would rely on the Neymann-Pearson Lemma, which asserts that the most powerful test of the hypothesis  $H_0$  at level  $0 < \alpha < 1$  upon observing  $\lambda_n = \lambda^*$  is to Reject  $H_0$  in favor of  $H_1$  if  $\lambda^* \geq c_\alpha$ , where

$$c_\alpha \equiv \inf \{c < \infty : P[\lambda_n \geq c | \theta_0] \leq \alpha\}$$

and otherwise fail to reject, or equivalently to reject if and only if the “P-value”

$$P_n \equiv P[\lambda_n \geq \lambda^* | \theta_0]$$

satisfies  $P_n \leq \alpha$ . This is “most powerful” in the sense that the “power”  $P[\text{Reject } H_0 | H_1] = P[\lambda_n \geq c_\alpha | \theta_1]$  is larger than for any other test of “size”  $P[\text{Reject } H_0 | H_0] = P[\lambda_n \geq c_\alpha | \theta_0] \leq \alpha$ .

Let’s explore how  $\lambda_n$  behaves as  $n \rightarrow \infty$ , and from that see what inference Likelihoodists, Bayesians, and Classical statisticians are likely to make.

### 1.1. Log Likelihood Ratio as a Random Walk

For any  $\theta \in \Theta$  the log likelihood ratio is a sum of independent random variables

$$\log \lambda_n = \sum_{j=1}^n \log [f_1(X_j)/f_0(X_j)],$$

the  $j^{\text{th}}$  of which has mean

$$\begin{aligned} \mu &= E[\log \frac{f_1(X)}{f_0(X)} | \theta] \\ &= \int_{\mathcal{X}} [\log f_1(x) - \log f_0(x)] f(x) dx \\ &= \int_{\mathcal{X}} [\log \frac{f_1(x)}{f(x)} - \log \frac{f_0(x)}{f(x)}] f(x) dx \\ &= K(f : f_0) - K(f : f_1) \end{aligned}$$

where  $K(f : g)$  is the “Kullback-Liebler Divergence” defined by

$$K(f : g) \equiv \int_{\mathcal{X}} -\log \frac{g(x)}{f(x)} f(x) dx.$$

This “divergence” satisfies  $K(f : f) = 0$  for all  $f(x)$  and  $K(f : g) > 0$  for all other  $g(x)$ , but is not symmetric and so is not (quite!) a distance metric. It has many interesting properties, some described in the course text by *Bickel & Doksum* (§2.2, 3.2) and others we will encounter below. By the Law of Large Numbers,  $(1/n) \log \lambda_n \rightarrow \mu$  as  $n \rightarrow \infty$  and hence  $\lambda_n \approx e^{n\mu}$ , so

$$\lim_{n \rightarrow \infty} \lambda_n = \lim_{n \rightarrow \infty} e^{n\mu} = \begin{cases} 0 & \text{if } K(f : f_0) < K(f : f_1), \\ \infty & \text{if } K(f : f_0) > K(f : f_1) \end{cases}$$

and, in particular, statisticians of all three paradigms will be lead to the right conclusion in the limit as  $n \rightarrow \infty$  if either  $H_0$  is true, in which case  $\mu = -K(f_0 : f_1) < 0$  and  $\lambda_n \rightarrow 0$ , or  $H_1$  is true, so  $\mu = K(f_1 : f_0) > 0$  and  $\lambda_n \rightarrow \infty$ .

Similarly we can compute

$$\sigma^2 = \mathbb{E}[(\log \frac{f_1(X)}{f_0(X)} - \mu)^2 | \theta] = \int_{\mathcal{X}} [\log f_1(x) - \log f_0(x) - \mu]^2 f(x) dx$$

and, by the Central Limit Theorem,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ \frac{\log \lambda_n - n\mu}{\sigma \sqrt{n}} \geq z \right] = \Phi(-z)$$

and hence the Likelihoodist, Bayesian, and Frequentist reports for large  $n$  and a true hypothesis  $H_0 : [\theta = \theta_0]$  will be approximately

$$\begin{aligned} \lambda_n &\approx e^{-nK(f_0:f_1)} \rightarrow 0 \\ \mathbb{P}[H_0 | x] &\approx \frac{\pi_0}{\pi_0 + \pi_1 e^{-nK(f_0:f_1)}} \rightarrow 1 \\ P_n &\approx \Phi \left( -\frac{\log \lambda_n + nK(f_0 : f_1)}{\sigma \sqrt{n}} \right) \sim \text{Un}(0, 1) \end{aligned}$$

while if  $H_1 : [\theta = \theta_1]$  is true they will be approximately

$$\begin{aligned} \lambda_n &\approx e^{nK(f_1:f_0)} \rightarrow \infty \\ \mathbb{P}[H_0 | x] &\approx \frac{\pi_0}{\pi_0 + \pi_1 e^{nK(f_1:f_0)}} \rightarrow 0 \\ P_n &\approx \Phi \left( \frac{-nK(f_1 : f_0) - nK(f_0 : f_1)}{\sigma \sqrt{n}} \right) \\ &= \Phi \left( -[K(f_1 : f_0) + K(f_0 : f_1)] \sqrt{n/\sigma^2} \right) \rightarrow 0 \end{aligned}$$

## 1.2. Martingales and Iterated Logarithms

Suppose that  $H_0 : [\theta = \theta_0]$  is true. At any fixed sample-size  $n$ , we have seen how Classical and Bayesian testing methods behave. What happens as we observe successively larger samples?

The Law of the Iterated Logarithm for random walks states that, with probability one, a random walk like  $\log \lambda_n$  with *i.i.d.* steps of mean  $\mu \in \mathbb{R}$  and variance  $\sigma^2 > 0$  will satisfy:

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{\log \lambda_n - n\mu}{\sigma \sqrt{2n \log \log n}} &= +1 \\ \liminf_{n \rightarrow \infty} \frac{\log \lambda_n - n\mu}{\sigma \sqrt{2n \log \log n}} &= -1 \end{aligned}$$

and, in particular, that if  $H_0$  is true (so that  $\mu = -K(f_0 : f_1) < 0$ ), then with probability one there will be infinitely many numbers  $n \in \mathbb{N}$  for which  $\log \lambda_n + n K(f_0 : f_1) > \sigma \sqrt{n \log \log n}$ , so that

$$\begin{aligned} P_n &\approx \Phi \left( -\frac{\log \lambda_n + n K(f_0 : f_1)}{\sigma \sqrt{n}} \right) \\ &< \Phi(-\sqrt{\log \log n}) \rightarrow 0. \end{aligned}$$

Thus the  $P$ -value will fall arbitrarily close to zero, as  $n \rightarrow \infty$ , lending arbitrarily strong evidence against even a true hypothesis. (Note:  $\sqrt{\log \log n}$  grows exceedingly slowly—for example, it only exceeds 1.645 for  $n > 3.2 \cdot 10^6$ , and only exceeds 1.96 for  $n > 1.7 \cdot 10^{20}$ ).

If  $H_0$  is true then  $\lambda_n$  is a positive martingale starting at  $\lambda_0 = 1$ , since

$$\mathbb{E}[\lambda_{n+1} \mid X_1, \dots, X_n] = \lambda_n \int_{\mathcal{X}} \frac{f_1(x)}{f_0(x)} f_0(x) dx = \lambda_n \int_{\mathcal{X}} f_1(x) dx = \lambda_n$$

and so, by Doob's Martingale Maximal Inequality, for any  $b < \infty$  and  $N > 0$ ,

$$\mathbb{P}[\sup_{n \leq N} \lambda_n > b] \leq \frac{\mathbb{E}[\lambda_0]}{b} = \frac{1}{b}$$

and hence if  $H_0$  is true the probability that the Bayesian posterior probability of  $H_0$  *ever* falls below any  $0 < p < 1$  is bounded above by

$$\begin{aligned} \mathbb{P} \left[ \inf_{n \leq N} \mathbb{P}[H_0 \mid \mathbf{x}_n] \leq p \right] &= \mathbb{P} \left[ \inf_{n \leq N} \frac{\pi_0}{\pi_0 + \pi_1 \lambda_n} \leq p \right] \\ &= \mathbb{P} \left[ \sup_{n \leq N} \lambda_n \geq \frac{\pi_0(1-p)}{\pi_1 p} \right] \leq \frac{\pi_1}{\pi_0} \frac{p}{1-p} \end{aligned}$$

as  $n \rightarrow \infty$  so, unlike Classical methods, Bayesian methods will not give arbitrarily strong evidence against a true hypothesis.

### 1.3. Martingales and Sample Size

If  $H_0$  is true we have seen that  $\lambda_n$  is a positive martingale starting at  $\lambda_0 = 1$ . For any  $0 < a < 1 < b$  let  $\tau_{ab}$  be the stopping time

$$\tau_{ab} \equiv \inf\{n \geq 0 : \lambda_n \notin (a, b)\}$$

(the first time  $\lambda_n$  either exceeds  $b > 1$  or falls below  $a < 1$ ) and denote by  $p_{ab}$  the probability that  $\lambda_n$  exceeds  $b$  before falling below  $a$ ,  $p_{ab} = \mathbb{P}[\lambda_{\tau_{ab}} \geq b]$ . Then Doob's theorem applied to the two martingales  $\lambda_n$  and  $\log \lambda_n - n\mu$  give us

$$\begin{aligned} \mathbb{E}[\lambda_{\tau_{ab}} \mid H_0] &\approx (1 - p_{ab})a + (p_{ab})b \\ &= 1 \implies \\ p_{ab} &= \frac{1 - a}{b - a} \\ \mathbb{E}[\log \lambda_{\tau_{ab}} - \mu\tau_{ab} \mid H_0] &\approx (1 - p_{ab})\log a + (p_{ab})\log b - \mu\mathbb{E}[\tau_{ab} \mid H_0] \\ &= \frac{(b - 1)\log a + (1 - a)\log b}{b - a} - \mu\mathbb{E}[\tau_{ab} \mid H_0] \\ &= 0 \implies \\ \mathbb{E}[\tau_{ab} \mid H_0] &= -\frac{(b - 1)\log a + (1 - a)\log b}{(b - a)K(f_0 : f_1)} \end{aligned}$$

and, in the limit as  $b \rightarrow \infty$ ,

$$\mathbb{E}[\inf\{n \geq 0 : \lambda_n < a\} \mid H_0] = \frac{\log(1/a)}{K(f_0 : f_1)}$$

so the sample-size needed to reach a posterior probability of  $\mathbb{P}[H_0 \mid \vec{x}_n] = \frac{\pi_0}{\pi_0 + \pi_1 \lambda_n} > 1 - \epsilon$  for a true hypothesis  $H_0 : \{x_j \sim f_0(x)\}$  has expectation

$$n \geq \frac{\log \frac{\pi_1}{\pi_0} + \log \frac{1-\epsilon}{\epsilon}}{K(f_0 : f_1)}$$

and, similarly, if  $H_0$  is false and  $x_j \sim f_1(x)$ , then

$$\mathbb{E}[\inf\{n \geq 0 : \lambda_n > b\} \mid H_1] = \frac{\log b}{K(f_1 : f_0)}$$

and the sample-size needed to achieve a posterior probability of  $P[H_0 \mid \vec{x}_n] = \frac{\pi_0}{\pi_0 + \pi_1 \lambda_n} < \epsilon$  for a false hypothesis is

$$n \geq \frac{\log \frac{\pi_0}{\pi_1} + \log \frac{1-\epsilon}{\epsilon}}{K(f_1 : f_0)}.$$

Evidently the sample-size needed varies directly with the logistic of the desired posterior probability, and inversely as the Kullback-Liebler discrepancy between  $f_0$  and  $f_1$ .

## 2. Kullback-Liebler and Fisher Information

For small  $\epsilon > 0$ , a second-order Taylor-series approximation of the KL divergence from  $f(x \mid \theta)$  to  $f(x \mid \theta + \epsilon)$  gives

$$\begin{aligned} K(f(x \mid \theta), f(x \mid \theta + \epsilon)) &= \int_{\mathcal{X}} -\log \frac{f(x \mid \theta + \epsilon)}{f(x \mid \theta)} f(x \mid \theta) dx \\ &\approx \int_{\mathcal{X}} \left[ -\epsilon \nabla \log f(x \mid \theta) - \frac{\epsilon^2}{2} \nabla^2 \log f(x \mid \theta) \right] f(x \mid \theta) dx \\ &= I(\theta) \epsilon^2 / 2 \end{aligned}$$

or, in  $q > 1$  dimensions,  $\epsilon^\top I(\theta) \epsilon / 2$ . This suggests a close link between KL divergence and the “Information Metric” notion of the distance between different distributions,

$$\begin{aligned} d_I(\theta_0, \theta_1) &= \int_{\theta_0}^{\theta_1} \sqrt{I(\theta)} d\theta \quad (\text{in } q = 1 \text{ dimension}) \\ &= \inf_{\gamma} \int_0^1 \sqrt{\dot{\gamma}_t^\top I(\gamma_t) \dot{\gamma}_t} dt \end{aligned}$$

where the infimum is over all differentiable paths  $\gamma$  from  $\gamma_0 = \theta_0$  to  $\gamma_1 = \theta_1$ ; evidently, for near-by  $\theta_0, \theta_1$ ,

$$\begin{aligned} K(f(x \mid \theta_0), f(x \mid \theta_1)) &\approx d_I(\theta_0, \theta_1)^2 / 2 \\ &\approx \frac{1}{2} (\theta_1 - \theta_0)' I(\theta) (\theta_1 - \theta_0). \end{aligned}$$