

Topic 5: Generalized Likelihood Ratio Test

### 5.1 Introduction

We already started studying *composite* hypothesis problems of the form:

$$\begin{aligned} H_0 : x &\sim p_0(x|\theta_0^*), \quad \theta_0^* \in \Theta_0, \\ H_1 : x &\sim p_1(x|\theta_1^*), \quad \theta_1^* \in \Theta_1. \end{aligned} \tag{5.1}$$

where the parameters  $\theta_0^*$  and/or  $\theta_1^*$  are unknown.

**Example 5.1.** In Example 3.3 we wanted to determine whether two meteorites came from the same asteroid in space using the difference  $x$  of their magnesium composition. In Example 3.4 we wanted to determine whether a certain gene was associated with a disease, using the difference  $x$  of the gene's average activation levels between healthy and sick people. Both of these problems can be modeled as

$$\begin{aligned} H_0 : x &\sim \mathcal{N}(0, \sigma^2) && \Rightarrow \text{same asteroid/gene unrelated to disease,} \\ H_1 : x &\sim \mathcal{N}(\mu_1^*, \sigma^2) && \Rightarrow \text{different asteroids/gene related to disease,} \end{aligned}$$

where  $\mu_1^*$  is unknown.

**Example 5.2.** In Example 4.1 we use an array of sensors to record small voltages generated by your brain and store them in a *signal* vector  $\mathbf{x} \in \mathbb{R}^D$ . A machine (phone, computer, server, etc.) should

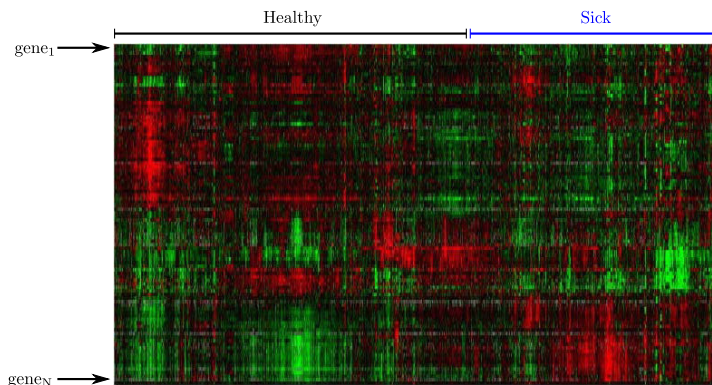


Figure 5.1: Gene microarrays are data matrices indicating gene *activation levels*. Each row corresponds to one gene, and each column corresponds to one individual. We want to know which genes are related to a disease. See Example 5.1.

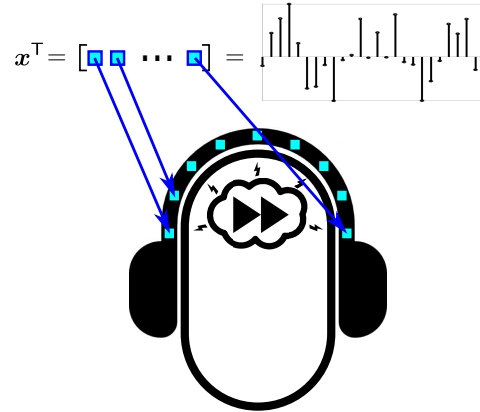


Figure 5.2: Sensors record small voltages generated by your brain and store them in a *signal* vector  $\mathbf{x} \in \mathbb{R}^D$ . A machine (phone, computer, server, etc.) should interpret  $\mathbf{x}$  and *skip* the song if that is what you thought about. See Example 5.2.

interpret  $\mathbf{x}$  and *skip* the song if that is what you thought about. This can be modeled as

$$\begin{aligned}
 H_0 : \mathbf{x} &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) && \Rightarrow \text{do nothing,} \\
 H_1 : \mathbf{x} &\sim \mathcal{N}(\boldsymbol{\mu}_\star, \sigma^2 \mathbf{I}) && \Rightarrow \text{skip song,}
 \end{aligned}$$

where  $\boldsymbol{\mu}_\star$  is unknown.

In Section 3.7 we already studied the hypothesis problems in Example 5.1, which led us to the intuitive answer of Wald’s test in Example 3.14. We now generalize this using our knowledge from estimation theory to obtain the *generalized likelihood ratio test* (GLRT).

**Definition 5.1** (Generalized likelihood ratio test (GLRT)). Consider a hypothesis problem as in (5.1), where  $\theta_0^\star$  and/or  $\theta_1^\star$  are unknown. The *generalized likelihood ratio statistic* is defined as:

$$\hat{\Lambda}(\mathbf{x}) := \frac{\max_{\theta_1 \in \Theta_1} p_1(\mathbf{x}|\theta_1)}{\max_{\theta_0 \in \Theta_0} p_0(\mathbf{x}|\theta_0)},$$

and the *generalized likelihood ratio test* is defined as  $\hat{\Lambda}(\mathbf{x}) \underset{H_0}{\overset{H_1}{\gtrless}} \tau$ .

The idea behind the GLRT is actually quite simple: if you don’t know a parameter, then first estimate it using maximum likelihood, and then use it as if it were the *true* parameter in a *likelihood ratio test* (LRT).

**Example 5.3** (Derivation of Wald’s test as a GLRT). Let us show that Wald’s test is just one particular case of the GLRT. Consider The setup in Example 5.1. Here  $\theta_0$  is known, and  $\theta_1^\star = \mu_1^\star$ , so  $\Theta_1 = \mathbb{R}$ .

Then

$$\begin{aligned} \arg \max_{\theta_1 \in \Theta_1} p_1(x|\theta_1) &= \arg \max_{\mu_1 \in \mathbb{R}} p_1(x|\mu_1) = \arg \max_{\mu_1 \in \mathbb{R}} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma}\right)^2} \\ &= \arg \max_{\mu_1 \in \mathbb{R}} \log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma}\right)^2} \right) = \arg \max_{\mu_1 \in \mathbb{R}} \underbrace{\log \frac{1}{\sqrt{2\pi}\sigma}}_{\text{constant}} - \frac{1}{2}\left(\frac{x-\mu_1}{\sigma}\right)^2 \\ &= \arg \max_{\mu_1 \in \mathbb{R}} - \underbrace{\frac{1}{2\sigma^2}}_{\text{constant}} (x-\mu_1)^2 = \arg \max_{\mu_1 \in \mathbb{R}} (\mu_1 - x)^2 \end{aligned}$$

To find this maximum we use our usual tricks: take derivative with respect to (w.r.t.)  $\mu_1$ , set to zero and solve for  $\mu_1$ .

$$\frac{\partial}{\partial \mu_1} (\mu_1 - x)^2 = 2(\mu_1 - x) = 0$$

Which yields  $\hat{\mu}_1 = x$ . Then

$$\hat{\Lambda}(x) = \frac{\max_{\theta_1 \in \Theta_1} p_1(x|\theta_1)}{\max_{\theta_0 \in \Theta_0} p_0(x|\theta_0)} = \frac{\max_{\mu_1 \in \mathbb{R}} p_1(x|\mu_1)}{p_0(x|\theta_0^*)} = \frac{p_1(x|\hat{\mu}_1)}{p_0(x|\theta_0^*)} = \frac{\cancel{\frac{1}{\sqrt{2\pi}\sigma}} e^{-\frac{1}{2}\left(\frac{x-\hat{\mu}_1}{\sigma}\right)^2}}{\cancel{\frac{1}{\sqrt{2\pi}\sigma}} e^{-\frac{1}{2}\left(\frac{x}{\sigma}\right)^2}} = \frac{e^{-\frac{1}{2}\left(\frac{x-x}{\sigma}\right)^2}}{e^{-\frac{1}{2}\left(\frac{x}{\sigma}\right)^2}} = e^{\frac{x^2}{2\sigma^2}},$$

and our GLRT is  $e^{\frac{x^2}{2\sigma^2}} \underset{H_0}{\overset{H_1}{\geq}} \tau$ . Taking log on both sides this becomes  $\frac{x^2}{2\sigma^2} \underset{H_0}{\overset{H_1}{\geq}} \log \tau$ , or equivalently

$$|x| \underset{H_0}{\overset{H_1}{\geq}} \tau',$$

where  $\tau' = \sqrt{2\sigma^2 \log \tau}$ , which is precisely Wald's test from Example 3.14.

## 5.2 Asymptotics

Recall that in hypothesis testing we often want to bound the probability  $p_{10}$  of deciding  $H_1$  when  $H_0$  was true, and so we select our threshold test  $\tau'$  accordingly. For instance, in Wald's test,  $p_{10} = P(|x| > \tau')$  given that  $H_0$  is true. Since  $H_0 : x \sim \mathcal{N}(0, \sigma^2)$ , if we want  $p_{10} < \alpha$ , all we need to do is find the  $\tau'$  such that the probability of the tails of a  $\mathcal{N}(0, \sigma^2)$  is  $\alpha$  (see Figure 5.3 to build some intuition). We can do this because we know the distribution of  $|x|$ . Similarly, if we had a test like:

$$\sum_{i=1}^N x_i^2 \underset{H_0}{\overset{H_1}{\geq}} \tau', \tag{5.2}$$

where  $H_0 : x_i \overset{iid}{\sim} \mathcal{N}(0, 1)$ , we would also know that  $\sum_{i=1}^N x_i^2 \sim \chi^2(N)$ . Again, if we wanted  $p_{10} < \alpha$ , all we would need to do is find the  $\tau'$  such that the probability of the tail of a  $\chi^2(N)$  is  $\alpha$  (see Figure 5.3 for more intuition).

However, we don't always know the distribution of our test. For example, if we had the same test in (5.2), but with  $H_0 : x_i \overset{iid}{\sim}$  Poisson, or Cauchy, or Weibull, then what is the distribution of  $\sum_{i=1}^N x_i^2$ ? The following theorem states that if  $N$  is large enough, we don't need to worry about it, because under  $H_0$ , the GLRT is

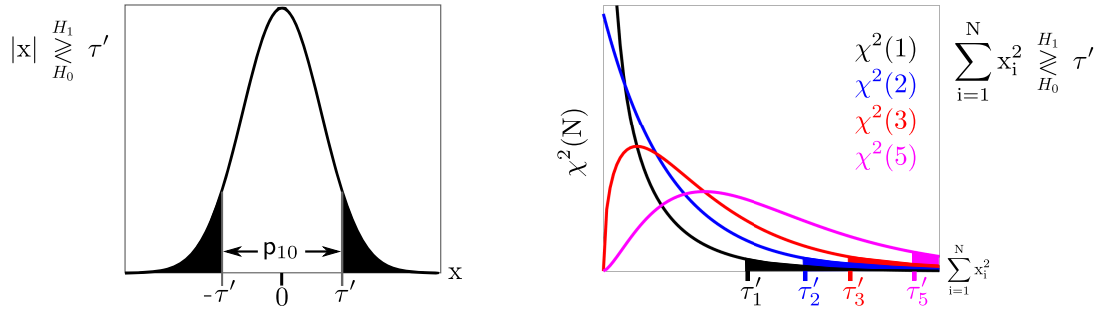


Figure 5.3: We want our test to have  $p_{10} < \alpha$ . **Left:** If our test is  $|x| \geq_{H_0}^{H_1} \tau'$  and  $x \sim \mathcal{N}(0, \sigma^2)$ , then we need to find the  $\tau'$  such that the probability of the tails of a  $\mathcal{N}(0, \sigma^2)$  is  $\alpha$ . **Right:** if our test is  $\sum_{i=1}^N x_i^2 \geq_{H_0}^{H_1} \tau'$  and  $x_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ , then we need to find the  $\tau'$  such that the probability that the tail of a  $\chi^2(N)$  is  $\alpha$ . We can do this because we know the distribution of our test. What happens if we don't? Wilks' Theorem gives us an answer.

asymptotically distributed  $\chi^2(N)$ .

**Theorem 5.1** (Wilks' Theorem — Asymptotic distribution of the GLRT). Consider a composite hypothesis problem of the form:

$$\begin{aligned}
 H_0 : x_1, \dots, x_N &\stackrel{iid}{\sim} p(x|\theta_0^*), \quad \theta_0^* \in \mathbb{R}^{K_0}, \\
 H_1 : x_1, \dots, x_N &\stackrel{iid}{\sim} p(x|\theta_1^*), \quad \theta_1^* \in \mathbb{R}^{K_1},
 \end{aligned}$$

where  $p$  has the same form in both hypotheses, and the  $K_0$  unknown parameters in  $\theta_0^*$  are a subset of the  $K_1$  unknown parameters in  $\theta_1^*$  (these are called *nested hypotheses*). Suppose that for every  $k, \ell \in \{1, \dots, K\}$ ,  $\frac{\partial p(x|\theta)}{\partial \theta_k}$  and  $\frac{\partial^2 p(x|\theta)}{\partial \theta_k \partial \theta_\ell}$  exist, and that  $E \left[ \frac{\partial \log p(x|\theta)}{\partial \theta_k} \right] = 0$  (this guarantees that the MLE  $\hat{\theta}_{ML}$  converges to the true  $\theta^*$ ). Then under  $H_0$ ,

$$2 \log \hat{\Lambda}(x) \stackrel{N \rightarrow \infty}{\rightsquigarrow} \chi^2(K_1 - K_0).$$

*Proof.* The proof follows as a consequence of the central limit theorem. For a detailed proof see Theorems 10.3.1 and 10.3.3 in *Statistical Inference* by George Casella and Roger L. Berger, second edition.  $\square$

**Example 5.4.** Consider

$$\begin{aligned}
 H_0 : x_1, \dots, x_N &\stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma_0), \\
 H_1 : x_1, \dots, x_N &\stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma_1^{*2}), \quad \sigma_1^* > 0.
 \end{aligned}$$

The MLE of  $\sigma_1^*$  is

$$\begin{aligned} \arg \max_{\theta_1 \in \Theta_1} p_1(\mathbf{x}|\theta_1) &= \arg \max_{\sigma_1 \in \mathbb{R}_+} p_1(\mathbf{x}|\sigma_1) = \arg \max_{\sigma_1 \in \mathbb{R}_+} \frac{1}{(\sqrt{2\pi}\sigma_1)^N} e^{-\frac{1}{2\sigma_1^2}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1})} \\ &= \arg \max_{\sigma_1 \in \mathbb{R}_+} \log \left( \frac{1}{(\sqrt{2\pi}\sigma_1)^N} e^{-\frac{1}{2\sigma_1^2}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1})} \right) \\ &= \arg \max_{\sigma_1 \in \mathbb{R}_+} \underbrace{\log \frac{1}{(\sqrt{2\pi})^N}}_{\text{constant}} - \log \sigma_1^N + \log e^{-\frac{1}{2\sigma_1^2}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1})} \\ &= \arg \max_{\sigma_1 \in \mathbb{R}_+} -\frac{N}{2} \log \sigma_1^2 - \frac{1}{2\sigma_1^2}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1}) \end{aligned}$$

Taking derivative w.r.t.  $\sigma_1^2$  and setting to zero we have:

$$\frac{1}{\sigma_1^4}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1}) - \frac{N}{\sigma_1^2} = \frac{1}{\sigma_1^2} \left( \frac{1}{\sigma_1^2}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1}) - N \right) = 0.$$

Solving for  $\sigma_1^2$  we obtain the MLE:

$$\hat{\sigma}_1^2 = \frac{1}{N}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1}) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2.$$

Then

$$\begin{aligned} \hat{\Lambda}(\mathbf{x}) &= \frac{\max_{\theta_1 \in \Theta_1} p_1(\mathbf{x}|\theta_1)}{\max_{\theta_0 \in \Theta_0} p_0(\mathbf{x}|\theta_0)} = \frac{\max_{\sigma_1 \in \mathbb{R}_+} p_1(\mathbf{x}|\sigma_1)}{p_0(\mathbf{x}|\sigma_0^*)} = \frac{p_1(\mathbf{x}|\hat{\sigma}_1)}{p_0(\mathbf{x}|\sigma_0^*)} = \frac{\frac{1}{(\sqrt{2\pi}\hat{\sigma}_1)^N} e^{-\frac{1}{2\hat{\sigma}_1^2}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1})}}{\frac{1}{(\sqrt{2\pi}\sigma_0^*)^N} e^{-\frac{1}{2\sigma_0^{*2}}(\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1})}} \\ &= \left( \frac{\sigma_0^*}{\hat{\sigma}_1} \right)^N e^{\left( \frac{1}{2\sigma_0^{*2}} - \frac{1}{2\hat{\sigma}_1^2} \right) (\mathbf{x}-\mu\mathbf{1})^\top(\mathbf{x}-\mu\mathbf{1})} = \left( \frac{\sigma_0^*}{\hat{\sigma}_1} \right)^N e^{\frac{N}{2} \left( \frac{\hat{\sigma}_1^2}{\sigma_0^{*2}} - 1 \right)}. \end{aligned}$$

Taking log on both sides we obtain:

$$\log \hat{\Lambda}(\mathbf{x}) = -\underbrace{\frac{N}{2} \log \left( \frac{\hat{\sigma}_1}{\sigma_0^*} \right)^2}_{H_0: \sim \chi^2(N)} + \underbrace{\frac{N}{2} \left( \frac{\hat{\sigma}_1}{\sigma_0^*} \right)^2}_{H_0: \sim \chi^2(N)} - \frac{N}{2}.$$

Now the question is: do you know the distribution of  $-\log \chi^2(N) + \chi^2(N)$ ? Sure as hell I don't! Luckily, there are no unknown parameters under  $H_0$ , and one unknown parameter under  $H_1$  (namely  $\sigma_1^*$ ), so these are nested hypotheses, and hence can use Theorem 5.1 to know that under  $H_0$ ,

$$2 \log \hat{\Lambda}(\mathbf{x}) \stackrel{N \rightarrow \infty}{\sim} \chi^2(1).$$

Hence, if we want  $p_{10} < \alpha$ , it suffices to select the threshold  $\tau$  for which the tail probability of a  $\chi^2(1)$  random variable is  $\alpha$ .

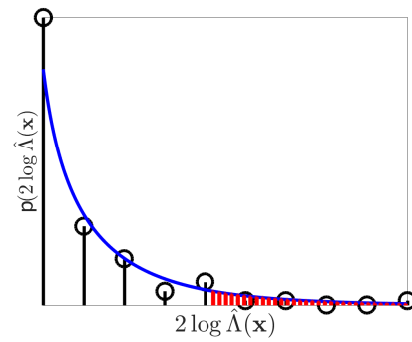


Figure 5.4: Results of the simulation of Example 5.4. In **black** is the histogram of  $2 \log \hat{\Lambda}(\mathbf{x})$  over  $T = 100$  trials. In **blue** is the asymptotic distribution of  $2 \log \hat{\Lambda}(\mathbf{x})$ . In **red** is the probability of error  $p_{10}$  set to  $\alpha = 0.05$ . The code for this simulation is in the appendix.

### 5.3 Experiments

Let us now run some simulations to verify Theorem 5.1 and our results from Example 5.4. We will generate  $x_1, \dots, x_N \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma_0^{*2})$  and we will compute  $2 \log \hat{\Lambda}(\mathbf{x})$  as in Example 5.4. We will repeat this  $T$  trials, plot a histogram, and compare it with the  $\chi^2(1)$  distribution predicted by Theorem 5.1. Figure 5.4 shows some results. The code for this simulation is in the appendix.

## A Code for Simulation of Example 5.4

```

1 clear all; close all; clc; warning('off','all');
2
3 % === Code to simulate the asymptotic distribution of log Lambda.hat ===
4 % ===== See Example 5.4.
5
6 T = 100;           % Number of trials.
7 N = 100;          % Number of samples.
8 mu = 5;           % Known mean under both hypotheses.
9 sigma0_star = 1;  % True standard deviation of X.i under H0.
10 sigma1_star = 2;  % True standard deviation of X.i under H1.
11 alpha = 0.05;     % Allowed probability of error (1|0).
12
13 log_Lambda_hat = zeros(T,1); %log_Lambda_hat over trials.
14 for t=1:T,
15     X = (randn(N,1)+mu) * sigma0_star;
16
17     log_Lambda_hat(t) = -N/2 * log( 1/N * sum(((X-mu)/sigma0_star).^2) ) ...
18         + 1/2 * sum(((X-mu)/sigma0_star).^2) ...
19         - N/2;
20 end
21
22 % Plot results.
23 figure(1);
24 axes('Box','on');
25 hold on;
26 [h,rangeh] = hist(2*log_Lambda_hat);
27 stem(rangeh,h/T,'k','LineWidth',4,'MarkerSize',20);
28
29 % Asymptotic distribution of 2*log Lambda.hat.
30 rangeh = min(rangeh):0.01:max(rangeh);
31 chi2 = chi2pdf(rangeh,1);
32 plot(rangeh,chi2,'b','LineWidth',4);
33
34 % Highlight the probability alpha.
35 tau = chi2inv(1-alpha,1);
36 rangeh = tau:0.125:max(rangeh);
37 chi2 = chi2pdf(rangeh,1);
38 stem(rangeh,chi2,'r','LineWidth',5,'MarkerSize',1);
39
40 % Make figure look sexy.
41 axis tight;
42 set(gca,'XTick',[],'xticklabel',[]);
43 set(gca,'YTick',[],'yticklabel',[]);
44 ylabel('$\mathsf{p}(2\log\hat{\Lambda})(\mathbf{x})$', 'Interpreter','latex','fontsize',25);
45 xlabel('$2\log\hat{\Lambda}(\mathbf{x})$', 'Interpreter','latex','fontsize',30);
46
47 % Save figure.
48 set(gcf, 'renderer','default');
49 figurename = 'wilks.pdf';
50 saveas(gcf,figurename);

```