

## Wald, Score and Likelihood Ratio tests

Consider a scalar or vector parameter  $\theta$  indexing a family of distributions and hence log-likelihood  $L(\theta)$ . The (log-)likelihood-ratio test (LRT) of  $\theta = \theta_0$  against unrestricted  $\theta \in \Theta$  is  $L(\hat{\theta}) - L(\theta_0)$  where  $\hat{\theta}$  is the MLE. Under regularity conditions (which include that  $\theta_0$  is not a boundary value in  $\Theta$ ), the twice the LRT is asymptotically distributed as  $\chi_q^2$ , where  $q$  is the number of restrictions imposed. (Here  $\theta_0$  can either be some fixed value or the MLE in  $\Theta_0 \subset \Theta$ .)

The LRT requires us to maximize the log-likelihood and also to evaluate it at  $\theta_0$ . If we wish to consider many possible restrictions (for example to test if many individual parameters are zero) the latter can be inconvenient and avoided by using Wald's approximation, which is to replace the log-likelihood  $L$  by a quadratic approximation  $L_W$  which agrees with  $L$  and its first and second derivatives at  $\hat{\theta}$ . Since  $\hat{\theta}$  is the MLE, the first derivative will be zero, and hence

$$L_W(\theta) = L(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^T H(\hat{\theta})(\theta - \hat{\theta})$$

where  $H$  is the Hessian<sup>1</sup> of  $L$ . This is the same log-likelihood we would obtain if we assumed that  $\hat{\theta} \sim N(\theta, J^{-1})$  where  $J = -H$  is the *observed information matrix*. Further, we approximate  $J$  by  $I(\hat{\theta})$ , the Fisher information at the MLE: this is what is most commonly known as the Wald test. (Occasionally  $I(\theta_0)$  is used.)

We can use a Wald test either by assuming that  $\hat{\theta} \sim N(\theta, I(\hat{\theta})^{-1})$  or via using one of the approximations as if it were the actual log-likelihood. For example, to test if a single parameter is zero via assuming the  $t$ -ratio is  $N(0, 1)$  is an application of a Wald test.

Alternatively, we might want to consider many extensions to a model, for example adding one out of many new explanatory variables. In that case we may wish to avoid maximizing  $L$  under each extended model. This is what Rao's score tests do, by approximating the log-likelihood from its shape at  $\theta_0$ . The usual form is to consider the score  $U$ , the derivative of  $L$  with respect to  $\theta$  at  $\theta_0$ : if  $\theta_0$  were the MLE this would be zero. To see how non-zero it is we compare it to a reference normal distribution, and under the null hypothesis this has covariance matrix the Fisher information  $I(\theta_0)$ . Thus the score statistic is  $U^T I(\theta_0)^{-1} U$ .

Score tests can also be derived by approximating  $L$  by fitting the quadratic which agrees with  $L$  and its first two derivatives at  $\theta_0$ , that is

$$L_S(\theta) = L(\theta_0) + (\theta - \theta_0)^T L'(\theta_0) + \frac{1}{2}(\theta - \theta_0)^T H(\theta_0)(\theta - \theta_0)$$

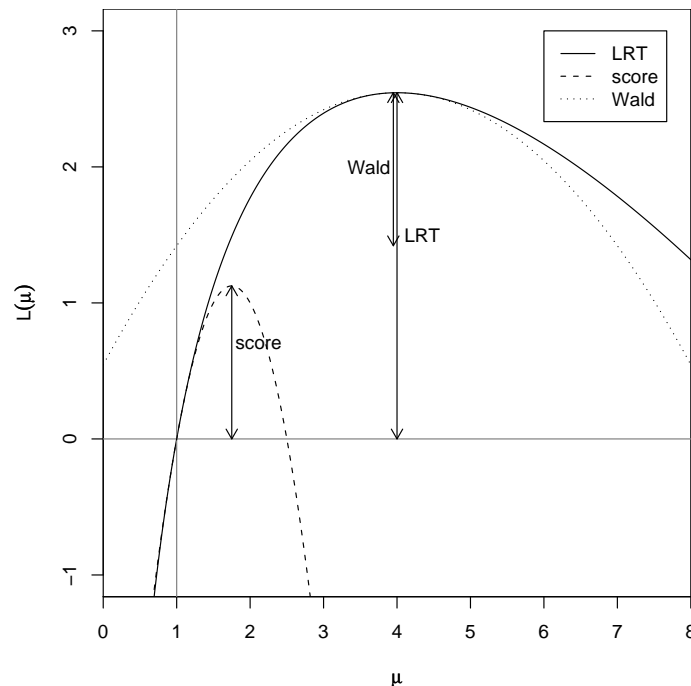
and expanding  $L'(\theta) \sim L'(\hat{\theta}) + H(\theta)(\theta - \hat{\theta})$  if we replace  $H(\theta_0)$  by  $-I(\theta_0)$ .

For regression with known  $\sigma^2$  the log-likelihood is quadratic and has known curvature (the Hessian is constant in both  $\theta$  and the data), so all these methods agree exactly. Under standard large-sample theory all these tests have an asymptotic  $\chi_q^2$  distribution, but in small-sample situations they can differ considerably (and the  $\chi_q^2$  distribution can be a poor approximation).

---

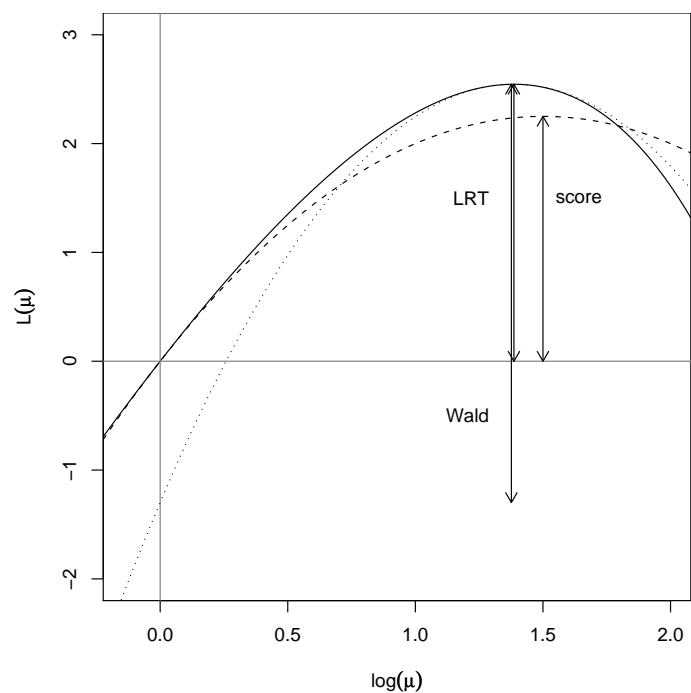
<sup>1</sup>the matrix of second derivatives.

As a simple example, consider the  $\text{Poisson}(\mu)$  family, and a test of  $\mu = 1$  with a single observation  $Y = 4$ : see figure 1. (This is an unusually large observation, but is not extreme having probability 1.5%.)



**Figure 1:** Tests for  $\mu = 1$  in a Poisson family with observation  $Y = 4$ . The exact log-likelihood and the two quadratic approximations are shown.

Note that whereas the LRT is invariant to a change in parametrization, the most of the approximations are not. It is natural to consider log-linear models for a Poisson family, so let us repeat the example for  $\log \mu$ : see figure 2. In this case the approximations are much better.



**Figure 2:** Tests for  $\mu = 1$  in a Poisson family with observation  $Y = 4$ , parameterized by  $\theta = \log \mu$ . Other details as figure 1.

## Deviance and Chi-squared tests

Now suppose we have  $k$  Poisson-distributed counts  $O_i$  with means  $\mu_i$  which are functions of the parameter(s)  $\theta$ . The log-likelihood is

$$L(\theta) = \sum_i O_i \log \mu_i(\theta) - \mu_i(\theta) - \log O_i!$$

and since the saturated model has  $\hat{\mu}_i = O_i$ , the deviance is

$$G^2 = 2 \sum_i \left[ O_i \log \{O_i / \mu_i(\hat{\theta})\} + (\hat{\mu}_i(\theta) - O_i) \right]$$

and for most models (including all log-linear model with an intercept)

$$\sum_i (\hat{\mu}_i(\theta) - O_i) = 0$$

so the  $G^2$  test as normally considered ignores this term. It is conventional to call  $\mu(\theta_i)$  the expected value  $E_i$ .

Now consider the test of the parametrized means within the general model of freely-specified means. Then

$$L(\mu) = \sum_i O_i \log \mu_i - \mu_i - \log O_i!$$

and so

$$L(\mu) - L(E) = \sum_i O_i \log(\mu_i / E_i)$$

provided  $\sum \mu_i = \sum E_i$ . The MLE is  $\hat{\mu}_i = O_i$  and the LRT is

$$G^2 = 2L(\hat{\mu}) - 2L(E) = 2 \sum_i O_i \log(O_i / E_i)$$

If we expand this about  $O_i$  we get

$$G^2 = -2 \sum_i O_i \log(E_i / O_i) \approx 2 \sum_i O_i \frac{(O_i - E_i)^2}{2O_i^2} = \sum_i \frac{(O_i - E_i)^2}{O_i}$$

Here the Hessian of  $L$  is  $-\text{diag}(O_i / \mu_i^2)$  so the Fisher information matrix is  $\text{diag}(1 / \mu_i)$  which when evaluated at the MLE is the same as the observed information matrix. So the two forms of the Wald approximation agree and give a variation on the Chi-squared test.

For the score test, we first find the first and second derivatives:

$$\frac{\partial L}{\partial \mu_i} = O_i / \mu_i - 1, \quad \frac{\partial^2 L}{\partial \mu_i \partial \mu_j} = -\delta_{ij} O_i / \mu_i^2, \quad I(\mu) = \text{diag}(1 / \mu_i)$$

Thus the score statistic is

$$U^T I(E)^{-1} U = [O_i / E_i - 1]^T I(E)^{-1} [O_i / E_i - 1] = \sum_i \left( \frac{O_i}{E_i} - 1 \right) \frac{1}{E_i} \left( \frac{O_i}{E_i} - 1 \right) = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

the Chi-squared statistic.

## The Mantel–Haenszel test

The Mantel–Haenszel test of conditional independence in multiple  $2 \times 2$  tables is a score test approximation to the LRT in a logistic regression. Suppose we have a  $2 \times 2 \times K$  table with counts  $n_{ijk}$  and with variables  $X$ ,  $Y$  and  $Z$ , and we are interested in whether

$$\pi_{ik} = P(Y = 2 | X = i, Z = k)$$

depends on  $i$ . The Mantel–Haenszel test statistic is

$$MH = \sum_k \frac{[n_{22k} - \mu_{22k}]^2}{\text{var}(n_{22k})}$$

The full logit model is

$$\text{logit } \pi_{ik} = \alpha + \beta x_i + \gamma_k, \text{ for } j = 1, 2 \text{ and } k = 1, \dots, K$$

where  $x_1 = 0$  and  $x_2 = 1$ , say. The null hypothesis of conditional independence of  $X$  and  $Y$  given  $Z$  corresponds to  $\beta = 0$ . Writing out the details of the score test (Day and Byar, *Biometrics*, 1979) of  $\beta = 0$  leads to  $MH$ .

## Code

For interest, figure 1 was computed in R (to allow mathematics in the labels) by

```
xx <- seq(0, 8, len = 500)
f <- function(mu) log(dpois(4, mu)/dpois(4,1))
plot(xx, f(xx), type = "l", ylim = c(-1, 3), xaxs = "i",
      xlab = expression(mu), ylab = expression(L(mu)) )
abline(h = 0, col = "gray50")
abline(v = 1, col = "gray50")
arrows(4, 0, 4, f(4), length = 0.1, code = 3)
text(4.1, 1.5, "LRT", adj = 0)

fW <- function(mu) f(4) - (mu-4)^2/8
lines(xx, fW(xx), lty = 3)
arrows(3.95, fW(1), 3.95, f(4), length = 0.1, code = 3)
text(3.9, 2, "Wald", adj = 1)

fS <- function(mu) 3*(mu-1) - 2*(mu-1)^2
lines(xx, fS(xx), lty = 2)
arrows(1.75, 0, 1.75, fS(1.75), length = 0.1, code = 3)
text(1.8, 0.7, "score", adj = 0)

legend(6, 3, c("LRT", "score", "Wald"), lty = 1:3)
```