

CHAPTER 7. HYPOTHESIS TESTING

7.1. THE GENERAL PROBLEM

It is often necessary to make a decision, on the basis of available data from an experiment (carried out by yourself or by Nature), on whether a particular proposition H_0 (theory, model, hypothesis) is true, or the converse H_1 is true. This decision problem is often encountered in scientific investigation. Economic examples of hypotheses are

- (a) The commodities market is efficient (i.e., opportunities for arbitrage are absent).
- (b) There is no discrimination on the basis of gender in the market for academic economists.
- (c) Household energy consumption is a necessity, with an income elasticity not exceeding one.
- (d) The survival curve for Japanese cars is less convex than that for Detroit cars.

Notice that none of these economically interesting hypotheses are framed directly as precise statements about a probability law (e.g., a statement that the parameter in a family of probability densities for the observations from an experiment takes on a specific value). A challenging part of statistical analysis is to set out maintained hypotheses that will be accepted by the scientific community as true, and which in combination with the proposition under test give a probability law. Deciding the truth or falsity of a null hypothesis H_0 presents several general issues: the cost of mistakes, the selection and/or design of the experiment, and the choice of the test.

7.2. THE COST OF MISTAKES

Consider a two-by-two table that compares the truth with the result of the statistical decision. For now, think of each of the alternatives H_0 and H_1 as determining a unique probability law for the observations; these are called simple alternatives. Later, we consider compound hypotheses or alternatives that are consistent with families of probability laws.

		Truth	
		H_0	H_1
Decision	H_0 Accepted	Cost = 0 Probability = $1 - \alpha$	Cost = C_{II} Probability = β
	H_0 Rejected	Cost = C_I Probability = α	Cost = 0 Probability = $\pi = 1 - \beta$

There are two possible mistakes, *Type I* in which a true hypothesis is rejected, and *Type II* in which a false hypothesis is accepted. There are costs associated with these mistakes -- let C_I denote the cost

associated with a Type I mistake, and C_{II} denote the cost associated with a Type II mistake. If the hypothesis is true, then there is a probability α that a particular decision procedure will result in rejection; this is also called the *Type I error probability* or the *significance level*. If the hypothesis is false, there is a probability β that it will be accepted; this is called the *Type II error probability*. The probability $\pi \equiv 1-\beta$ is the probability that the hypothesis will be rejected when it is false, and is called the *power* of the decision procedure.

This table is in principle completely symmetric between the states H_0 and H_1 : You can call your favorite theory H_0 and hope the evidence leads to it being accepted, or call it H_1 and hope the evidence leads to H_0 being rejected. However, classical statistical analysis is oriented so that α is chosen by design, and β requires a sometimes complex calculation. Then, the Type I error is easier to control. Thus, in classical statistics, it is usually better to assign your theory between H_0 and H_1 so that the more critical mistake becomes the Type I mistake. For example, suppose you set out to test your favorite theory. Your study will be convincing only if your theory passes a test which it would have a high (and known) probability of failing if it is in fact false. You can get such a stringent test by making your theory H_1 and selecting a null and a decision procedure for which α is known and small; then your theory will be rejected in favor of H_0 with large known probability $1-\alpha$ if in fact H_0 rather than H_1 is true. (This will not work if you pick a "straw horse" for the null that no one thinks is plausible.) Conversely, if you set out to do a convincing demolition of a theory that you think is false, then make it the null, so that there is a small known probability α of rejecting the hypothesis if it is in fact true.

A common case for hypothesis testing is that the null hypothesis H_0 is simple, but the alternative hypothesis H_1 is compound, containing a family of possible probability laws. Then, the probability of a Type II error depends on which member of this family is true. Thus, the power of a test is a function of the specific probability law in a compound alternative. When both the null hypothesis and alternative are compound, the probability of a Type I error is a function of which member of the family of probability laws consistent with H_0 is true. In classical statistics, the significance level is always defined to be the "worst case": the largest α for any probability law consistent with the null.

Given the experimental data available and the statistical procedure adopted, there will be a trade off between the probabilities of Type I and Type II errors. When the cost C_I is much larger than the cost C_{II} , a good decision procedure will make α small relative to β . Conversely, when C_I is much smaller than C_{II} , the procedure will make α large relative to β . For example, suppose the null hypothesis is that a drug is sufficiently safe and effective to be released to the market. If the drug is critical for treatment of an otherwise fatal disease, then C_I is much larger than C_{II} , and the decision procedure should make α small. Conversely, a drug to reduce non-life-threatening wrinkles should be tested by a procedure that makes β small.

7.3. DESIGN OF THE EXPERIMENT

One way to reduce the probability of Type I and Type II errors is to collect more observations by increasing sample size. One may also by clever design be able to get more information from a given sample size, or more relevant information from a given data collection budget. One has the

widest scope for action when the data is being collected in a laboratory experiment that you can specify. For example, the Negative Income Experiments in the 1960's and 1970's were able to specify experimental treatments that presented subjects with different trade offs between wage and transfer income, so that labor supply responses could be observed. However, even in investigations where only natural experiments are available, important choices must be made on what events to study and what data to collect. For example, if a survey of 1000 households is to be made to determine the income elasticity of the demand for energy, one can get more precision by oversampling high income and low income households to get a greater spread of incomes.

There is an art to designing experiments or identifying natural experiments that allow tests of a null hypothesis without confounding by extraneous factors. For example, suppose one wishes to test the null hypothesis that Japanese cars have the same durability as Detroit cars. One might consider the following possible experiments:

- (a) Determine the average age, by origin, of registered vehicles.
- (b) Sample the age/make of scrapped cars as they arrive at junk yards.
- (c) Draw a sample of individual new cars, and follow them longitudinally until they are scrapped.
- (d) Draw a sample of individual new cars, and operate them on a test track under controlled conditions until they fail.

Experiment (a) is confounded by potential differences in historical purchase patterns; some of this could be removed by econometric methods that condition on the number of original purchases in earlier years. Experiments (a)-(c) are confounded by possible variations in usage patterns (urban/rural, young/old, winter roads/not). For example, if rural drivers who stress their cars less tend to buy Detroit cars, this factor rather than the intrinsic durability of the cars might make Detroit cars appear to last longer. One way to reduce this factor would be to assign drivers to car models randomly, as might be done for example for cars rented by Avis in the "compact" category. The ideal way to do this is a "double blind" experiment in which neither the subject nor the data recorder knows which "treatment" is being received, so there is no possibility that bias in selection or response could creep in. Most economic experimental treatments are obvious to aware subjects, so that "double blind" designs are impossible. This puts an additional burden on the researcher to carefully randomize assignment of treatments and to structure the treatments so that their form does not introduce factors that confound the experimental results.

Economists are often confronted with problems and data where a designed experiment is infeasible and Nature has not provided a clean "natural experiment", and in addition sample frames and protocols are not ideal. It may nevertheless be possible to model the data generation process to take account of sampling problems, and to use multivariate statistical methods to estimate and test hypotheses about the separate effects of different factors. This exercise can provide useful insights, but must be used cautiously and carefully to avoid misattribution and misinterpretation. Econometricians should follow the rule "Do No Harm". When a natural experiment or data are not adequate to resolve an economic hypothesis, econometric analysis should stop, and not be used to dress up propositions that a righteous analysis cannot support. Every econometric study should

consider very carefully all the possible processes that could generate the observed data, candidly discuss alternative explanations of observations, and avoid unsupportable claims..

7.4. CHOICE OF THE DECISION PROCEDURE

Suppose one thinks of hypothesis testing as a statistical decision problem, like the problem faced by Cab Franc in Chapter 1, with a prior p_0 that H_0 is true and $p_1 = 1 - p_0$ that H_1 is true. Let $f(\mathbf{x}|H_0)$ denote the likelihood of \mathbf{x} if H_0 is true, and $f(\mathbf{x}|H_1)$ denote the likelihood if H_1 is true. Then, the posterior likelihood of H_0 given \mathbf{x} is, by application of Bayes Law, $q(H_0|\mathbf{x}) = f(\mathbf{x}|H_0)p_0/[f(\mathbf{x}|H_0)p_0 + f(\mathbf{x}|H_1)p_1]$. The expected cost of rejecting H_0 given \mathbf{x} is then $C_1q(H_0|\mathbf{x})$, and the expected cost of accepting H_0 given \mathbf{x} is $C_0q(H_1|\mathbf{x})$. The optimal decision rule is then to *reject* H_0 for \mathbf{x} in the *critical region* C where $C_1q(H_0|\mathbf{x}) < C_0q(H_1|\mathbf{x})$. This inequality simplifies to $C_1f(\mathbf{x}|H_0)p_0 < C_0f(\mathbf{x}|H_1)p_1$, implying

$$\mathbf{x} \in C \text{ (i.e., reject } H_0) \text{ if and only if } f(\mathbf{x}|H_1)/f(\mathbf{x}|H_0) > k \equiv C_0p_0/C_1p_1.$$

The expression $f(\mathbf{x}|H_1)/f(\mathbf{x}|H_0)$ is termed the *likelihood ratio*. The optimal criterion is then to reject H_0 if and only if the likelihood ratio exceeds a threshold k . The larger C_1 or p_0 , the larger this threshold.

A classical statistical treatment of this problem will also pick a *critical region* C of \mathbf{x} for which H_0 will be rejected, and will do so by maximizing power $\pi = \int_C f(\mathbf{x}|H_1)d\mathbf{x}$ subject to the constraint

$$\alpha = \int_C f(\mathbf{x}|H_0)d\mathbf{x}. \text{ But this is accomplished by picking } C = \{\mathbf{x} | f(\mathbf{x}|H_1)/f(\mathbf{x}|H_0) > k\}, \text{ where } k \text{ is}$$

a constant chosen so the constraint is satisfied. To see why observe that if C contains a little rectangle $[\mathbf{x}, \mathbf{x} + \delta \mathbf{1}]$, where δ is a tiny positive constant, then this rectangle contributes $f(\mathbf{x}|H_0)\delta^n$ to meeting the constraint and $f(\mathbf{x}|H_1)\delta^n$ to power. The ratio $f(\mathbf{x}|H_1)/f(\mathbf{x}|H_0)$ then gives the rate at which power is produced per unit of type I error probability used up. The optimal critical region will start where this rate is the highest, and keep adding to C until by decreasing the rate threshold k until the type I error constraint is met.

The optimal decision rule for various prior probabilities and costs and the classical statistical test procedure trace out the same families of procedures, and will coincide when the critical likelihood ratio k in the two approaches is the same. In more general classical hypothesis testing situations where the alternative is compound, there is no longer an exact coincidence of the classical and statistical decision theory approaches to decisions. However, the likelihood ratio often remains a useful basis for constructing good test procedures. In many cases, a "best" test by some classical statistical criterion and a test utilizing the likelihood ratio criterion will be the same or nearly the same.

In general, we will consider DGP which we maintain are members of a family $f(\mathbf{x}, \theta)$ indexed by a parameter θ . The null hypothesis is that the true value θ_0 of θ is contained in a set N , and the

alternative is that it is contained in a set \mathbf{A} , with \mathbf{A} and \mathbf{N} partitioning the universe Θ of possible values of θ_0 . The value θ_e of θ that maximizes $f(\mathbf{x}, \theta)$ over $\theta \in \Theta$ is the *maximum likelihood estimator*. The theory of the maximum likelihood estimator given in Chapter 6 shows that it will have good statistical properties in large samples under mild regularity conditions. The value θ_{oe} that maximizes $f(\mathbf{x}, \theta)$ over $\theta \in \mathbf{N}$ is called the *constrained maximum likelihood estimator* subject to the null hypothesis. When the null hypothesis is true, the constrained maximum likelihood estimator will also have good statistical properties. Intuitively, the reason is that when the null hypothesis is true, the true parameter satisfies the hypothesis, and hence the maximum value of the constrained likelihood will be at least as high as the value of the likelihood at the true parameter. If an identification condition is met, the likelihood at the true parameter converges in probability to a larger number than the likelihood at any other parameter value. Then, the constrained maximum likelihood must converge in probability to the true parameter. A rigorous proof of the properties of constrained estimators is given in Chapter 22.

A likelihood ratio critical region for the general testing problem is usually defined as a set of the form

$$\mathbf{C} = \{\mathbf{x} \mid \sup_{\theta \in \mathbf{A}} f(\mathbf{x}, \theta) / \sup_{\theta \in \mathbf{N}} f(\mathbf{x}, \theta) > k\}.$$

The likelihood ratio in this criterion is less than or equal to one when the maximum likelihood estimator of θ_0 falls in \mathbf{N} , and otherwise is greater than one. Then a critical region defined for some $k > 1$ will include the observed vectors \mathbf{x} that are the least likely to have been generated by a DGP with a parameter in \mathbf{N} . The significance level of the test is set by adjusting k . An equivalent way to characterize the likelihood ratio critical region \mathbf{C} is in terms of the log likelihood function,

$$\mathbf{C} = \{\mathbf{x} \mid \sup_{\theta \in \Theta} \log f(\mathbf{x}, \theta) - \sup_{\theta \in \mathbf{N}} \log f(\mathbf{x}, \theta) > \kappa = \log k\}.$$

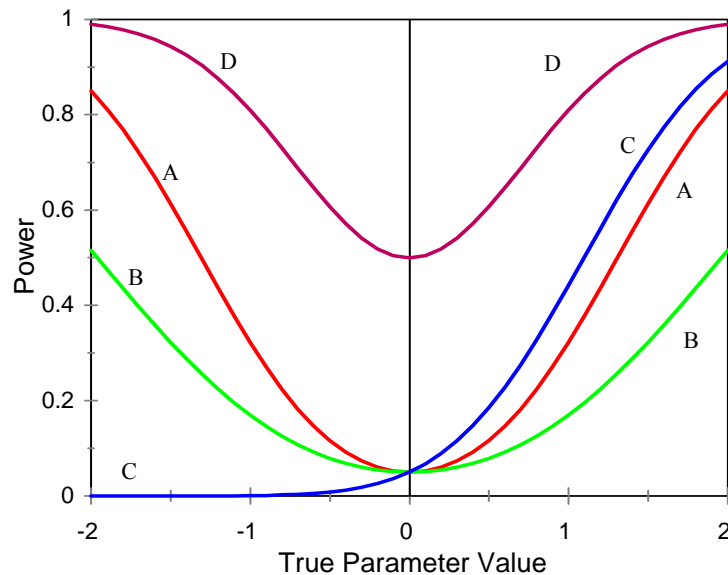
Clearly, the log ratio in this expression equals the difference in the log likelihood evaluated at the maximum likelihood estimator and the log likelihood evaluated at the constrained maximum likelihood estimator. This difference is zero if the maximum likelihood estimator is in \mathbf{N} , and is otherwise positive.

The analyst will often have available alternative testing procedures in a classical testing situation. For example, one procedure to test a hypothesis about a location parameter might be based on the sample mean, a second might be based on the sample median, and a third might be based on the likelihood ratio. Some of these procedures may be better than others in the sense of giving higher power for the same significance level. The ideal, as in the simple case, is to maximize the power given the significance level. When there is a compound alternative, so that power is a function of the alternative, one may be able to tailor the test to have high power against alternatives of particular importance. In a few cases, there will be a single procedure that will have uniformly best power against a whole range of alternatives. If so, this will be called a *uniformly most powerful test*.

The figure below shows power functions for some alternative test procedures. The null hypothesis is that a parameter θ is zero. Power curves A and B equal 0.05 when $H_0: \theta = 0$ is true. Then, the significance level of these three procedures is $\alpha = 0.05$. The significance level of D is

much higher, 0.5. Compare the curves A and B. Since A lies everywhere above B and has the same significance level, A is clearly the superior procedure. A comparison like A and B most commonly arises when A uses more data than B; that is, A corresponds to a larger sample. However, it is also possible to get a picture like this when A and B are using the same sample, but B makes poor use of the information in the sample.

Compare curves A and C. Curve C has significance level $\alpha = 0.05$, and has lower power than A against alternatives less than $\theta = 0$, but better power against alternatives greater than $\theta = 0$. Thus, A is a better test if we want to test against all alternatives, while C is a better test if we are mainly interested in alternatives to the right of $\theta = 0$ (i.e., we want to test $H_0: \theta \leq 0$). Compare curves A and D. Curve D has high power, but at the cost of a high probability of a Type I error. Thus, A and D represent a trade off between Type I and Type II errors.



Finally, suppose we are most interested in the alternative $H_1: \theta = 1.5$. The procedure giving curve A has power 0.61 against this alternative, and hence has a reasonable chance of discriminating between H_0 and H_1 . On the other hand, the procedure B has power 0.32, and much less chance of discriminating. We would conclude that the procedure A is a moderately satisfactory statistical test procedure, while B is of limited use.

7.5. HYPOTHESIS TESTING IN NORMAL POPULATIONS

This section provides a summary of hypothesis test calculations for standard setups involving data drawn from a normal population, including power calculations. Assume that we start from a simple random sample of size n , giving i.i.d. observations x_1, \dots, x_n . Recall from Chapter 6.3 that the log likelihood of a normal random sample is

$$L(\mathbf{x}, \mu, \sigma^2) = - \frac{n}{2} \cdot \text{Log}(2\pi) - \frac{n}{2} \cdot \text{Log} \sigma^2 - \frac{1}{2} \cdot \sum_{i=1}^n (x_i - \mu)^2 / \sigma^2$$

$$= - \frac{n}{2} \cdot \text{Log}(2\pi) - \frac{n}{2} \cdot \text{Log} \sigma^2 - \frac{1}{2} \cdot \frac{(n-1)s^2}{\sigma^2} - \frac{n}{2\sigma^2} (\bar{x} - \mu)^2 \equiv L(\bar{x}, s^2, \mu, \sigma^2).$$

where the sample mean $\bar{x} = \sum_{i=1}^n x_i$ and the sample variance $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ are unbiased estimators of μ and σ^2 , respectively. If \mathbf{N} denotes the set of parameter values (μ, σ^2) consistent with a null hypothesis, then a likelihood ratio critical set for this hypothesis will take the form

$$C = \{(\bar{x}, s^2) \mid \sup_{\theta \in \Theta} L(\bar{x}, s^2, \mu, \sigma^2) - \sup_{\theta \in \mathbf{N}} L(\bar{x}, s^2, \mu, \sigma^2) > \kappa\}.$$

We consider a sequence of hypotheses and conditions. See Chapter 3.7 for the densities and other properties of the distributions used in this section, and Chapter 6.3 for the relation of these distributions to data from a normal population. The following table gives the statistical functions that are available in many econometrics software packages; the specific notation is that used in the Statistical Software Tools (SST) package. Tables at the end of most statistics texts can also be used to obtain values of the central versions of these distributions. The direct functions give the CDF probability for a specified argument, while the inverse functions give the argument that yields a specified probability.

Distribution	CDF	Inverse CDF
Normal	cumnorm(x)	invnorm(p)
Chi-Square	cumchi(x,df)	invchi(p,df)
F-Distribution	cumf(x,df1,df2)	invf(p,df1,df2)
T-Distribution	cumt(x,df)	invt(p,df)
Non-central Chi-Square	cumchi(x,df, δ)	NA
Non-central F-Distribution	cumf(x,df1,df2, δ)	NA
Non-Central T-Distribution	cumt(x,df, λ)	NA

In this table, df denotes degrees of freedom, and λ and δ are non-centrality parameters. Inverse CDF's are not available for non-central distributions in most packages, and are not needed. In most statistical packages, values of these functions can either be printed out or saved for further calculations. For example, in SST, the command “calc cumnorm(1.7)” will print out the probability

that a standard normal random variable is less than 1.7, the command “calc p = cumnorm(1.7)” will store the result of this calculation in the variable p for further use, and a subsequent command “calc p” will also print out its value.

Problem 1: Testing the mean of a normal population that has known variance

Suppose a random sample of size n from a normal population with an unknown mean μ and a known variance σ^2 . The null hypothesis is $H_0: \mu = \mu_0$, and the alternative is $H_1: \mu \neq \mu_0$. Verify that the likelihood ratio, $\text{Max}_{\mu} n(\bar{x}, \mu, 1) / n(\bar{x}, \mu_0, 1)$, is an increasing function of $(\bar{x} - \mu_0)^2$. Hence, a test equivalent to a likelihood ratio test can be based on $(\bar{x} - \mu_0)^2$. From Chapter 6.3(8), one has the result that under the null hypothesis, the statistic $n(\bar{x} - \mu_0)^2 / \sigma^2$ is distributed χ_1^2 . Alternately, from Chapter 6.3(5), the square root of this expression, $n^{1/2}(\bar{x} - \mu_0) / \sigma$, has a standard normal distribution.

Using the Chi-Square form of the statistic, the critical region will be values exceeding a critical level z_c , where z_c is chosen so that the selected significance level α satisfies $\chi_1^2(z_c) = 1 - \alpha$. For example, taking $\alpha = 0.05$ yields $z_c = 3.84146$. This comes from a statistical table, or from the SST command “calc invchi(1 - α , k)”, where α is the significance level and k is degrees of freedom. The test procedure rejects H_0 whenever

$$(1) \quad n(\bar{x} - \mu_0)^2 / \sigma^2 > z_c = 3.84146.$$

Consider the power of the Chi-square test against an alternative such as $\mu = \mu_1 \neq \mu_0$. The non-centrality parameter is

$$(2) \quad \delta = n(\mu_1 - \mu_0)^2 / \sigma^2.$$

For example, if $\mu_1 - \mu_0 = 1.2$, $\sigma^2 = 25$, and $n = 100$, then $\delta = 1.44 \cdot 100 / 25 = 5.76$. The power is calculated from the non-central Chi-square distribution (with 1 degree of freedom), and equals the probability that a random draw from this distribution exceeds z_c . This probability π is readily calculated using the SST command “calc 1 - cumchi(z_c , k, δ)”. In the example, $\pi = \text{calc } 1 - \text{cumchi}(3.84146, 1, 5.76) = 0.67006$. Then, a test with a five percent significance level has power of 67 percent against the alternative that the true mean is 1.2 units larger than hypothesized.

An equivalent test can be carried out using the standard normal distributed form $n^{1/2}(\bar{x} - \mu_0) / \sigma$. The critical region will be values of this expression that in magnitude exceed a critical level w_c , where w_c is chosen for a specified significance level α so that a draw from a standard normal density has probability $\alpha/2$ of being below $-w_c$, and symmetrically a probability $\alpha/2$ of being above $+w_c$. One can find w_c from statistical tables, or by using a SST command “calc invnorm(1 - $\alpha/2$)”. For example, if $\alpha = 0.05$, then $w_c = \text{calc invnorm}(0.975) = 1.95996$. The test rejects H_0 whenever

$$(3) \quad n^{1/2}(\bar{x} - \mu_0) / \sigma < -w_c \text{ or } n^{1/2}(\bar{x} - \mu_0) / \sigma > w_c.$$

For example, if $n = 100$, $\sigma = 5$, and $\mu_0 = 0$, the critical region for a test with significance level $\alpha = 0.05$ is $10 \bar{x} / 5 < -1.95996$ or $10 \bar{x} / 5 > +1.95996$. Note that $w_c^2 = z_c$, so this test rejects exactly when the Chi-square test rejects. The power of the test above against the alternative $\mu = \mu_1 \neq \mu_0$ is the

probability that the random variable $n^{1/2}(\bar{x} - \mu_0)/\sigma$ lies in the critical region when $\bar{x} \sim N(\mu_1, \sigma^2)$. In this case, $n^{1/2}(\bar{x} - \mu_1)/\sigma \equiv Y$ is standard normal, and therefore $n^{1/2}(\bar{x} - \mu_0)/\sigma \equiv Y + \lambda$, where

$$(4) \quad \lambda = n^{1/2}(\mu_1 - \mu_0)/\sigma.$$

Note that $\lambda^2 = \delta$, where δ is given in (2). The probability of rejection in the left tail is $\Pr(n^{1/2}(\bar{x} - \mu_0)/\sigma < -w_c | \mu = \mu_1) = \Pr(Y < -w_c - \lambda)$. For the right tail, $\Pr(n^{1/2}(\bar{x} - \mu_0)/\sigma > w_c | \mu = \mu_1) = \Pr(Y > w_c - \lambda)$. Using the fact that the standard normal is symmetric, we then have

$$(5) \quad \pi = \Phi(-w_c - \lambda) + 1 - \Phi(w_c - \lambda) \equiv \Phi(-w_c - \lambda) + \Phi(-w_c + \lambda).$$

This can be calculated using the SST command

$$\pi = \text{calc cumnorm}(-w_c - \lambda) + \text{cumnorm}(-w_c + \lambda).$$

For example, $\sigma = 5$, $N = 100$, $\mu_1 - \mu_0 = 1.2$, $w_c = 1.95996$ give $\delta = 2.4$ and power $\pi = \text{calc cumnorm}(-w_c - 2.4) + \text{cumnorm}(-w_c + 2.4) = 0.670$. Note this is the same as the power of the Chi-square version of the test.

Suppose that instead of testing the null hypothesis $H_0: \mu = \mu_0$ against the alternative $H_1: \mu \neq \mu_0$, you want to test the one-sided hypothesis $H_0: \mu \leq \mu_0$ against the alternative $H_1: \mu > \mu_0$. The likelihood ratio in this case is $\text{Sup}_{\mu > \mu_0} n(\bar{x}, \mu, \sigma^2) / \text{Sup}_{\mu \leq \mu_0} n(\bar{x}, \mu, \sigma^2)$, which is constant for $\bar{x} \leq \mu_0$

and is monotone increasing in $(\bar{x} - \mu_0)$ for $\bar{x} > \mu_0$. Hence, a test that rejects H_0 for $\bar{x} - \mu_0$ large appears desirable. This suggests using a test based on the statistic $n^{1/2}(\bar{x} - \mu_0)/\sigma$, which is normal with variance one, and has a non-positive mean under the null. Pick a critical level $w_c > 0$ such that

$$\text{Sup}_{\mu \leq \mu_0} \text{Prob}(n^{1/2}(\bar{x} - \mu)/\sigma > w_c) = \alpha.$$

Note that the sup is taken over all the possible true μ consistent with H_0 , and that α is the selected significance level. The maximum probability of Type I error is achieved when $\mu = \mu_0$. (To see this, note that $\text{Prob}(n^{1/2}(\bar{x} - \mu_0)/\sigma > w_c) \equiv \Pr(Y \equiv n^{1/2}(\bar{x} - \mu)/\sigma > w_c + n^{1/2}(\mu_0 - \mu)/\sigma)$, where μ is the true value. Since Y is standard normal, this probability is largest over $\mu \leq \mu_0$ at $\mu = \mu_0$.) Then, w_c is determined to give probability α that a draw from a standard normal exceeds w_c . For example, if $n = 100$, $\alpha = 0.05$, $\sigma = 5$, and H_0 is that $\mu \leq 0$, then $w_c = \text{calc invnorm}(0.95) = 1.64485$. The power of the test of $\mu \leq \mu_0 = 0$ against the alternative $\mu = \mu_1 = 1.2$ is given by

$$(6) \quad \begin{aligned} \pi &= \Pr(n^{1/2}(\bar{x} - \mu_0)/\sigma > w_c | \mu = \mu_1) \equiv \Pr(Y \equiv n^{1/2}(\bar{x} - \mu_1)/\sigma > w_c - \lambda) \\ &\equiv 1 - \Phi(w_c - \lambda) \equiv \Phi(-w_c + \lambda) \equiv \text{calc cumnorm}(-w_c + \lambda), \end{aligned}$$

where λ is given in (4). In the example, $\pi = \text{calc cumnorm}(-1.64485 + 2.4) = 0.775$. Hence, a test which has a probability of at most $\alpha = 0.05$ of rejecting the null hypothesis when it is true has power 0.775 against the specific alternative $\mu_1 = 1.2$.

Problem 2. Testing the Mean of a Normal Population with Unknown Variance

This problem is identical to Problem 1, except that σ^2 must now be estimated. Use the estimator s^2 for σ^2 in the Problem 1 test statistics. From Chapter 6.3(8), the Chi-square test statistic with σ replaced by s , $F = n(\bar{x} - \mu_0)^2/s^2$, has an F-distribution with degrees of freedom 1 and $N-1$. Hence, to test $H_0: \mu = \mu_0$ against the alternative $H_1: \mu \neq \mu_0$, find a critical level z_c such that a specified significance level α equals the probability that a draw from $F_{1,n-1}$ exceeds z_c . The SST function $\text{calc } z_c = \text{invf}(1-\alpha, 1, n-1)$ gives this critical level; it can also be found in standard tables. For $n = 100$ and $\alpha = 0.05$, the critical level is $z_c = 3.93694$.

The power of the test against an alternative μ_1 is the probability that the statistic F exceeds z_c . Under this alternative, F has a non-central F-distribution (from Chapter 3.9) with the non-centrality parameter $\delta = n(\mu_1 - \mu_0)^2/\sigma^2$ given in (2). Then, the power is given by

$$(7) \quad \pi = \text{calc } 1 - \text{cumf}(z_c, 1, n-1, \delta).$$

In the example with $\mu_1 - \mu_0 = 1.2$ and $\sigma^2 = 25$, one has $\delta = 144/25$, and the power is

$$(8) \quad \pi = \text{calc } 1 - \text{cumf}(3.93694, 1, 99, 144/25) = 0.662.$$

The non-centrality parameter is defined using the true σ^2 rather than the estimate s^2 . Calculating power at an estimated non-centrality parameter $\delta_e = n(\mu_1 - \mu_0)^2/s^2$ introduces some error -- you will evaluate the power curve at a point somewhat different than you would like. For most practical purposes, you do not need an exact calculation of power; you are more interested in whether it is 0.1 or 0.9. Then, the error introduced by this approximation can be ignored. In particular, for large sample sizes where the power against economically interesting alternatives is near one, this error is usually negligible. Note that $\delta/\delta_e = s^2/\sigma^2$, so $(n-1)\delta/\delta_e$ is distributed $\chi^2(n-1)$. For the rare application where you really need to know how precise your power calculation is, you can form a confidence interval as follows: Given a "significance level" α , compute $z_1 = \text{calc } \text{invchi}(\alpha/2, n-1)$ and $z_2 = \text{calc } \text{invchi}(1-\alpha/2, n-1)$. Then, with probability α , $\delta_1 \equiv z_1\delta_e/(n-1) < \delta < z_2\delta_e/(n-1) \equiv \delta_2$. The power π_1 calculated at δ_1 and the power π_2 calculated at δ_2 give a α -level confidence bound on the exact power. For example, $\alpha = 0.5$, $n = 100$, $\mu_1 - \mu_0 = 1.2$, and $s^2 = 25$ imply $\delta_e = 144/25$, $z_1 = \text{calc } \text{invchi}(.25, 99) = 89.18$, $\delta_1 = 5.189$, and $\pi_1 = \text{calc } 1 - \text{cumf}(3.93694, 1, 99, 5.189) = 0.616$. Also, $z_2 = \text{calc } \text{invchi}(.75, 99) = 108.093$, $\delta_2 = 6.289$, and $\pi_2 = \text{calc } 1 - \text{cumf}(3.93694, 1, 99, 6.289) = 0.700$. Then, with probability 0.5, the exact power for the alternative $\mu_1 - \mu_0 = 2$ is in the interval $[0.616, 0.700]$.

The test of $H_0: \mu = \mu_0$ can be carried out equivalently using

$$(9) \quad T = n^{1/2}(\bar{x} - \mu_0)/s,$$

which by Chapter 6.3(7) has a t-distribution with $n-1$ degrees of freedom under $H_0: \mu = \mu_0$. For a significance level α , choose a critical level w_c , and reject the null hypothesis when $|T| > w_c$. The value of w_c satisfies $\alpha/2 = t_{n-1}(-w_c)$, and is given in standard tables, or in SST by $w_c = \text{invt}(1-\alpha/2, n-1)$. For the example $\alpha = 0.05$ and $n = 100$, this value is $w_c = \text{calc } \text{invt}(.975, 99) = 1.9842$.

The power of the test is calculated as in Problem 1, replacing the normal distribution by the non-central t-distribution: $\pi = t_{n-1,\lambda}(-w_c) + 1 - t_{n-1,\lambda}(w_c)$, where $\lambda = n^{1/2}(\mu_1 - \mu_0)/\sigma$ as in equation (4). Points of the non-central t are not in standard tables, but are provided by a SST function, $\pi = \text{cumt}(-w_c, n-1, \lambda) + 1 - \text{cumt}(w_c, n-1, \lambda)$. For the example $\alpha = 0.05$, $N = 100$, $\sigma = 5$, and $\mu_1 - \mu_0 = 1.2$ imply $\lambda = 2.4$, and this formula gives $\pi =$

The T-statistic (9) can be used to test the one-sided hypothesis $H_0: \mu \leq \mu_0$. The hypothesis will be rejected if $T > w_c$, where w_c satisfies $\alpha = t_{n-1}(-w_c)$, and is given in standard tables, or in SST by $w_c = \text{invt}(1-\alpha, n-1)$. The power of the test is calculated in the same way as the one-sided test in Problem 1, with the non-central t-distribution replacing the normal: $\pi = 1 - \text{cumt}(w_c, n-1, \lambda)$.

Problem 3. Testing the Variance of a Normal Population with Unknown Mean

Suppose $H_0: \sigma^2 = \sigma_0^2$ versus the alternative H_1 that this equality does not hold.. Under the null, the statistic $X \equiv (n-1)s^2/\sigma_0^2$ is distributed $\chi^2(n-1)$. Then, a test with significance level α can be made by rejecting H_0 if $X < z_{c1}$ or $X > z_{c2}$, where z_{c1} and z_{c2} are chosen so the probability is $\alpha/2$ that a draw from $\chi^2(n-1)$ is less than z_{c1} , and $\alpha/2$ that it is greater than z_{c2} . These can be calculated using $z_{c1} = \text{calc invchi}(\alpha/2, n-1)$ and $z_{c2} = \text{calc invchi}(1-\alpha/2, n-1)$. To calculate the power of the test against the alternative $H_1: \sigma^2 = \sigma_1^2$, note that in this case $(n-1)s^2/\sigma_1^2 = X\sigma_0^2/\sigma_1^2 \equiv Y$ is $\chi^2(n-1)$. Then,

$$\begin{aligned} \pi &= 1 - \Pr(z_{c1} \leq X \leq z_{c2} | \sigma^2 = \sigma_1^2) = 1 - \Pr(z_{c1}\sigma_0^2/\sigma_1^2 \leq Y \leq z_{c2}\sigma_0^2/\sigma_1^2) \\ &= \text{calc cumchi}(z_{c1}\sigma_0^2/\sigma_1^2, n-1) + 1 - \text{cumchi}(z_{c2}\sigma_0^2/\sigma_1^2, n-1). \end{aligned}$$

Problem 4. Testing the Equality of Unknown Variances in Two Populations

Suppose independent random samples of sizes n_i are drawn from normal populations with means μ_i and variances σ_i^2 , respectively, for $i = 1, 2$. The null hypothesis is $H_0: \sigma_1^2 = \sigma_2^2$, and the alternative is $\sigma_1^2 \neq \sigma_2^2$. For each population, we know from 3.6 that $(n_i-1)s_i^2/\sigma_i^2$ has a Chi-square distribution with n_i-1 degrees of freedom. Further, we know that the ratio of two independent Chi-square distributed random variables, each divided by its degrees of freedom, has an F-distribution with these respective degrees of freedom. Then, $R = s_1^2/s_2^2$ is distributed $F(n_1-1, n_2-1)$ under H_0 . One can form a critical region $C = \{R | R < c_L \text{ or } R > c_U\}$ that has significance level α by choosing the lower and upper tails c_L and c_U of the F-distribution so that each has probability $\alpha/2$.

Under alternatives to the null, the ratio s_1^2/s_2^2 , multiplied by the ratio σ_2^2/σ_1^2 , has a central $F(n_1-1, n_2-1)$ -distribution, and the power of the test is

$$\begin{aligned} \pi &= 1 - \text{Prob}(c_L \leq R \leq c_U) = 1 - \text{Prob}(c_L \sigma_2^2/\sigma_1^2 \leq R \sigma_2^2/\sigma_1^2 \leq c_U \sigma_2^2/\sigma_1^2) \\ &= F(c_L \sigma_2^2/\sigma_1^2, n_1-1, n_2-1) + 1 - F(c_U \sigma_2^2/\sigma_1^2, n_1-1, n_2-1). \end{aligned}$$

Problem 5. Testing the Equality of Unknown Means in Two Populations with a Common Unknown Variance

Suppose independent random samples of sizes n_i are drawn from normal populations with means μ_i for $i = 1, 2$ and a common variance σ^2 . The null hypothesis is $H_0: \mu_1 = \mu_2$, and the alternative is $\mu_1 \neq \mu_2$. Then $\bar{x}_1 - \bar{x}_2$ is normally distributed with mean $\mu_1 - \mu_2$ and variance $\sigma^2(n_1^{-1} + n_2^{-1})$. Further,

$(n_1-1)s_1^2/\sigma^2$ is chi-square with n_1-1 degrees of freedom, $(n_2-1)s_2^2/\sigma^2$ is chi-square with n_2-1 degrees of freedom, and all three random variables are independent. Then $((n_1-1)s_1^2 + (n_2-1)s_2^2)/\sigma^2$ is chi-square with $n_1 + n_2 - 2$ degrees of freedom. It follows that

$$s^2 = (n_1^{-1} + n_2^{-1}) \cdot ((n_1-1)s_1^2 + (n_2-1)s_2^2) / (n_1 + n_2 - 2)$$

is an unbiased estimator of $\sigma^2(n_1^{-1} + n_2^{-1})$, with $(n_1 + n_2 - 2)s^2/\sigma^2(n_1^{-1} + n_2^{-1})$ distributed Chi-square with $n_1 + n_2 - 2$ degrees of freedom. Therefore, the statistic

$$(\bar{x}_1 - \bar{x}_2)/s = (\bar{x}_1 - \bar{x}_2) / [(n_1^{-1} + n_2^{-1}) \cdot ((n_1-1)s_1^2 + (n_2-1)s_2^2) / (n_1 + n_2 - 2)]^{1/2}$$

is distributed under the null hypothesis with a T-distribution with $n_1 + n_2 - 2$ degrees of freedom. The power against an alternative $\mu_1 \neq \mu_2$ is calculated exactly as in Problem 2, following (9), except the degrees of freedom is now $n_1 + n_2 - 2$ and the non-centrality parameter is

$$\lambda = (\mu_1 - \mu_2) / \sigma(n_1^{-1} + n_2^{-1})^{1/2}.$$

7.6. HYPOTHESIS TESTING IN LARGE SAMPLES

Consider data $\mathbf{x} = (x_1, \dots, x_n)$ obtained by simple random sampling from a population with density $f(\mathbf{x}, \theta_0)$, where θ_0 is a $k \times 1$ vector of unknown parameters contained in the interior of a set Θ . The sample DGP is $f(\mathbf{x}, \theta_0) = \prod_{i=1}^n f(x_i, \theta_0)$ and log likelihood is $L_n(\mathbf{x}, \theta) = \sum_{i=1}^n l(x_i, \theta)$, where $l(x, \theta) = \log f(x, \theta)$ is the log likelihood of an observation. Consider the maximum likelihood estimator $T_n(\mathbf{x})$, given by the value of θ that maximizes $L_n(\mathbf{x}, \theta)$. Under general regularity conditions like those given in Chapter 6.4, the maximum likelihood estimator is consistent and asymptotically normal. This implies specifically that $n^{1/2}(T_n(\mathbf{x}) - \theta_0) \rightarrow_d Z_0$ with $Z_0 \sim N(0, J^{-1})$ and J the Fisher information in an observation, $J = \mathbf{E} [\nabla_{\theta} l(\mathbf{x}, \theta_0)] [\nabla_{\theta} l(\mathbf{x}, \theta_0)]'$. The Chapter 3.1.18 rule for limits of continuous transformations implies $n^{1/2} \cdot J^{1/2}(T_n(\mathbf{x}) - \theta_0) \rightarrow_d N(0, I)$, and hence that the quadratic form $W(\mathbf{x}, \theta_0) \equiv n \cdot (T_n(\mathbf{x}) - \theta_0)' J (T_n(\mathbf{x}) - \theta_0) \equiv (T_n(\mathbf{x}) - \theta_0)' \mathbf{V}(T_n(\mathbf{x}))^{-1} (T_n(\mathbf{x}) - \theta_0) \rightarrow_d \chi^2(k)$, the Chi-square distribution with k degrees of freedom. When $k = 1$, this quadratic form equals the square of the difference between $T_n(\mathbf{x})$ and θ_0 , divided by the variance $\mathbf{V}(T_n(\mathbf{x}))$ of $T_n(\mathbf{x})$. The square root of this expression, $(T_n(\mathbf{x}) - \theta_0) / (\mathbf{V}(T_n(\mathbf{x})))^{1/2}$, converges in distribution to a standard normal.

Consider the null hypothesis $H_0: \theta = \theta_0$. When this null hypothesis is true, the quadratic form $W(\mathbf{x}, \theta_0)$ has a limiting Chi-square distribution with k degrees of freedom. Then, a test of the hypothesis with a significance level α can be carried out by choosing a critical level c from the upper tail of the $\chi^2(k)$ distribution so that the tail has probability α , and rejecting H_0 when $W(\mathbf{x}, \theta_0) > c$. We term $W(\mathbf{x}, \theta_0)$ the *Wald statistic*.

Suppose an alternative $H_1: \theta = \theta_1$ to the null hypothesis is true. The power of the Wald test is the probability that the null hypothesis will be rejected when H_1 holds. But in this case, $n^{1/2} \cdot J^{1/2}(T_n(\mathbf{x}) - \theta_0) = n^{1/2} \cdot J^{1/2}(T_n(\mathbf{x}) - \theta_1) + n^{1/2} \cdot J^{1/2}(\theta_1 - \theta_0)$, with the first term converging in distribution to

$N(0, I)$. For fixed $\theta_1 \neq \theta_0$, the second term blows up. This implies that the probability that $n^{1/2} \cdot J^{1/2} (T_n(\mathbf{x}) - \theta_0)$ is small enough to accept the null hypothesis goes to zero, and the power of the test goes to one. A test with this property is called *consistent*, and consistency is usually taken to be a minimum requirement for a hypothesis testing procedure to be statistically satisfactory. A closer look at the power of a test in large samples is usually done by considering what is called *local power*. Suppose one takes a sequence of alternatives to the null hypothesis that get closer and closer to the null as sample size grows. Specifically, consider $H_1: \theta = \theta_0 + \lambda/n^{1/2}$. For this sequence of alternatives, the term $n^{1/2} \cdot J^{1/2} (\theta_1 - \theta_0) = J^{1/2} \delta$ is a constant, and we have the result that $n^{1/2} \cdot J^{1/2} (T_n(\mathbf{x}) - \theta_0) \rightarrow_d N(J^{1/2} \lambda, I)$. This implies that $(T_n(\mathbf{x}) - \theta_0)' (nJ) (T_n(\mathbf{x}) - \theta_0)$, the Wald statistic, converges in distribution to a noncentral Chi-square distribution with k degrees of freedom and a noncentrality parameter $\lambda' J \lambda$. The local power of the test is the probability in the upper tail of this distribution above the critical level c for the Wald statistic. The local power will be a number between zero and one which provides useful information on the ability of the test to distinguish the null from nearby alternatives. In finite sample applications, the local power approximation can be used for a specific alternative θ_1 of interest by taking $\lambda = n^{1/2} \cdot (\theta_1 - \theta_0)$ and using the noncentral Chi-square distribution as described above.

In practice, we do not know the Fisher Information J exactly, but must estimate it from the sample by

$$(10) \quad J_{en} = \mathbf{E}_n[\nabla_{\theta} l(\mathbf{x}, T_n)] [\nabla_{\theta} l(\mathbf{x}_i, T_n)]' \equiv n^{-1} \sum_{i=1}^n [\nabla_{\theta} l(\mathbf{x}_i, T_n)] [\nabla_{\theta} l(\mathbf{x}_i, T_n)]'.$$

The expression in (10) is termed the *outer product of the score* $\nabla_{\theta} l(\mathbf{x}_i, T_n)$ of an observation. When there is a single parameter, this reduces to the square of $\nabla_{\theta} l(\mathbf{x}_i, T_n)$; otherwise, it is a $k \times k$ array of squares and cross-products of the components of $\nabla_{\theta} l(\mathbf{x}_i, T_n)$. From the theorem in Chapter 6.4, $J_{en} \rightarrow_p J$, and the rule 1.17 in Chapter 4 implies that replacing J by J_{en} in the Wald test statistic does not change its asymptotic distribution.

In the discussion of maximum likelihood estimation in Chapter 6.4 and the proof of its asymptotic normality, we established that when θ_0 is the true parameter,

$$(11) \quad n^{1/2} \cdot (T_n(\mathbf{x}) - \theta_0) = J^{-1} \cdot \nabla_{\theta} L_n(\mathbf{x}, \theta_0) / n^{1/2} + o_p(1);$$

that is, the difference of the maximum likelihood estimator from the true parameter, normalized by $n^{1/2}$, equals the normalized score of the likelihood at θ_0 , transformed by J^{-1} , plus asymptotically negligible terms. If we substitute (11) into the Wald statistic, we obtain $LM(\mathbf{x}, \theta_0) = W(\mathbf{x}, \theta_0) + o_p(1)$, where

$$(12) \quad LM(\mathbf{x}, \theta_0) = [\nabla_{\theta} L(\mathbf{x}, \theta_0)]' (nJ)^{-1} [\nabla_{\theta} L(\mathbf{x}, \theta_0)].$$

The statistic (12) is called the *Lagrange Multiplier (LM) statistic*, or the *score statistic*. The name Lagrange Multiplier comes from the fact that if we maximize $L_n(\mathbf{x}, \theta)$ subject to the constraint $\theta_0 - \theta = 0$ by setting up the Lagrangian $L_n(\mathbf{x}, \theta) + \lambda(\theta_0 - \theta)$, we obtain the first order condition $\lambda = \nabla_{\theta} L_n(\mathbf{x}, \theta)$ and hence $LM(\mathbf{x}, \theta_0) = \lambda' (nJ)^{-1} \lambda$. Because $LM(\mathbf{x}, \theta_0)$ is *asymptotically equivalent* to the Wald statistic, it will have the same asymptotic distribution, so that the same rules apply for determining critical

levels and calculating power. The Wald and LM statistics will have different numerical values in finite samples, and sometimes one will accept a null hypothesis when the other rejects. However, when sample sizes are large, their asymptotic equivalence implies that most of the time they will either both accept or both reject, and that they have the same power. In applications, J in (11) must be replaced by an estimate, either J_{en} from (10), or $J_{oen} = \mathbf{E}_n[\nabla_{\theta} l(\mathbf{x}, \theta_o)] [\nabla_{\theta} l(\mathbf{x}, \theta_o)]'$ in which the score is evaluated at the hypothesized θ_o . Both converge in probability to J , and substitution of either in (12) leaves the asymptotic distribution of the LM statistic unchanged. A major advantage of the LM form of the asymptotic test statistic is that it does not require that one compute the estimate $T_n(\mathbf{x})$. Computation of maximum likelihood estimates can sometimes be difficult. In these cases, the LM statistic avoids the difficulty.

The *generalized likelihood ratio criterion* was suggested in a number of simple tests of hypotheses as a good general procedure for obtaining test statistics. This method rejects H_o if

$$(13) \quad \kappa < \max_{\theta} L_n(\mathbf{x}, \theta) - L_n(\mathbf{x}, \theta_o),$$

where κ is a constant that is adjusted to give the desired significance level for the test. A Taylor's expansion of $L_n(\mathbf{x}, \theta_o)$ about $T_n(\mathbf{x})$ yields

$$(14) \quad L_n(\mathbf{x}, T_n(\mathbf{x})) - L_n(\mathbf{x}, \theta_o) = \nabla_{\theta} L_n(\mathbf{x}, T_n(\mathbf{x})) \cdot (T_n(\mathbf{x}) - \theta_o) - (T_n(\mathbf{x}) - \theta_o)' \nabla_{\theta\theta} L_n(\mathbf{x}, \theta_{en}) (T_n(\mathbf{x}) - \theta_o),$$

where θ_{en} is between θ_o and θ_n . But $\nabla_{\theta} L_n(\mathbf{x}, T_n(\mathbf{x})) = 0$. Under the regularity conditions in Chapter 6.4, $\nabla_{\theta\theta} L_n(\mathbf{x}, \theta_{en})/n \rightarrow_p J$. (To make the last statement rigorous, one needs to either establish that the convergence in probability of $\nabla_{\theta\theta} L_n(\mathbf{x}, \theta)/n$ to $J(\theta)$ is uniform in θ , or expand $\nabla_{\theta\theta} L_n(\mathbf{x}, \theta_{en})/n$ to first order about θ_o and argue that the first term goes in probability to $-J$ and the second term goes in probability to zero.) Then, $LR(\mathbf{x}, \theta_o) = 2 \cdot [L_n(\mathbf{x}, T_n(\mathbf{x})) - L_n(\mathbf{x}, \theta_o)]$, termed the *likelihood ratio statistic*, satisfies

$$(15) \quad LR(\mathbf{x}, \theta_o) = (T_n(\mathbf{x}) - \theta_o)' (nJ) (T_n(\mathbf{x}) - \theta_o) + o_p(1) \equiv W(\mathbf{x}, \theta_o) + o_p(1),$$

and the LR statistic is asymptotically equivalent to the Wald statistic. Therefore, the LR statistic will be asymptotically distributed Chi-square with k degrees of freedom, where k is the dimension of θ_o , and its local power is the same as that of the Wald statistic, and calculated in the same way.

The major advantage of the LR statistic is that its computation requires only the values of the log likelihood unrestricted and with the null imposed; it is unnecessary to obtain an estimate of J or perform any matrix calculations. We conclude that the trinity consisting of the Wald, LM, and LR statistics are all asymptotically equivalent, and provide completely substitutable ways of testing a hypothesis using a large sample approximation.