

1

Introduction

1.1 PREAMBLE

Stochastics prevail in every walk of life and science, amidst uncertainty, unpredictability, unknowability, and improbability to shear diversity on one hand, and inequality, inexplicability, and incompatibility to impossibility on the other hand. Stochastics distorts transparency of inherent *determinacy* so that the Trojan horses of *statistical inference* are often invoked for drawing valid and interpretable conclusions from pertinent experimental or observational data. Fathoming the depth of stochastics along with the latent chaos or chance elements (i.e., uncontrollable variation) is a challenging task that statistical inference shares intricately with *statistical decision theory* (SDT), for theoretical foundation, and *statistical methodology* for effective applications. In real-life problems, often, there is an additional layer of complexity: *Constraints* (latent or not) of diverse types mar the underlying determinacy (even if that exists), and, hence, there may be a genuine need to sort out the underlying determinacy from its superimposed stochastics; we need to focus on such logical restraints, and incorporate them mathematically as well as statistically in contemplated decision making processes. *Constrained statistical inference* (CSI) has its roots in this complex.

With the ever-expanding horizons of the domain of statistical inference, from the tenuous simple parametric models (such as the binomial, Poisson, normal, Laplace, exponential laws) to more complex *functional parametric* models (such as in many *stochastic processes*), to *semiparametrics* to the omnibus *nonparametrics*, it became necessary to recharge its arsenals with the battery of compatible developments in SDT, and thereby regulate the needed methodology to meet the demand of the vast

interdisciplinary fields of applications. Within this diversity, often, constraints crop up either as (partial, (non-)linear, total, or even implicit) *ordering* on the so-called parameter space(s), or as (linear or not) *inequalities* on contemplated parameters or functionals; they may also show up on the observations or experimental outcomes, termed the *responses*. In this sense, statistical constraints have their abode in the *decision* as well as the *sample* spaces. Albeit, such constrained universes could be portrayed in matching abstractions from the mathematical universe with SDT tones added for stochastic resolutions, for a more down-to-earth, comprehensive treatise of CSI, we shall adopt a pedestrian's path through the alleys of simple parametrics to functional parametrics to the vast domain, *beyond parametrics*. We specifically focus here on the constraints in the form of ordering and inequality restraints of diverse types, and keep in mind the vast field of applications. Perhaps it would be more motivating, and convincing too, to portrait first some simple and yet genuine inequality and order constrained inference problems as illustrative examples, and then step out to a more general prescription for a broader class of problems encountered in CSI. This is advocated in Section 1.2. The final section in this chapter attempts to capture the basic organization and coverage of the present treatise of CSI.

1.2 EXAMPLES

To briefly indicate an important feature of constrained statistical inference, let us first consider the simple example of testing $\mu = 0$ based on a random sample from the univariate normal distribution, $N(\mu, 1)$. Suppose that μ is known to be nonnegative. Then, we know from elementary statistics that a one-sided test with the sample mean as the test statistic is more powerful than a two-sided one. Intuitively, we would expect to do better if we know that $\mu \geq 0$ and if such a constraint is incorporated in the statistical inference procedure. We would expect a similar trend in multiparameter situations as well. In other words, if we know that parameters are restricted by some constraints then it is reasonable to expect that we should be able to do better by incorporating such additional information than by ignoring them. While this appears reasonable, it may or may not be an easy task.

The main theme of this book is incorporation of inequality and order restrictions on the parameters of the statistical model, and our objective is to provide a systematic development of the relevant statistical theory, motivated and illustrated by examples. In the rest of this section, we provide a broad range of motivating examples.

Example 1.2.1 Ordered treatment means in one-way layout

An experiment was conducted to evaluate the effect of exercise on the age at which a child starts to walk. Let Y denote the age (in months) at which a child starts to walk; the data on Y are given in Table 1.1. (The original experiment consisted of another treatment; however, here we consider only three treatments for simplicity. The data and analyses involving all four treatments will be given in a later chapter.)

The first treatment group received a special walking exercise for 12 minutes per day beginning at age 1 week and lasting 7 weeks. The second group received daily

Table 1.1 The age at which a child first walks

Treatment (i)		Age (in months)			n_i	\bar{y}_i	μ_i
1		9.00	9.50	9.75	10.00	9.50	μ_1
2		11.00	10.00	11.75	10.50	6	μ_2
3		13.25	11.50	12.00	13.50	1.50	μ_3

Reprinted from: *Science*, Volume 176, P. R. Zelazo, N. A. Zelazo, and S. Kolb, Walking in the newborn, Pages 314–315, Copyright (1972), with permission from AAAS.

$\mu_i =$ Mean age (in months) at which a child starts to walk.

In the traditional analysis of variance (ANOVA), one would usually test

$$H_0 : \mu_1 = \mu_2 = \mu_3 \quad \text{against} \quad H_2 : \mu_1, \mu_2, \text{ and } \mu_3 \text{ are not all equal.}$$

However, let us suppose that the researcher was prepared to assume that the walking exercises would not have the negative effect of increasing the mean age at which a child starts to walk, and it was desired that this additional information be incorporated to improve on the statistical analysis. For illustrative purposes, let us suppose that the researcher wishes to incorporate the information $\mu_1 \leq \mu_2 \leq \mu_3$. In this case, the testing problem is

$$H_0 : \mu_1 = \mu_2 = \mu_3 \text{ vs } H_1 : \mu_1 \leq \mu_2 \leq \mu_3 \text{ and } \mu_1, \mu_2, \text{ and } \mu_3 \text{ are not all equal.}$$

Clearly, the usual ANOVA where one would test

$$H_0 : \mu_1 = \mu_2 = \mu_3 \text{ against } H_2 : \mu_1, \mu_2, \text{ and } \mu_3 \text{ are not all equal,}$$

fails to incorporate the additional information $\mu_1 \leq \mu_2 \leq \mu_3$. Therefore, one would expect that we should be able to do better than the traditional *F*-test. In the next chapter, we illustrate a method that is specifically designed for this testing problem; the method to be illustrated, which we shall call the \bar{F} -test, is easy to understand and use because it uses essentially the same idea as that which underlies the traditional *F*-test and modifies it to incorporate additional information such as $\mu_1 \leq \mu_2 \leq \mu_3$.

Let us make a few remarks about this example. If the objective of the experiment was to establish that the special walking exercise results in a reduction in the mean age at which a child starts to walk, the testing problem needs to be formulated differently; for example, it may be formulated as test of

$$H_a : \mu_1 < \mu_2 < \mu_3 \text{ does not hold} \quad \text{vs} \quad H_b : \mu_1 < \mu_2 < \mu_3.$$

This is a different statistical inference problem and it will be discussed in later examples.

To find out which treatments were effective and which were not would involve multiple tests and require Bonferroni-type adjustment to control overall type-I error; this is also likely to require large samples to achieve moderate power.

A test of H_0 against H_2 in one-way ANOVA, for example the F -test, is also a valid procedure for testing H_0 against H_1 as well. Here valid means that a 5% level test of H_0 against H_2 is also a 5% level test of H_0 against H_1 . This is because the critical value and the size of a test are computed under the null hypothesis, and the null hypothesis for H_0 vs H_1 and for H_0 vs H_2 are the same. Therefore, the validity of the test is not an issue in this example.

The ordering on the μ_i 's in H_1 is an assumption. Even if the true values of $\{\mu_1, \mu_2, \mu_3\}$ do not satisfy the order restrictions in H_1 , it is still possible for some test of H_0 against H_1 to reject H_0 . In this case, while rejecting H_0 would be the correct statistical decision, accepting H_1 as true would be a wrong conclusion.

In general, when we apply a test of hypothesis, it is assumed that either H_0 or H_1 is true. For this reason, $H_0 \cup H_1$ can be called the *maintained hypothesis*; in this example, the maintained hypothesis is $\mu_1 \leq \mu_2 \leq \mu_3$. It will be seen in the next chapter, that tests of H_0 against H_1 can be carried out easily.

Example 1.2.2 Relationship between *El Niño* and hurricanes

El Niño refers to unusually warm ocean currents in the Pacific that appear around Christmas time and last for several months. Monsoon rains in the central Pacific and droughts and forest fires in Indonesia and Australia have been linked to *El Niño*. A hypothesis concerning *El Niño* is the following (Kitchens (1998) page 812):

H : Warm phase of *El Niño* suppresses hurricanes and cold phase encourages.

The data¹ in Table 1.2 provide information on the numbers of hurricanes from 1950 to 1995. In this context, different types of hypothesis testing problems may arise depending on how the question of interest is formulated.

Let us first formulate the basic model as a simple one-way classification. Let *El Niño* be the factor with three levels : Cold ($i = 1$), Neutral ($i = 2$), and Warm ($i = 3$). Let Y_{ij} denote the number of hurricanes and let

$$Y_{ij} = \mu_i + e_{ij},$$

where μ_i is the expected number of hurricanes ($i = 1, 2, 3$).

Let us suppose that the hypothesis H is based on various conjectures and scientific reasoning, and it is of interest to know whether or not there is any evidence against H . In this case, the null and alternative hypotheses can be stated as

$$H_0 : \mu_1 \geq \mu_2 \geq \mu_3 \quad \text{and} \quad H_1 : \mu_1 \geq \mu_2 \geq \mu_3 \text{ does not hold.}$$

A feature of this inference problem is that inequality constraints are present in the null hypothesis; this is different from the usual ANOVA problem where the null hypothesis would typically be of the form $\mu_1 = \mu_2 = \mu_3$.

Table 1.2 Effect of *El Niño* on hurricanes

	Y	50	51	52	53	54	55	56	57	58	59	60	61
H	Y	11	8	6	8	9	4	3	7	7	7	4	8
E	H	c	w	n	w	c	n	w	n	n	n	n	c
E	E	c	c	w	w	n	n	n	c	c	c	w	c
E	H	3	7	6	4	7	6	4	12	5	6	3	4
E	E	n	n	c	w	n	c	n	c	c	c	w	c
E	H	4	6	6	5	5	5	9	7	2	3	5	7
E	E	c	c	w	w	n	n	n	w	w	w	n	c
E	H	4	3	5	7	8	4	3	4	3	3	10	10
E	E	w	w	c	c	n	w	w	w	w	w	w	c

Yi (64 refers to 1964); H: number of hurricanes; E: *El Niño*; c: cold; w: warm; n: neutral.

Source: The original source of these data is the National Hurricane Center. The data in this table were extracted from the table on page 812 in Kitchens (1998).

On the other hand, if the interest is to establish that there is sufficient evidence to support the hypothesis H , then the testing problem is

$$H_0 : \mu_1 > \mu_2 > \mu_3 \text{ does not hold} \quad \text{vs} \quad H_1 : \mu_1 > \mu_2 > \mu_3 \text{ holds.}$$

The one-way model $Y_{ij} = \mu_i + e_{ij}$ used in the foregoing two examples is a simple one with no other covariates. If \mathbf{x}_{ij} is a vector of covariates, then the model could be of the form²

$$Y_{ij} = \mu_i + \mathbf{x}_{ij}^T \boldsymbol{\beta} + e_{ij};$$

the hypotheses of interest are still the same. These inference problems come under the general framework of the linear model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{E}.$$

Within the context of such a linear model, three testing problems of interest are:

1. $H_0 : \mathbf{R}\boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \mathbf{R}\boldsymbol{\theta} \geq \mathbf{0}, \mathbf{R}\boldsymbol{\theta} \neq \mathbf{0}$
2. $H_0 : \mathbf{R}\boldsymbol{\theta} \geq \mathbf{0}$ against $H_1 : \mathbf{R}\boldsymbol{\theta} \not\geq \mathbf{0}$
3. $H_0 : \mathbf{R}\boldsymbol{\theta} > \mathbf{0}$ does not hold against $H_1 : \mathbf{R}\boldsymbol{\theta} > \mathbf{0}$

where \mathbf{R} is a given fixed matrix that does not depend on $\boldsymbol{\theta}$. The first two types of problems will be discussed in Chapters 3 and 4, and the third type will be discussed in later chapters under the topic *sign testing*. ■

¹The authors are grateful to the National Hurricane Center for help with these data.

²Vectors and matrices are shown in bold and the superscript T denotes transpose.

Example 1.2.3 Comparison of two treatments when the response is multivariate (Pocock et al. (1987))

Seventeen patients with asthma or chronic obstructive airways disease entered a trial to evaluate an inhaled active drug versus placebo. The experiment was designed as a randomized, double-blind, cross-over trial. The main purpose was to study the possibly harmful effects of drugs on lung mucociliary clearance. Analysis of those results produced no evidence of harm. In addition, standard respiratory function measures were taken at the end of both treatment periods. These were peak respiratory flow rates (*PEFR*), forced expiratory volume (*FEV*), and forced vital capacity (*FVC*), the latter two being expressed as a percentage of the predicted value for that patient's age, gender, and height in the normal population. Let us define the following variables for the changes (i.e., drug - placebo) in the three measures:

$$X_1 = \text{change in FEV},$$

$$X_2 = \text{change in FVC},$$

$$\text{and } X_3 = \text{change in PEFR}.$$

Let

$$X = (X_1, X_2, X_3)^T \quad \text{and} \quad \mu = (\mu_1, \mu_2, \mu_3)^T$$

where $\mu_i = E(X_i)$ for $i = 1, 2, 3$. A question of interest was whether the addition of the inhaled drug could further improve respiratory function. For each measure, there was no sign of period or carry-over effects. Therefore, the univariate analyzes of drug versus placebo was performed using paired *t*-test (see Table 1.3).

The estimated correlation matrix of $(\bar{X}_1, \bar{X}_2, \bar{X}_3)$ reported in Pocock et al. (1987) is

$$\begin{pmatrix} 1 & 0.95 & 0.219 \\ 0.95 & 1 & 0.518 \\ 0.219 & 0.518 & 1 \end{pmatrix}.$$

All three measures showed a mean improvement on active drug but none achieved statistical significance at the 5% level. Thus, Bonferroni correction would lead to a conclusion of no evidence of improvement for the active drug. In any case, such a procedure is not the most powerful method in this setting. What is required is a

multivariate statistical procedure to evaluate the collective evidence of the overall benefit of the drug. Therefore, the null and alternative hypotheses take the form

$$H_0 : \mu = 0 \quad \text{and} \quad H_1 : \mu \geq 0,$$

respectively, where

$$\mu \geq 0 \text{ means that } \mu_i \geq 0 \text{ for every } i \text{ and } \mu \neq 0.$$

It would be helpful to point out an important difference between testing $H_0 : \mu = 0$ against $H_1 : \mu \geq 0$ and performing the three tests,

$$H_{01} : \mu_1 = 0 \text{ vs } H_{11} : \mu_1 > 0,$$

$$H_{02} : \mu_2 = 0 \text{ vs } H_{12} : \mu_2 > 0,$$

$$\text{and } H_{03} : \mu_3 = 0 \text{ vs } H_{13} : \mu_3 > 0$$

with Bonferroni correction. For illustrative purposes let us suppose that the standard errors, $se(\bar{X}_i)$, in the three columns of Table 1.3 are 1.0, 10.0, and 0.2, respectively. In this case, the three Bonferroni corrected tests would conclude that the effects on *FEV* and *PEFR* are significant, but not that on *FVC* at 5% level of significance. A test of $H_0 : \mu = 0$ against $H_1 : \mu \geq 0$ would also reject H_0 and accept H_1 ; however, this does not say which of *FEV*, *PEFR*, and *FVC* is significant and which is not. Therefore, the two testing procedures answer different questions.

Note that the parameter space for μ has $2^3 = 8$ quadrants, and the Hotelling's *T*² looks for departure from the null value $\mu = 0$ into all of the 2^3 quadrants, while a test against $H_1 : \mu \geq 0$ looks for departure into just one quadrant and hence would be more targeted to the specific objective of this study.

In a later chapter we will develop the theory for testing

$$H_0 : R\mu = 0 \text{ against } H_1 : R\mu \neq 0, R_1\mu \geq 0$$

based on a sample from $N(\mu, V)$ where R_1 is a submatrix of R and V is unknown. These results can be used to deal with problems similar to those in this example. ■

Example 1.2.4 Testing the validity of Liquidity Preference Hypothesis [an ordered hypotheses on the mean of a multivariate normal] (Richardson et al. (1992))

Let $H(r, t)$ denote the rate of return during the time period t to $t + 1$ from a bill with maturity r (i.e., at time t , the maturity date of the investment/bill is $t + r$). At time t , the term premium on a bill with maturity r is defined as

$$\theta_r = E\{H(r, t) - H(1, t)\}$$

where the expectation is taken with respect to the information available to the investors at time t . Thus, the term premium is the additional rate of return that the investors expect at time t for investing in a longer-term bill rather than in a short-term one that would mature at the end of next time point. The *Liquidity Preference Hypothesis*

Table 1.3 Effect of drug on chronic respiratory disease

Variable	X_1	X_2	X_3
Sample Mean (i.e. \bar{X}_i)	7.56	4.81	2.29
s.e. (\bar{X}_i)	18.53	10.84	8.51
t-statistic	1.63	1.77	1.11

Reprinted from: *Biometrics*, Vol 43, S. J. Pocock, N. L. Geller, and A. A. Tsiatis, The analysis of multiple endpoints in clinical trials, Pages 487–498, Copyright (1987), with permission from Blackwell Publishing.

Table 1.4 Term premia of T-bills with up to 11 months to maturity for Aug/1964-Nov/1990

Average premium	t -statistics for $\theta_r - \theta_{r-1}$	Average premium	t -statistics for $\theta_r - \theta_{r-1}$
0.030	$\theta_3 - \theta_2$ 0.056	6.78 -0.12	θ_7 θ_8
0.055	$\theta_4 - \theta_3$ $\theta_5 - \theta_4$	0.070 0.093	$\theta_8 - \theta_7$ $\theta_9 - \theta_8$
0.070	$\theta_6 - \theta_5$	0.80 -0.80	$\theta_{10} - \theta_9$
0.073	$\theta_7 - \theta_6$	0.071 0.077	$\theta_{11} - \theta_{10}$ $\theta_{11} - \theta_9$

Reprinted from *Journal of Financial Economics*, Volume 31, M. Richardson, P. Richardson, and T. Smith, The monotonicity of the term premium, pp 97–105, Copyright(1992), with permission from Elsevier.

(*LPH*), which plays an important role in economics, says that the term premium, θ_r , is a nonincreasing function of r . The basic argument supporting the hypothesis is that longer term bills are more risky than short term ones, and hence investors expect higher return for higher risk. There has been a series of publications in which the validity of this hypothesis was investigated (see Richardson et al. (1992) and the references therein). Table 1.4 provides summary data and some relevant statistics that have been used in the debates concerning the validity of *LPH*.

According to *LPH*, $\theta_r - \theta_s$ is nonnegative for $r > s$. Yet the estimates of $\theta_4 - \theta_3, \theta_7 - \theta_6, \theta_{10} - \theta_9$, and $\theta_{11} - \theta_9$ are negative and the last two are large compared with 5% level critical values from a standard normal distribution (see Table 1.4). This led to interesting debates about the validity of *LPH*. If we were to examine all possible pairwise differences then there is a good chance that some “large” negative estimates of $\theta_r - \theta_{r-1}$ would be observed. In any case, multiple tests of all possible pairwise differences in the θ_r s is not quite the appropriate statistical method for this problem, although the t -statistics in Table 1.4 are useful. The problem needs to be formulated as a test of

$$H_0 : \theta_{11} \geq \theta_{10} \geq \dots \geq \theta_2 \quad \text{against} \quad H_1 : H_0 \text{ does not hold.} \quad (1.1)$$

As in Example 1.2.2, a feature of this testing problem is that inequalities are present in the null hypothesis.

In a later chapter, we will develop the theory for testing $\theta \geq 0$ against $\theta \not\geq 0$ based on observations from \mathbf{X} where \mathbf{X} has the multivariate normal distribution, $N(\boldsymbol{\theta}, \mathbf{V})$; \mathbf{V} is either known or an estimate of it is available. These results can be used for testing (1.1). A test of (1.1) would say whether or not there is sufficient evidence against

$$H_0 : \theta_{11} \geq \theta_{10} \geq \dots \geq \theta_2;$$

when there is such evidence against H_0 , the test would *not* provide information as to which of the inequalities in the null hypothesis fails to hold. On the other hand, multiple tests on the pairwise differences with Bonferroni adjustment, attempt to identify which of the pairwise differences are positive; thus the later approach attempts to answer a much deeper question, one that is not really required for testing the validity of *LPH*. Therefore, in this example, it would be unnecessary to carry out multiple tests on pairwise differences; a set of multiple tests is likely to have low

power compared with one that is specifically designed for testing H_0 against H_1 in (1.1). If there are a large number of observations, then multiple tests can be used for the testing problem (1.1).

In the foregoing discussions, the problem was formulated as

$$\mathbf{X} \sim N(\boldsymbol{\beta}, \mathbf{V}) \text{ where } \boldsymbol{\beta} = (\theta_2, \dots, \theta_{11})^T,$$

and the estimate of $\boldsymbol{\beta}$ was the sample mean of 313 observations on \mathbf{X} . It is quite possible that in a more elaborate model, $(\theta_2, \dots, \theta_{11})$ may appear as some parameters in a regression type model. The theory for such problems will also be developed in a later chapter. ■

Example 1.2.5 Test against a set of one-sided hypotheses in a linear regression model (Wolak (1989))

Consider the following double-log demand function:

$$\log Q_t = \alpha + \beta_1 \log PE_t + \beta_2 \log PG_t + \beta_3 \log I_t + \gamma_1 D_{1t} + \gamma_2 D_{2t} + \gamma_3 D_{3t} + \epsilon_t,$$

where

- Q_t = aggregate electricity demand,
- PE_t = average price of electricity to the residential sector,
- PG_t = price of natural gas to the residential sector,
- I_t = income per capita,
- and D_{1t}, D_{2t} , and D_{3t} are seasonal dummy variables.

The error terms $\{\epsilon_t\}$ may or may not be independent; for example, they may satisfy an AR(1) process. Prior knowledge suggests that $(-\beta_1, \beta_2, \beta_3) \geq (0, 0, 0)$. A typical model selection question that arises is whether or not the foregoing model provides a better fit than the simpler model,

$$\log Q_t = \alpha + \gamma_1 D_{1t} + \gamma_2 D_{2t} + \gamma_3 D_{3t} + \epsilon_t.$$

This would require a test of

$$H_0 : (\beta_1, \beta_2, \beta_3) = (0, 0, 0) \text{ against } H_1 : (-\beta_1, \beta_2, \beta_3) \geq (0, 0, 0).$$

Strictly speaking, we should write the alternative hypothesis as $(-\beta_1, \beta_2, \beta_3) \geq (0, 0, 0)$. However, in what follows, we shall always interpret H_0 vs H_1 as H_0 vs $H_1 - H_0$ (i.e., $H_1 \setminus H_0$), so that the two hypothesis are disjoint. A test of

$$H_0 : (\beta_1, \beta_2, \beta_3) = (0, 0, 0) \text{ vs } H_2 : (-\beta_1, \beta_2, \beta_3) \neq (0, 0, 0)$$

is not the best for this testing/model-selection problem, although this is the one used by most standard statistical software packages. ■

Example 1.2.6 Multivariate ANOVA (Dietz (1989))

Vinylidene fluoride is suspected of causing liver damage. An experiment was carried out to evaluate its effects. Four groups of 10 male Fischer-344 rats received, by inhalation exposure, one of several dosages of vinylidene fluoride. Among the response variables measured on the rats were three serum enzymes: SDH, SGOT, and SGPT. Increased levels of these serum levels are often associated with liver damage. It is of interest to test whether or not these enzyme levels are affected by vinylidene fluoride. The data are given in Table 1.5 [the authors gratefully acknowledge receiving generous help from the National Toxicology Program of NIEHS regarding these data].

Let $\mathbf{Y}_{ij} = (Y_{ij1}, Y_{ij2}, Y_{ij3})^T$ denote the observations on the three enzymes for ith subject ($i = 1, \dots, 10$) in treatment j ($j = 1, \dots, 4$). Let θ_{jk} denote the mean (or median) response for j^{th} treatment (i.e., dose) and k^{th} variable and let $\boldsymbol{\theta}_j = (\theta_{j1}, \theta_{j2}, \theta_{j3})^T$ for $j = 1, \dots, 4$. Now, one formulation of the null and alternative hypotheses is

$$H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}_3 = \boldsymbol{\theta}_4 \quad \text{and} \quad H_1 : \boldsymbol{\theta}_1 \leq \boldsymbol{\theta}_2 \leq \boldsymbol{\theta}_3 \leq \boldsymbol{\theta}_4.$$

For this type of problem, general procedures for multivariate linear models can be used; these include methods based on likelihood, M -estimators, and ranks. The latter two methods have some robustness properties against nonnormal error distributions compared with the normal theory likelihood approach, similar to the case in the location scale model. ■

Just as in the linear regression model, inequality constrained hypotheses arise in more general models such as logistic regression models, proportional hazards models, loglinear models for categorical data, Generalized Linear Models, time-series models, and threshold models. Some examples of this type are provided below.

Example 1.2.7 Test against an ordered hypothesis on the regression parameters of a binomial response model (Piegorsch (1990), Silvapulle (1994))

An assay was carried out with the bacterium *E. coli* strain 343/358(+) to evaluate the genotoxic effects of 9-aminoacridine (9-AA) and potassium chromate (KCr). The data are presented in Table 1.6.

An objective of the study was to evaluate the synergistic effects of KCr when it is administered simultaneously with 9-AA. The null model of interest is known as the *Simple Independent Action* (SIA) model; the reasons for choosing this particular model for the null hypothesis are discussed in Piegorsch (1990) and the references therein. They show that the complementary loglinear model,

$$-\log(1 - \pi_{ij}) = \mu + \alpha_i + \beta_j + \eta_{ij}$$

with $\alpha_1 = \beta_1 = \eta_{1j} = \eta_{i1} = 0$ for every (i, j) , provides an appropriate framework to test for/against synergism. When $\eta_{ij} = 0 \forall (i, j)$, the SIA hypothesis holds. If η_{ij} are nonnegative and at least one of them is positive, then there is departure away from SIA towards synergism. Thus, the null and the alternative hypotheses take the forms

$$H_0 : \boldsymbol{\eta} = \mathbf{0} \quad \text{and} \quad H_1 : \boldsymbol{\eta} \geq \mathbf{0}, \quad (1.2)$$

Table 1.5 Serum enzyme levels in rats

Dosage	Enzyme	Rat within dosage								
		1	2	3	4	5	6	7	8	9
0	SDH	18	27	16	21	26	22	17	27	26
	SGOT	101	103	90	98	101	92	123	105	92
	SGPT	65	67	52	58	64	60	66	63	68
1500	SDH	25	21	24	19	21	22	20	25	24
	SGOT	113	99	102	144	109	135	100	95	89
	SGPT	65	63	70	73	67	66	58	53	58
5000	SDH	22	21	22	30	25	21	29	22	24
	SGOT	88	95	104	92	103	96	100	122	102
	SGPT	54	56	71	59	61	57	61	59	63
15000	SDH	31	26	28	24	33	23	27	24	28
	SGOT	104	123	105	98	167	111	130	93	99
	SGPT	57	61	54	56	45	49	57	51	48

Serum enzyme levels are in international units/liter; dosage of vinylidene in parts/million.

Reprinted from: Thirteen-week study of Vinylidene fluoride in F344 male rats, with permission from The National Toxicology Program, NIEHS; the data appeared in a report prepared by Linton Bionetics Inc. (1984). These data also appeared in Dietz (1989), where they were used to illustrate constrained inference.

9-AA(μm)	KCr(μm)			
	0	1.4	2.9	4.3
0	0/192	9/192	3/192	5/7
40	49/102	92/192	19/92	3/192

Note: 9/192 means 9 positive responses out of 192 trials.

Reprinted from: *Mutation Research*, Vol 171, La Velle, Potassium chromate potentiates frameshift mutagenesis in *E. coli* and *S. Typhimurium*, Pages 1-10, Copyright (1986), with permission from Elsevier.

This type of issue arises in other areas. For example, noncarcinogenic toxic reactions to therapeutic agents have been caused by interactions of compounds that were normally safe when given alone (see Piegorsch (1990)). A formulation of this phenomenon also involves sign constraints on the interaction terms similar to (1.2). In Chapter 4, we will consider general statistical models that include the generalized linear model in which the distribution of y_i belongs to an exponential family and that for some suitable link function, g , we have that

$$g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\eta}$$

where μ_i is the mean of y_i . The loglikelihood functions of such generalized linear models (and some partial likelihoods) are well behaved, for example, they are concave and smooth. Consequently, we can (for example, see Silvapulle (1994)) obtain elegant results for testing

$$H_0 : \mathbf{R}\boldsymbol{\eta} = 0 \text{ against } H_1 : \mathbf{R}\boldsymbol{\eta} \geq 0.$$

This collection of results can be used for the foregoing testing problem in the binomial response model. ■

Example 1.2.8 Testing that a treatment is better than the control when the responses are ordinal (Grove (1980), Silvapulle (1994))

The results of a trial involving two treatments for ulcer are given in Table 1.7. The objective is to test the hypothesis that Treatment B is better than Treatment A.

This testing problem can be formulated in several ways. One possible approach is to say that Treatment B is better than Treatment A, if the local odds ratios are all greater than or equal to 1 with at least one of them being greater than one; recall that the local odds ratios are defined as follows:

$$\phi_1 = \pi_{11}\pi_{22}/\{\pi_{21}\pi_{12}\}, \phi_2 = \pi_{12}\pi_{23}/\{\pi_{22}\pi_{13}\}, \text{ and } \phi_3 = \pi_{13}\pi_{24}/\{\pi_{23}\pi_{14}\}$$

where π_{ij} is the probability that the response of a given individual in the i th row falls in column j . Therefore, the null and alternative hypotheses are

$$H_0 : \boldsymbol{\phi} = \mathbf{1} \quad \text{and} \quad H_1 : \boldsymbol{\phi} \geq \mathbf{1}, \quad (1.3)$$

where $\boldsymbol{\phi} = (\phi_1, \phi_2, \phi_3)^T$ and $\mathbf{1} = (1, 1, 1)^T$. The restriction $\boldsymbol{\phi} \geq \mathbf{1}$ is known as the *likelihood ratio order*.

Table 1.7 Response of two groups to treatments to an ulcer

Treatment group	Response			Total
	Worse	Good	Very good	
A (Control)	12	10	4	6
B (Treatment)	5	8	8	11

Reprinted from: *Lancet*, Vol 259, R. Doll and F. Pygott, Factors influencing the rate of healing of gastric ulcers, Pages 171-175, Copyright (1952), with permission from Elsevier.

Table 1.8 The 2×2 tables formed by merging adjacent columns of Table 1.7

Treatment group	1		2		3
	Col 1	Col 2,3,4	Col 1,2	Col 3,4	
A	12	20	22	10	26
B	5	27	13	19	11

Odds ratio : γ_1 / γ_2

Reprinted from: *Biometrics*, Vol 50, M. J. Silvapulle, On tests against one-sided hypotheses in some generalized linear models, Pages 853-858, Copyright (1994), with permission from Blackwell Publishing.

Grove (1980) argues that (1.3) may be a rather stringent formulation of the hypothesis, and suggests the following formulation of the alternative hypothesis:

$$H_1 : \left(\sum_{j=1}^q \pi_{1j} \right) \left(\sum_{j=q+1}^C \pi_{2j} \right) \left\{ \left(\sum_{j=1}^q \pi_{2j} \right) \left(\sum_{j=q+1}^C \pi_{1j} \right) \right\}^{-1} \geq 1, q = 1, \dots, C - 1 \quad (1.4)$$

where $C = 4$ is the number of columns. To interpret this, consider the collection of 2×2 tables (see Table 1.8) that may be formed by merging the adjacent columns of Table 1.7. Let $\gamma_1, \gamma_2, \gamma_3$ denote the odds ratios for the three 2×2 tables. Now the alternative hypothesis in (1.4) is the same as

$$H_1 : \gamma_1 \geq 1, \gamma_2 \geq 1, \gamma_3 \geq 1;$$

this is not as stringent as the H_1 in (1.3). A detailed discussion of inference when the response variable is ordinal is provided in Chapter 6 where other formulations of the alternative hypothesis such as *hazard rate order*, *continuation odds ratio order*, and *stochastic order* will also be studied. ■

Example 1.2.9 A group sequential study (*Lipids and cardiovascular diseases*)

The National Heart, Lung and Blood Institute (USA) planned a multicenter clinical trial in 1971 to investigate the impact of lowering the blood cholesterol level on the risk for cardiovascular diseases. The study involved 3952 males, between the ages 35 and 60, chosen from a large pool of people having no previous heart problems but having an elevated blood cholesterol level. They were randomly divided into a placebo and a treatment group, the treatment group receiving some drug to lower the cholesterol level to about 160 or less. Based on cost and other medical considerations, a 12-year (1972-1984) study was planned, and the data coordination task was given to the University of North Carolina, Chapel Hill. The primary response variable was the failure time, while there were a large number of explanatory and covariables. Medical ethics prompted that if at any time during the study, the treatment is judged superior, the trial should be terminated, and all surviving persons should be switched to the treatment protocol. This clearly set the tone for *interim analysis*, looking at the accumulating data set regularly. It was decided to have a multiple one-sided

hypothesis testing setup, in a time-sequential framework: every quarter, there would be a statistical test, one-sided alternative, and if at any time a significance is observed, the trial would be concluded along with the rejection of the null hypothesis of no treatment effect (i.e., the cholesterol level in blood has no impact on heart problems). These test statistics were not independent, nor could they be based on simple normal or exponential laws. So, the first task was to develop a time-sequential hypothesis testing procedure for one-sided alternatives. For some of these details, we refer to Sen (1981, Ch.11) and Sen (1999a). Much of these technicalities have become a part of standard statistical interim analysis, and we shall refer to that again in Chapter 5. The restraints on the parameters (survival functions) were compatible with the medical and statistical considerations, and a conclusive result was announced in the 1984 issue of the *Journal of the American Medical Association*. ■

Example 1.2.10 Testing for/against the presence of random effects (Stram and Lee (1994))

Consider the well-known growth curve data set of Pothoff and Roy (1964). The data consists of measurements of the distance (in mm) from the center pituitary to the piterymaxillary fissure for 27 children (11 girls and 16 boys). Every subject has four measurements, taken at ages 8, 10, 12, and 14 years. The object is to model the growth in distance as a function of age and sex. To allow for individual variability in the growth function, the model being considered is

$$\mathbf{y}_i = (1, \text{Age}, \text{Sex})\boldsymbol{\alpha} + (1, \text{Age})\boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i$$

where \mathbf{y}_i is 4×1 for the 4 measurements of the i^{th} subject, $\text{Age} = (8, 10, 12, 14)^T$, $\text{Sex} = (a, a, a, a)^T$ where $a = 0$ for boys and 1 for girls. The parameter $\boldsymbol{\alpha}(4 \times 1)$ captures the fixed effects and $\boldsymbol{\beta}_i(2 \times 1)$ captures the random effect of subject i . Assume that

$$\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \boldsymbol{\beta}_i \sim N(\mathbf{0}, \Psi)$$

for some 2×2 positive semi-definite matrix Ψ , and that $\boldsymbol{\epsilon}_i$ and $\boldsymbol{\beta}_i$ are independent. More generally, the regression model can be written as

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\alpha} + \mathbf{Z}_i \boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i.$$

In the process of analyzing these data, one would be interested to know whether or not some random effects are present; equivalently, would a smaller model fit as well as a larger one? In this case, the parameter of interest is $\Psi = \text{cov}(\boldsymbol{\beta}_i)$, and the following testing problems are of interest:

1. $H_0 : \Psi = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$ vs $H_1 : \Psi$ is positive semi-definite.
 2. $H_0 : \Psi = \begin{pmatrix} \psi_{11} & 0 \\ 0 & 0 \end{pmatrix}$ and $\psi_{11} \geq 0$ vs $H_1 : \Psi$ is positive semi-definite.
- If $\Psi = \mathbf{0}$ then there are no random effects. If $\Psi = (\psi_{11}, 0 | 0, 0)$ then the individual variability is captured by a random intercept, and hence the regression planes are parallel for different subjects. If Ψ is positive definite then individual variability results in not only a random intercept but also a random slope parameter for Age.

More generally, consider a generalized linear model with linear predictor of the form

$$\mathbf{X}_i \boldsymbol{\alpha} + U_i \boldsymbol{\gamma}_i + \mathbf{Z}_i \boldsymbol{\beta}_i$$

where $\boldsymbol{\gamma}_i$ and $\boldsymbol{\beta}_i$ are random effects, for example

$$\boldsymbol{\gamma}_i \sim N(\mathbf{0}, D_1) \quad \text{and} \quad \boldsymbol{\beta}_i \sim N(\mathbf{0}, D_2).$$

Suppose that we wish to test

$$H_0 : D_1 = \mathbf{0} \quad \text{vs} \quad H_1 : D_1 \text{ is positive semi-definite.}$$

These hypotheses can be tested using Likelihood Ratio and Local Score tests, among others; the theory for this will be developed in Chapter 4. ■

Example 1.2.11 Sign testing problem: testing whether an identified treatment is the best

The problem here is to test whether or not a particular treatment, say A , is better than each of K other treatments. To illustrate an example where this arises, let A denote a drug formed by combining drugs B and C ; in the medical/pharmaceutical literature, A is known as *combination drug/therapy*. Since A is a combination of drugs, it may have a higher level of risk of adverse effects. Therefore, a standard requirement of regulatory authorities to approve a combination drug, such as A , is that it must be shown that the combination drug is *better* than each of its components (see Laska and Meisner, 1989). Let X_0, X_1 and X_2 be three random variables that characterize the effects of the drugs A, B , and C , respectively. Suppose that the distribution of X_i is associated with a parameter, μ_i , and assume that larger values of μ_i are more desirable. Let

$$\theta_1 = \mu_0 - \mu_1 \text{ and } \theta_2 = \mu_0 - \mu_2.$$

Now, the null and alternative hypotheses are

$$H_0 : \theta_1 \leq 0 \text{ or } \theta_2 \leq 0, \quad \text{and} \quad H_1 : \theta_1 > 0 \text{ and } \theta_2 > 0$$

respectively. In each of the previous examples, the alternative parameter space included its boundary as well. By contrast, in this example, the alternative parameter space does not include the boundary, and the null parameter space is the union of the second, third, and fourth quadrants, and the boundary of the first quadrant. These types of testing problems can be solved using *intersection-union tests*.

In this example, we assumed that there is only one primary variable of interest and hence X_i was univariate. There are clinical trials in which the assessment involves several variables (called, *multiple end points*); for example, interest may be on several symptoms simultaneously. Therefore, X_i and μ_i may be vectors. ■

Example 1.2.12 Combination drug for low-back pain and spasm (Snappin (1987))

Dolobid(D) and Flexeril(F) are approved drugs for spasm and pain relief respectively. It was conjectured that the new drug DF , a combination of D and F , is likely to be better. As was indicated in Example 1.2.11, a requirement of the regulatory authorities for the approval of a new combination drug such as DF is that it must be better than its components D and F . A study was conducted to evaluate the efficacy of DF relative to D only, F only, and Placebo (P). Patients entering the study were randomly assigned to one of the four treatments, DF , D , F , and P . At the end of the study period, each patient provided an evaluation of the treatment on a five-point ordinal scale ranging from Marked Improvement to Worsening of the condition (see Table 1.9). The objective is to test the claim that the combination drug is better than each of the other three treatments.

Let

$$DF \succ F \text{ denote } DF \text{ is better than } F;$$

this can be formulated using different odds ratios as in Example 1.2.8. Now, the inference problem may be formulated as test of

$$H_0 : \text{Not } H_1 \quad \text{vs} \quad H_1 : DF \succ F, DF \succ D, \text{ and } DF \succ P.$$

The inference problem can also be formulated in the following simpler form. Let μ_0, μ_1, μ_2 , and μ_3 denote parameters that capture the effects of the four treatments in Table 1.9; such parameters can be defined by introducing various parametric assumptions (see Snappin (1987)). Let $\theta_i = \mu_0 - \mu_i$ for $i = 1, 2, 3$. Now, one way of formulating the inference problem is test of

$$H_0 : \theta_1 \leq 0, \text{ or } \theta_2 \leq 0, \text{ or } \theta_3 \leq 0, \text{ vs } H_1 : \theta_1 > 0, \theta_2 > 0, \text{ and } \theta_3 > 0.$$

This problem will be discussed under intersection union tests. ■

Table 1.9 Comparison of Dolobid/Flexeril with Dolobid, Flexeril, and placebo.

Treatment	Degree of improvement				Total	
	Marked	Mod	Mild	No Change		
Dolobid/Flexeril	50	29	19	11	5	114
Dolobid alone	41	28	23	22	2	116
Flexeril alone	38	27	18	27	3	113
Placebo	35	20	23	37	4	119
Total	69	80	125			2067.7

Reprinted from: *Statistics in Medicine*, Vol 6, Snappin, Evaluating the efficacy of a combination therapy, Pages 656–665, Copyright (1987), with permission from John Wiley and Sons.

Example 1.2.13 Cross-over interaction in a multicenter clinical trial (Ciminera et al. (1993b))

A multicenter clinical trial was conducted to compare a drug with a placebo; the drug was a treatment to reduce blood pressure (see Table 1.10). This is an unpaired design (i.e., parallel group) within each center, and there were 12 centers. In the context of this example, we say that there is *cross-over interaction* between the treatment and center if the drug is beneficial in at least one center and harmful in at least one other. The question of interest is whether or not there is any evidence of crossover interaction.

The problem may be formulated as follows. Let δ_i denote the effect of the treatment in Center i ($i = 1, \dots, 12$). Then the null and the alternative hypotheses may take the form

$$H_0 : \delta \geq 0 \text{ or } \delta \leq 0, \quad \text{and} \quad H_1 : \delta \not\geq 0 \text{ and } \delta \not\leq 0,$$

respectively. In this formulation, the null hypothesis says that there is no cross-over interaction. It is also possible to interchange the roles of H_0 and H_1 in the testing problem so that the null and alternative hypotheses are H_1 and H_0 , respectively. In all of these, it may be of interest to replace δ_i by $\delta_i - \epsilon_i$ for some given ϵ_i so that $\delta_i - \epsilon_i > 0$ corresponds to a treatment effect that is clinically significant; such a modification through the introduction of small ϵ_i s helps to accommodate the fact that a cross-over interaction may not be important if it is not clinically significant. ■

Table 1.10 Comparison of a drug with a placebo in a multicenter trial ^a

i	n_{1i}	n_{2i}	df	\bar{z}_{1i}	\bar{z}_{2i}	$\hat{\delta}_i$	ss_i	$\hat{\delta}_i/se(\hat{\delta}_i)$
1	3	3	4	-7.0	-15.3	8.3	176.67	1.53
2	4	6	8	-13.0	-6.0	-7.0	180.00	-2.29
3	7	8	13	-3.0	-15.9	12.9	1492.87	2.33
4	7	8	13	-13.6	-11.0	-2.6	693.71	-0.69
5	6	8	12	-3.0	-11.8	8.8	345.50	3.04
6	7	6	11	2.7	-4.3	7.0	236.76	2.71
7	5	8	11	-2.2	-5.0	2.8	472.80	0.75
8	7	9	14	-8.3	-17.6	9.3	1080.93	2.10
9	8	8	14	-7.4	-12.2	4.8	507.38	1.59
10	8	9	15	-7.4	-6.7	-0.7	711.87	-0.21
11	4	4	6	2.0	-8.5	10.5	519.00	1.60
12	3	3	4	0.7	-12.0	12.7	116.67	2.88
Total	69	80	125					

^aThe within center sample means for treatments 1 and 2 are denoted by \bar{z}_{1i} and \bar{z}_{2i} , respectively; $\hat{\delta}_i = \bar{z}_{1i} - \bar{z}_{2i}$; ss_i = within group sum of squares for group i ; $se(\cdot)$ = standard error.

Reprinted from: *Statistics in Medicine*, Vol 12, Ciminera, Heyse, Nguyen, and Tukey, Tests for qualitative treatment-by-centre interaction using a "pushback" procedure, pages 1033–1045, Copyright (1993), with permission from John Wiley and Sons.

Example 1.2.14 Testing against the presence of ARCH effect in ARCH or ARCH-M models (Silvapulle and Silvapulle (1995) and Beg et al. (2001))

The AutoRegressive Conditional Heteroscedasticity (ARCH) model, introduced by Engle (1982), is one of the widely used econometric models, particularly in Financial Economics. This example involves two models in this family.

Let y_t denote the rate of return from an investment (portfolio) at time t ; here we are thinking of investment in stocks or bonds. The basic form of the ARCH model is

$$y_t = \mu + \beta^T x_t + \epsilon_t, \quad (1.5)$$

$$\epsilon_t | \mathcal{F}_{t-1} \sim N(0, h_t^2), \quad h_t^2 = \alpha_0 + \psi_1 \epsilon_{t-1}^2 + \dots + \psi_p \epsilon_{t-p}^2 \quad (1.6)$$

where \mathcal{F}_{t-1} is the collection of information up to time $(t-1)$ and x_t is a vector of exogenous variables. A main feature of this model is that periods of small (respectively, large) random fluctuations in y tend to be clustered. This is a frequently observed phenomenon in many high-frequency financial time series.

Another frequently observed phenomenon is that large random fluctuations in stock price would be seen as periods of high risk and this would in turn affect the stock price as investors demand higher returns for higher risk. To incorporate this, Engle et al. (1987) suggested that the specification in (1.5) be modified as

$$y_t = \mu + \beta^T x_t + \phi h_t + \epsilon_t, \quad (1.7)$$

Note that according to (1.7), the conditional variance h_t^2 has a direct effect on the mean of y_t ; in other words, the degree of risk, measured by h_t , affects the level of stock price and rate of return. The model (1.6)-(1.7) is known as the ARCH-in-Mean (ARCH-M) model.

In many empirical applications involving the ARCH and ARCH-M models, a question of basic interest is whether or not the conditional variance h_t^2 is a constant over time; or more generally, are some of the ψ parameters equal to zero. Since h_t^2 cannot be negative, ψ_1, \dots, ψ_p are all nonnegative. Therefore, if we were to test against the alternative that h_t^2 depends on time, then the testing problem takes the form

$$H_0 : \psi = 0 \text{ against } H_1 : \psi_1 \geq 0, \dots, \psi_p \geq 0. \quad (1.8)$$

For (1.5), the method that is usually applied by econometric software packages is a score test that ignores the constraints $\psi \geq 0$ in the alternative hypothesis. However, it is possible to improve on this by applying a method that incorporates the constraints. The theory for this will be developed in Chapter 4; one of the methods that will be developed is a simple generalization of Rao's Score test, which we shall call a local score test (see Silvapulle and Silvapulle (1995)). For (1.7), there are some additional theoretical issues that need to be taken care of because a parameter becomes unidentified under the null hypothesis; this is a nonstandard problem (see Beg, Silvapulle and Silvapulle (2001), and Andrews (2001)). ■

Now, we provide a few more examples that fall into the categories of the previous acceptance sampling, pages 295–300.

Example 1.2.15 Acceptance sampling (Berger (1982))

Table 1.11 lists specifications for upholstery fabric. Similar specifications are available for other products.

It is required to develop an acceptance sampling scheme for a batch that consists of a large number of units. Suppose that a random sample of n units are to be chosen from a batch and tested. The objective is to decide whether or not to accept the batch after testing the sample of n units. Note that some of the tests would destroy or change the fabric, and therefore once a piece of fabric has been used for one test, it may not be possible to use it for another test. The inference problem can be formulated as test of

$$H_0 : \theta \notin \Theta \quad \text{vs} \quad H_1 : \theta \in \Theta \quad (1.6)$$

where

$$\begin{aligned} \Theta &= \{\theta : \theta_1 \geq 50, \theta_2 \geq 6, \theta_3 \leq 5, \theta_4 \leq 2, \theta_5 \in A_1, \\ &\quad \theta_6 \in B_4, \theta_7 \in C_4, \theta_8 \in D_3, \theta_9 \in E_4, \theta_{10} \in P\}. \end{aligned}$$

First, suppose that each unit in the sample can be subjected to every test in the table of specifications. Typically, this would be a reasonable assumption if the properties of the product are uniform within each sampling unit and it can be divided into subunits for different tests. In this case, the statistical inference problem reduces to elementary inference on a single population proportion.

A different approach is required if a sampling unit cannot be subjected to all the different tests in the table of specifications. In this case, the *intersection union test*

Table 1.11 Standard specification for woven upholstery fabric—plain, tufted, or flocked

Variable	Requirement	Minimum standard formulation
Breaking strength	50 pounds	$\theta_1 \geq 50$
Tear strength	6 pounds	$\theta_2 \leq 6$
Dimensional change	5% shrinkage	$\theta_3 \leq 5$
	2% gain	$\theta_4 \leq 2$
		$\theta_5 \in A_1$
Surface abrasion (heavy duty)	15000 cycles	$\theta_6 \in B_4$
Colorfastness to Water	class 4	$\theta_7 \in C_4$
Crocking	class 4	$\theta_8 \in D_3$
Dry	class 3	$\theta_9 \in E_4$
Wet	class 4	
Light-40 AATCCF (fading units)		
Flammability	Pass	$\theta_{10} \in P$

Source: The original source of the information contained in this table is the Annual Book of ASTM Standards (1978), American Society for Testing and Materials, Philadelphia, part 32, page 717. This table is reproduced from *Technometrics*, Volume 24, R. L. Berger, Multiparameter hypothesis testing and acceptance sampling, pages 295–300.

provides a convenient solution to this problem. Essentially, a test at level- α accepts the batch if *each* specification is passed by a level- α test; there are no Bonferroni type adjustment because this does not come within the framework of what is usually known as multiple testing. ■

Example 1.2.16 Analysis of $2 \times 2 \times 2$ tables (Cohen *et al.* (1983a))

A random sample of 15 individuals who have been vaccinated against a certain disease, and another independent random sample of 15 individuals who have not been vaccinated against the same disease were chosen. The individuals were cross-classified with respect to the binary variables, Vaccination (yes, no), Immune Level (low, high), and Infection (well, sick). The data are given in Table 1.12.

Table 1.12 Effect of vaccination and immune level on infection

Infection	Vaccinated			& high immun. & low immun.	& high immun.
	Not vaccinated & low immun.	Vaccinated 3	Vaccinated 3		
Well	4	3	3	7	2
Sick	7	1	3	3	1

Reprinted from: *Journal of the American Statistical Association*, Vol 78, A. Cohen, C. Gatsionis, and J. Marden, Hypothesis testing for marginal probabilities in a $2 \times 2 \times 2$ contingency table with conditional independence, Pages 920-929, Copyright (1983), with permission from the American Statistical Association.

Assume that vaccination affects the risk of infection only through its effect on the immune level. It is required to test the validity of the claim

H_0 : vaccination improves immune level, which in turn reduces the risk of infection;

H_1 : vaccination *reduces* immune level, which in turn improves the risk of infection.

The terms *improves* and *reduces* lead to multiparameter one-sided hypotheses. This type of inference problem can be formulated as test of

$$H_0 : \alpha \notin A \text{ or } \beta \notin B \quad \text{vs} \quad H_1 : \alpha \in A \text{ and } \beta \in B$$

for some parameters α and β and suitably chosen sets A and B . Cohen *et al.* (1983a) developed a methodology for this type of problem and illustrated their results using these data. This example will be studied in Chapter 6. ■

Example 1.2.17 Meta-analysis (Cutler *et al.* (1991), Follmann (1996a))

A recent meta-analysis summarized the effect of sodium reduction on blood pressure from published clinical trials. An objective of the study was to provide an overall probability statement on the joint effect of sodium reduction on reducing both the systolic blood pressure (X_1) and diastolic blood pressure (X_2).

Table 1.13 provides the data for the meta analysis. In fact, the table provides considerable information that would lend to more extensive analyses. Let X be

Table 1.13 Randomized trials of sodium reduction

Sample size	Study length (months)	Reduction in urinary sodium SBP (se)	Reduction in BP (mm Hg) DBP (se)
Cross-over trials; hypertensive subjects			
15	1	98	6.7 (3.76)
19	1	76	10 (3.06)
18	1	56	0.5 (1.5)
12	1	105	5.2 (4.1)
40	1.5	72	0.8 (1.8)
20	1	82	8.0 (2.6)
9	1	76	9.7 (4.33)
88	2	67	3.6 (0.7)
Parallel studies with treatment and control groups; hypertensive subjects			
31/31	24	27	1.5 (4.6)
6/6	2	98	*
6/6	2	78	*
10/15	12	53	8.7 (10.23)
15/19	1.5	117	-1.8 (4.14)
15/15	2.3	89	13.3 (5.46)
18/12	3	171	2.0 (4.96)
44/50	6	58	-2 (2.84)
48/52	3	54	2.7 (2.2)
37/38	6	32	3.4 (1.7)
17/17	3	59	5.1 (1.42)
50/53	2	71	4.2 (0.9)
20/21	12	*	18.3 (4.35)
Cross-over trials; normotensive subjects			
20	0.5	170	2.7 (2.36)
113	2	69	0.6 (0.84)
35	0.9	74	1.4 (0.74)
31	1	60	0.5 (0.82)
172	1	130	3.5 (1.25)
Parallel studies with treatment and control groups; normotensive subjects			
19/19	0.5	117	3.0 (2.03)
174/177	36	16	-0.1 (1.00)

Reprinted from: *Hypertension*, Vol 17, Cutler, Follmann, Elliott, and Suh, An overview of randomized trials of sodium reduction and blood pressure, Pages 27-33, Copyright (1991), with permission from Wolters Kluwer Health.

the bivariate vector that consists of X_1 and X_2 . Let the mean of (X_1, X_2) be denoted by (θ_1, θ_2) . Now, the statistical inference problem can be formulated as test of $H_0 : (\theta_1, \theta_2) = (0, 0)$ against $H_1 : (\theta_1, \theta_2) \geq (0, 0)$. Different rows in Table 1.13 correspond to studies with different sample sizes. Therefore, we have independent observations X_1, \dots, X_n with a common mean θ but with different covariance matrices. ■

Example 1.2.18 Distribution-free analysis of covariance (Boyd and Sen (1986))

The nature of the problem in this example is that there are several treatments and some covariates. We are interested to test for equality against an ordered alternative of the treatment effects. However, there was adequate evidence to suggest that the error distribution is unlikely to be normal and, further, the effects of the covariates on the response variable are unlikely to be linear. Thus, some kind of distribution-free method is required.

A trial was conducted (by A.H. Robins Company) to compare the efficacy of three treatments for the reversal of anesthesia. For confidentiality purposes, the treatments were identified as 1, 2, and 3. The response variable, Y , is time elapsed from administration of treatment to completion of reversal. The covariates are $Depth$ = “depth of neuromuscular block at time of reversal (time treatment administered)” and Age = “the age (in years) of the patient”, the data are given in Table 1.14.

Based on various diagnostics such as plot of the residuals, there is adequate evidence that the distribution of Y at a fixed value of ($Depth$, Age) is not normal. Further, within each treatment, the dependence of Y on ($Depth$, Age) is unlikely to be linear. Therefore, the usual ANCOVA based on normal theory linear model is inappropriate.

The problem may be formulated as follows. Let $F(y - \theta_k | x)$ denote the distribution of Y for Treatment k at ($Depth$, Age) = x . The functional form of F and the form of dependence of a location parameter of Y on x are unknown. The null and alternative hypotheses are $H_0 : \theta_1 = \theta_2 = \theta_3$ and $H_1 : \theta_1 \leq \theta_2 \leq \theta_3$. Some general results based on ranks can be applied for this type of problems. ■

1.3 COVERAGE AND ORGANIZATION OF THE BOOK

Statistical Decision Theory (SDT) attempts to bring both wings of statistical inference under a common shade. In that process, it often becomes more abstract than what could be more easily routed to fruitful applications. Faced with this dilemma, we plan to treat the estimation theory and hypothesis testing more or less on a complementary basis, although at the end, in Chapter 8, we would attempt to bind them together by SDT strings. Constrained statistical inference (CSI) arises in many areas of applications, and some of these have been illustrated in the preceding section. From pure theoretical perspectives, optimal statistical decision procedures generally exist only under very specific regularity assumptions that are not that likely to be tenable in constrained stochastic environments. As such, in the sequel, we shall first trace the rather tenuous domain of CSI problems where finite-sample optimal decisions exist. In a way to get out of this restricted scope of CSI, we would spring up in a beyond parametric scenario and illustrate the variety of competing methodologies that attempt to capture the glimpses of CSI with due emphasis on robustness, validity, and efficacy considerations in a far more general setup (where exact optimality may not percolate). In the rest of this section, we attempt to provide a motivating coverage of the book.

Ours is not the first book or monograph in this specified field. We would like to acknowledge with profound thanks the impetus we have had from the earlier two precursors: Barlow et al. (1972), to be referred to in the sequel as *B4* (1972), and Robertson et al. (1988), those two being 16 years apart and otherwise sequentially very much related. We reverently regard *B4* as the foundation of the CSI at a time when it was quite difficult to impart the much needed SDT to support the basic axioms. Sixteen years later, the foundation has been fortified further in Robertson et al. (1988) with annexation of some research published after 1970. In the past 16 years, there has been a phenomenal growth of research literature in CSI, not only covering the groundwork reported in the earlier two volumes, but much more so, in novel areas that have been annexed to statistical inference more intensively during the past two decades. As such, with enthusiasm, we have attempted to bring out an updated coverage of the entire integrated area of CSI, strangely after a lapse of another 16 years!

Given our motivation to integrate the CSI across the parametrics to all the way to beyond parametrics, our task has not been very easy or routine. Based on our contemplated intermediate level of presentation (for the convenience of upper-undergraduate- and lower-graduate-level students as well as researchers in applied areas), we have therefore recast the setup as follows.

Treatment	Y	$Depth$	Age	Treatment	Y	$Depth$	Age
1	4.8	17	27		2	5.2	61
1	13.2	6	41		2	6.6	34
1	5.8	30	27		2	2.7	41
1	4.6	40	33		2	5.4	29
1	6	30.5	34		2	8.2	21
1	2.9	48	21		2	16.4	15
1	5.2	50	43		3	6.7	41
1	5.6	20	26		3	6.7	32
1	3.9	40	36		3	6.7	32
1	5.6	24.4	23		3	7.9	15
2	6	30	25		3	6	32
2	9.6	25	58		3	19.4	0
2	15.5	25	55		3	19	15.7
2	8.7	56	43		3	2.8	23
2	7.9	51	36		3	6.6	26
					3	10.4	21

Table 1.14 Patient response to treatments for reversal of anesthesia

In Chapters 2, 3, and 4, we consider the so called exact CSI in the parametric setting with due emphasis on the role of the classical *likelihood principle* (LP) in CSI. Chapter 2 provides an introduction to the traditional order-restricted tests in the one-way classification model. An objective of this chapter is to illustrate that, for comparing several populations in the usual one-way classification, it is easy to use a natural generalization of the traditional F -test to incorporate order restrictions on the population means. It is hoped that these methods will be incorporated into standard statistical software so that they could be used easily. Chapter 3 provides a rigorous theoretical development of the likelihood approach for testing hypotheses about the mean of a multivariate normal and the regression parameter in the normal theory linear model. These results lay the foundation on a firm footing for CSI in general statistical models that may be parametric, semiparametric, or nonparametric. In Chapter 4, we broaden the multivariate normal framework of the previous chapter and consider general parametric models which includes the generalized linear models and nonlinear time-series models. In this setting, very little is known about the exact distribution of statistics (estimators/test statistics); this is to be expected because even when there are no order restrictions on parameters, most of the available results are asymptotic in nature. Chapter 4 is devoted to developing the essentials for large sample statistical inference based mainly on the likelihood when there are constraints on parameters. Chapters 3 and 4 form the core of constrained statistical inference and open the door to a vast literature of current research interest. The remaining chapters assume familiarity with the results in these two chapters.

It would become more evident as we make progress in the subsequent chapters that in CSI, LP is often encountered with some challenges, particularly, in not so simple models, and there could be some alternative approaches that may have greater flexibility in this respect. Chapter 5 brings out this issue with the exploration of the Roy (1953) *union-intersection principle* (UIP), which has been more conveniently adapted in beyond parametrics scenarios, too. Chapters 5, 6, and 7 highlight these developments with genuine nonparametrics as well as some other semiparametric models. Functional parameters as well as more complex statistical functionals arising in survival and reliability analysis and multiple comparisons have been focused in this context. Chapter 8 deals with the essentials of SDT and attempts to tie-up some loose ends in theory that were relegated from earlier chapters to this one. In particular, shrinkage estimation and Bayes tests in CSI are unified in this coverage. Some miscellaneous CSI problems, we found hard to integrate with the presentations in earlier chapters, are relegated to the concluding one. In this way, our contention has been to present the theory and methodology at a fairly consistent and uniform (intermediate) level, albeit, for technical reasons, in latter chapters, we might have some occasional detours into higher level of abstractions, to a small extent. We therefore strive to provide more up-to-date and thorough coverage of the bibliography, particularly the post-1987 era, so that serious and advanced readers could find their desired way out through the long bibliography appended.

2

Comparison of Population Means and Isotonic Regression

The need to compare several treatments/groups/populations with respect to their means or medians or some other location parameters arises frequently in many areas of applications. Typically, we have k populations with means μ_1, \dots, μ_k and independent samples from the k populations. One of the statistical inference problems that often arises in this context is test of

$$H_0 : \mu_1 = \dots = \mu_k \quad \text{vs} \quad H_2 : \mu_1, \dots, \mu_k \text{ are not all equal.}$$

Standard methods for this include the F -test in analysis of variance (ANOVA) and nonparametric methods based on ranks. However, when several treatments are being compared, one is usually prepared to assume that certain treatments are not worse than another, for example, $\mu_1 \leq \mu_2 \leq \mu_3$. In this setup, it is of interest to incorporate the prior information $\mu_1 \leq \mu_2 \leq \mu_3$ for the purposes of improving the statistical analysis. In other words, one would like to target the analysis to detect a departure away from $\mu_1 = \mu_2 = \mu_3$ but in the direction of $\mu_1 \leq \mu_2 \leq \mu_3$. The usual unrestricted F -test in ANOVA is not particularly suitable for this because it does not incorporate prior knowledge such as $\mu_1 \leq \mu_2 \leq \mu_3$ and hence it is not specifically targeted to detect departures in any particular direction of (μ_1, μ_2, μ_3) . This chapter provides an introduction to an approach for dealing with this type of statistical inference problems.

Most of the results used in this chapter are also discussed in the following chapters under more general settings. Therefore, advanced readers may start with the next chapter and refer back to this chapter when necessary, or read up to Section 2.3 and then turn to the next chapter.

In the statistics literature, a constraint/prior-knowledge of the form $\mu_i \leq \mu_j$ is called an *order restriction* or an *inequality constraint* on the parameters. When an

order restriction appears in a hypothesis, as in the foregoing example, it is called an *ordered hypothesis* or more generally a *constrained hypothesis*.

As an example, consider Example 1.2.1 (page 2) on the comparison of two exercise programs with a control. In this example, the response variable is the age at which a child starts to walk. Let μ_i denote the population mean of the response variable for Treatment i , $i = 1, 2, 3$. An inference problem that is of interest in this context is test of

$$H_0 : \mu_1 = \mu_2 = \mu_3 \quad \text{vs} \quad H_2 : \mu_1, \mu_2, \text{ and } \mu_3 \text{ are not all equal.}$$

For illustrative purposes let us suppose that the exercise program was expected not to have detrimental effects, and that the researcher was interested to ensure that the prior knowledge, $\mu_1 \leq \mu_2 \leq \mu_3$, is incorporated in the statistical analysis when testing H_0 against H_2 . This chapter introduces what we call the \bar{F} -test, a natural generalization of the usual F -test, to incorporate prior knowledge in the form of order restrictions. A classical test known as the \bar{E}^2 -test is also introduced. For most practical purposes this chapter is to illustrate that it is fairly easy to apply these tests to incorporate prior information such as $\mu_1 \leq \mu_2 \leq \mu_3$. To this end, only the essential results to implement the test are presented here; a detailed study of the distribution theory in a general setting will be presented in the next chapter.

Section 2.1 provides an introduction to test of $H_0 : \mu_1 = \dots = \mu_k$ against

$$H_1 : \mu_1, \dots, \mu_k \text{ are not all equal and they satisfy several order restrictions.}$$

Example 1.2.2 indicated a scenario where order restrictions appear in the null hypothesis. These type of problems are discussed in Section 2.2. More specifically, this section develops a test of

$$H_1 : \{\mu_1, \dots, \mu_k\} \text{ satisfy several order restrictions} \quad \text{vs} \quad H_2 : \text{not } H_1.$$

It will be shown that the tests developed in Section 2.1 for H_0 vs H_1 and those in Section 2.2 for H_1 vs H_2 can be implemented by using a simulation approach. Further, only a short computer program is required to implement this simulation approach. Thus, an important conclusion of the first two sections is that test of hypothesis involving μ_1, \dots, μ_k when there are order restrictions in the null or the alternative hypothesis can be implemented easily. Such a simulation approach overcomes the computational difficulties encountered in the classical approach of computing the exact critical/p-values and/or bounds for them. Consequently, the statistical results become clearer and easier to use.

Section 2.3 provides an introduction to *isotonic regression*. The foregoing example involving three treatments in which the treatment means are ordered is an example of topics studied under isotonic regression. The origin of the term isotonic regression is that the basic model used is the *regression* model $y_{ij} = \mu_i + \epsilon_{ij}$, and the parameters μ_1, \dots, μ_k are required to have the *same tone* (*iso tone*) as the order imposed on the treatments, for example,

where Treatment 1 \preceq Treatment 2 may be read as “the mean for Treatment 2 is at least as large as that for Treatment 1”; more generally, it could also be read as “Treatment 2 is at least as good as Treatment 1.”

The results in Sections 2.1 and 2.2 for testing $H_0 : \mu_1 = \dots = \mu_k$ against an order restriction are based on the linear model $y_{ij} = \mu_i + \epsilon_{ij}$. Most of these results have natural extensions to the more general linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}$ where \mathbf{E} has mean zero and covariance $\sigma^2 \mathbf{U}$, σ is unknown and \mathbf{U} is known. These will be discussed in the next chapter.

These have also been extended to incorporate robust tests based on M -estimators

and those based on ranks; see Silvapulle (1992b, 1992c, 1985), Puri and Sen (1985), and Geyer (1994). The finite sample distribution theory for normal error distribution will be derived in the next chapter, and the large sample theory will be derived in the following chapter. The results have been extended to bounded influence tests in general parametric models; see Silvapulle (1997c). Section 2.3 provides an introduction to isotonic regression; most of this is related to the results in Section 2.1.

2.1 ORDERED ALTERNATIVE HYPOTHESES

In this section we adopt the same context as that is usually adopted for one-way ANOVA. Suppose that there are k treatments to be compared. Let y_{ij} denote the j th observation for Treatment i , and assume that $y_{ij} \sim N(\mu_j, \sigma^2)$ and that the observations are independent ($j = 1, \dots, n_i$, $i = 1, \dots, k$); see Table 2.1. Later we shall relax the requirement that y_{ij} be normally distributed. In the next subsection, we consider the walking exercise example (see Example 1.2.1) to motivate some of the ideas that underlie the \bar{F} -test and to illustrate its application when the number of order restrictions of the form $\mu_i \leq \mu_j$ is less than four. Then in the following subsection, we consider the more general case with an arbitrary number of treatments and test of

$$H_0 : \mu_1 = \dots = \mu_k \quad \text{vs} \quad H_1 : \mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B \quad (2.1)$$

where

$$B \subset \{(i, j) : i, j = 1, \dots, k\}. \quad (2.2)$$

This particular form of the alternative hypothesis allows only pairwise contrasts; more general contrasts will be discussed in the next chapter.

Table 2.1 Normal theory one-way layout

Treatment	Independent observations	Sample mean	Population distribution
1	y_{11}, \dots, y_{1n_1}	\bar{y}_1	$N(\mu_1, \sigma^2)$
\vdots	\vdots	\vdots	\vdots
k	y_{k1}, \dots, y_{kn_k}	\bar{y}_k	$N(\mu_k, \sigma^2)$

2.1.1 Test of $H_0 : \mu_1 = \mu_2 = \mu_3$ Against an Order Restriction

In this subsection, we consider the special case when there are only three treatments and the null hypothesis is $H_0 : \mu_1 = \mu_2 = \mu_3$; this simple case is convenient to introduce the basic ideas. In this special case, there are three possible ordered alternatives that may be of interest in the context of one-way classification; they are

$$(1) \quad \mu_1 \leq \mu_2 \leq \mu_3, \quad (2) \quad \mu_1 \leq \mu_2 \text{ and } \mu_1 < \mu_3, \quad \text{and} \quad (3) \quad \mu_1 < \mu_2.$$

For illustrative purposes, we shall consider the first; the results for the second and third are similar and are also stated in this subsection.

Consider Example 1.2.1, and assume that the setting in Table 2.1 holds. If there is no prior knowledge in the form of an order among the treatment means, then the null and alternative hypotheses would be

$$H_0 : \mu_1 = \mu_2 = \mu_3 \quad \text{and} \quad H_2 : \mu_1, \mu_2, \text{ and } \mu_3 \text{ are not all equal,}$$

respectively; the notation H_2 , rather than H_1 , is used here for the alternative hypothesis to suggest that it is “two-sided” or unrestricted. For testing H_0 against H_2 , the classical F -test is usually preferred.

The purpose of the study was to test the claim that walking exercises are associated with a reduction in the mean age at which children start to walk. Let us suppose that the alternative hypothesis is

$$H_1 : \mu_1 \leq \mu_2 \leq \mu_3, \text{ and } \{\mu_1, \mu_2, \mu_3\} \text{ are not all equal.}$$

As indicated in the previous chapter, we shall abbreviate the statement of this hypothesis to

$$H_1 : \mu_1 \leq \mu_2 \leq \mu_3.$$

More generally, we shall adopt the following convention throughout this book:

Notation: When we refer to test of “ H_0 against H_1 ,” it should be read as “ H_0 against $H_1 \setminus H_0$ ”; in the literature, $H_1 \setminus H_0$ is also written as $H_1 - H_0$.

In view of this notation, test of $\mu = 0$ vs $\mu \geq 0$ and test of $\mu = 0$ vs $\mu > 0$ are exactly the same testing problems, where μ is a scalar.

To test H_0 against H_1 , it is perfectly valid to apply the usual F -test for H_0 against H_2 on $(2, \nu)$ degrees of freedom, where ν is the error degrees of freedom. The validity is justified by the fact that both testing problems have the same null hypothesis. However, one would expect that the standard F -test for H_0 against H_2 is unlikely to have good power properties for testing H_0 against H_1 because it does not make use of the additional information $\mu_1 \leq \mu_2 \leq \mu_3$ and hence is not specifically targeted to detect departures in the direction of H_1 .

Recall that the standard F -statistic is defined as

$$F = \{RSS(H_0) - RSS(H_2)\}(k-1)^{-1}/S^2 \quad (2.3)$$

where S^2 is the error mean square, k is the number of treatments, $RSS(H)$ is the abbreviation for Residual Sum of Squares under the hypothesis H ,

$$RSS(H_0) = \inf_{H_0} \sum_{i,j} (y_{ij} - \mu_i)^2 = \sum_{i,j} (y_{ij} - \bar{y})^2$$

$$RSS(H_2) = \inf_{H_2} \sum_{i,j} (y_{ij} - \mu_i)^2 = \sum_{i,j} (y_{ij} - \bar{y})^2$$

and \bar{y} is the grand mean. The term $\{RSS(H_0) - RSS(H_2)\}$ in the numerator of F (see (2.3)) is a measure of the discrepancy between the null H_0 and the alternative H_2 . The denominator, S^2 , acts as a scaling factor so that the null distribution of the test statistic does not depend on the unknown scale parameter σ for the error; S^2 does not contribute to quantify the discrepancy between H_0 and H_2 .

Remark: Strictly speaking, $RSS(H_2)$ should have been written as $RSS(H_0 \cup H_2)$. Since H_0 is on the boundary of H_2 , this does not affect the value of $RSS(H_2)$.

These observations suggest that to test H_0 against H_1 , a reasonable test statistic is obtained by modifying the foregoing F -statistic as follows:

$$\bar{F} = \{RSS(H_0) - RSS(H_1)\}/S^2 \quad (2.4)$$

where

$$RSS(H_1) = \min_{H_1} \sum_{i,j} (y_{ij} - \mu_i)^2 = \sum_{i,j} (y_{ij} - \tilde{\mu}_i)^2$$

and $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3)$ is the point at which the sum of squares $\sum_{i,j} (y_{ij} - \mu_i)^2$ is minimized subject to the constraint in H_1 . Thus, $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3)$ is an estimate of (μ_1, μ_2, μ_3) under H_1 and $RSS(H_1)$ is the corresponding sum of squares of the residuals; in fact, since the errors are iid as $N(0, \sigma^2)$, it can be shown that $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3)$ is the mle of (μ_1, μ_2, μ_3) under H_1 [see Problem 2.3]. Thus, the numerator $\{RSS(H_0) - RSS(H_1)\}$ of \bar{F} is a measure of the discrepancy between H_0 and the restricted alternative H_1 , just as the numerator of F is a measure of the discrepancy between H_0 and the unrestricted alternative H_2 . The constant $(k-1)^{-1}$ in the numerator of the usual F statistic was not included in the definition of \bar{F} , but it could have been included without affecting the essential nature of any results.

To implement the foregoing procedure, we need to compute the restricted estimator $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3)$ under $H_1 : \mu_1 \leq \mu_2 \leq \mu_3$. Since

$$\sum_i \sum_j (y_{ij} - \mu_i)^2 = \sum_i \sum_j (y_{ij} - \bar{y}_i)^2 + \sum_i (\bar{y}_i - \mu_i)^2 n_i \quad (2.5)$$

it suffices to minimize $\sum_i (\bar{y}_i - \mu_i)^2 n_i$ subject to $H_1 : \mu_1 \leq \mu_2 \leq \mu_3$ for the purposes of computing $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3)$.

If the alternative hypothesis is any other order restriction, say H_1^* , then the foregoing discussions hold with trivial modifications; the main modification is that $RSS(H_1^*)$ is $\sum_i \sum_j (y_{ij} - \tilde{\mu}_i)^2$ where $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3)$ is the point at which $\sum_i \sum_j (y_{ij} - \mu_i)^2$ reaches its minimum subject to the constraints in H_1^* . Now we have the following result; this is a consequence of more general results in the next chapter (see (3.25)).

Proposition 2.1.1 Let the setting be as in Table 2.1 with $k = 3$ and let the null hypothesis be $H_0 : \mu_1 = \mu_2 = \mu_3$. Then, for testing against any order restriction of the form (2.1) on page 27, the null distribution of \bar{F} is

$$pr(\bar{F} \leq c | H_0) = w_0 + w_1 pr(F_{1,\nu} \leq c) + w_2 pr(F_{2,\nu} \leq c/2), \quad (c > 0) \quad (2.6)$$

Table 2.2 Formulae for determining $\{w_1, w_2, w_3\}$ in (2.6)

H_1	ρ in (2.7)
$\mu_1 \leqslant \mu_2 \leqslant \mu_3$,	$\rho = -[n_1 n_3 / \{(n_1 + n_3)(n_2 + n_3)\}]^{1/2}$
$\mu_1 \leqslant \mu_2$ and $\mu_1 \leqslant \mu_3$	$\rho = [n_2 n_3 / \{(n_1 + n_2)(n_1 + n_3)\}]^{1/2}$
$\mu_1 \leqslant \mu_2$	1.0 [i.e. $(w_0, w_1, w_2) = (0, 0.5, 0.5)$]

where

$$w_1 = 0.5, \quad w_2 = (0.5 - q), \quad w_0 + w_1 + w_2 = 1, \quad q = (2\pi)^{-1} \cos^{-1}(\rho) \quad (2.7)$$

for some ρ . The formulas for determining ρ are given in Table 2.2. ■

We chose the notation \bar{F} for the statistic in (2.4) because it is related to the unrestricted F -ratio and its null distribution is a weighted average of probabilities associated with F distributions. Some authors have used \bar{F} for a different statistic. It follows from (2.6) that the p -value for the \bar{F} -test is given by

$$p\text{-value} = w_1 \text{pr}(F_{1,\nu} \geqslant \bar{f}_{obs}) + w_2 \text{pr}(2F_{2,\nu} \geqslant \bar{f}_{obs}) \quad (2.8)$$

where \bar{f}_{obs} is the sample value of \bar{F} .**Example 2.1.1**

Now let us consider the exercise program example discussed at the beginning of this section. Suppose that we are interested to test

$$H_0 : \mu_1 = \mu_2 = \mu_3 \quad \text{vs} \quad H_1 : \mu_1 \leqslant \mu_2 \leqslant \mu_3.$$

The unrestricted estimate of (μ_1, μ_2, μ_3) , obtained by minimizing $\sum (\bar{y}_i - \mu_i)^2 n_i$ is $(\bar{y}_1, \bar{y}_2, \bar{y}_3) = (10.125, 11.375, 12.35)$. Since this estimate satisfies the constraints in H_1 , it follows that the estimate of (μ_1, μ_2, μ_3) subject to the constraint in H_1 is also equal to the unconstrained estimate. Therefore, $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3) = (\bar{y}_1, \bar{y}_2, \bar{y}_3) = (10.125, 11.375, 12.35)$ and, by simple substitution, we have that $\bar{F} = 5.978$. Now, let us assume that the setting in Table 2.1 holds. Then, it follows from (2.8) that for testing $H_0 : \mu_1 = \mu_2 = \mu_3$ vs $H_1 : \mu_1 \leqslant \mu_2 \leqslant \mu_3$, we have that

$$p\text{-value} = 0.5 \text{pr}(F_{1,14} \geqslant 5.989) + \{(0.5 - 0.329)\} \text{pr}(F_{2,14} \geqslant 2.989) = 0.028.$$

By contrast, for the unrestricted F -statistic for testing

$$H_0 : \mu_1 = \mu_2 = \mu_3 \quad \text{vs} \quad H_2 : \{\mu_1, \mu_2, \mu_3\} \text{ are not all equal}$$

$p\text{-value} = \text{pr}(F_{2,14} \geqslant 2.989) = 0.083$. Note that the p -value for the F -test is larger than that for the \bar{F} -test. If the alternative hypothesis is $\mu_1 \leqslant \mu_2 \leqslant \mu_3$ and the sample means satisfy the corresponding inequality $\bar{y}_1 \leqslant \bar{y}_2 \leqslant \bar{y}_3$, then it can be shown using (2.6) that the p -value for \bar{F} would be smaller than that for F . Of course, it is

reasonable to expect that the \bar{F} -test, which is specifically targeted for testing against $\mu_1 \leqslant \mu_2 \leqslant \mu_3$, would provide stronger evidence to reject H_0 than the unrestricted F -test when the sample means satisfy $\bar{y}_1 \leqslant \bar{y}_2 \leqslant \bar{y}_3$; it is reassuring to note that this, in fact, is the case. ■

In this example, the calculations were simple because there were only two order restrictions and the sample means satisfied the constraints that correspond to those in the alternative hypothesis. If the alternative hypothesis has up to three order restrictions, then closed-form expressions similar to (2.6)-(2.7) are available for computing the p -value. If the number of order restrictions is four or more then the null distribution of \bar{F} is a weighted sum similar to (2.6), but in general it is quite inconvenient to use it for computing the p -value exactly. A simulation approach offers a simple and practically feasible method of computing the p -value sufficiently precisely irrespective of the number of order restrictions, for any error distribution. These are discussed in the next subsection. The derivations of the null distribution of \bar{F} and that of the likelihood ratio test statistic will be provided in the next chapter.

2.1.2 Test of $H_0 : \mu_1 = \dots = \mu_k$ Against an Order Restriction

Suppose that there are k treatments. Let the independent observations $\{y_{ij}\}$ be as in Table 2.3, where F is a cumulative distribution function. Clearly, μ_i is a location parameter for treatment i , but it does not need to be the mean ($i = 1, \dots, k$); for example, it could be the median. Let the null and alternative hypotheses be

$$H_0 : \mu_1 = \dots = \mu_k \quad \text{and} \quad H_1 : \mu_i - \mu_j \geqslant 0 \text{ for } (i, j) \in B, \quad (2.9)$$

respectively, for some $B \subset \{(i, j) : i, j = 1, \dots, k\}$. Let $\boldsymbol{\mu} = (\mu_1, \dots, \mu_k)^T$ and H denote a hypothesis concerning $\boldsymbol{\mu}$, for example, H can be H_0 or H_1 . Now, we define the residual sum of squares under H as

$$RSS(H) = \inf_{\boldsymbol{\mu} \in H} \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \mu_i)^2. \quad (2.10)$$

It is easily verified that this definition leads to the familiar expressions

$$RSS(H_0) = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 \quad \text{and} \quad RSS(H_2) = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2 \quad (2.11)$$

for the residual sum of squares under $H_0 : \mu_1 = \dots = \mu_k$ and under $H_2 : \mu_1, \dots, \mu_k$ are not restricted, respectively. In this sense, (2.10) is a natural extension of the two

Table 2.3 Comparison of k means

Treatment	Independent observations	Sample mean	Population distribution (cdf)
1	y_{11}, \dots, y_{1n_1}	\bar{y}_1	$F\{(t - \mu_1)/\sigma^2\}$
\vdots	\vdots	\vdots	\vdots
k	y_{k1}, \dots, y_{kn_k}	\bar{y}_k	$F\{(t - \mu_k)/\sigma^2\}$

familiar expressions in (2.11) for residual sums of squares. Let

$$\tilde{\mu} = (\tilde{\mu}_1, \dots, \tilde{\mu}_k)^T = \arg \min_{\mu \in H_1} \sum \sum (y_{ij} - \mu_i)^2.$$

Thus, by (2.10), we have that

$$RSS(H_1) = \sum \sum (y_{ij} - \tilde{\mu}_i)^2 \quad (2.12)$$

for the residual sum of squares under H_1 . If the error distribution is normal then $(\tilde{\mu}_1, \dots, \tilde{\mu}_k)$ is also the *mle* of (μ_1, \dots, μ_k) under H_1 . By motivations similar to those leading to the \bar{F} in (2.4), let us define

$$\bar{F} = \{RSS(H_0) - RSS(H_1)\}/S^2, \quad (2.13)$$

where $S^2 = \nu^{-1} \sum \sum (y_{ij} - \bar{y}_i)^2$ is the mean square for error and $\nu = (n_1 + \dots + n_k - k)$ is the error degrees of freedom. As in the previous subsection, the \bar{F} in (2.13) can be considered as a modification of the usual F -statistic to incorporate the order restrictions in H_1 . The exact finite sample null distribution of \bar{F} for normal errors will be discussed at the end of this subsection. The p -value of \bar{F} can be estimated by simulation. To this end, the following result is important.

Theorem 2.1.2 Consider the setting in Table 2.3 and the testing problem (2.9). Suppose that the null hypothesis, $H_0 : \mu_1 = \dots = \mu_k$, holds; let μ denote the common value of μ_1, \dots, μ_k . Then the exact finite sample null distribution of \bar{F} depends on the functional form of F but not on (μ, σ) . Further, the asymptotic null distribution of \bar{F} does not depend on (F, μ, σ) , where the limit is taken as $(n_1 + \dots + n_k) \rightarrow \infty$ such that $0 < \lim_{n \rightarrow \infty} n_i/(n_1 + \dots + n_k) < 1$ for $i = 1, \dots, k$.

Proof: Let $y_{ij}^* = (y_{ij} - \mu)/\sigma$ for every (i, j) . Let $RSS^*(H_0)$, $RSS^*(H_1)$, S^* and \bar{F}^* denote the values corresponding to $\{y_{ij}^*\}$. Then $RSS^*(H_0) = \sigma^{-2} RSS(H_0)$, $RSS^*(H_1) = \sigma^{-2} RSS(H_1)$, $(S^*)^2 = \sigma^{-2} S^2$, and hence

$$\bar{F} = \{RSS(H_0) - RSS(H_1)\}/S^2 = \{RSS^*(H_0) - RSS^*(H_1)\}/(S^*)^2 = \bar{F}^*.$$

Under H_0 , the distribution of y_{ij}^* is $F(t)$, which does not depend on (μ, σ) . Now, since \bar{F}^* is a function of $\{y_{ij}^*\}$ only and $\bar{F} = \bar{F}^*$, it follows that the distribution of \bar{F} does not depend on (μ, σ) . The proof of the second part will be discussed in the next chapter. ■

This result suggests the following simulation method for computing the p -value.

Computation of the exact p -value for the \bar{F} -test:

Suppose that H_0 and H_1 are as in (2.9), \bar{F} is as in (2.13), the setting is as in Table 2.3 and the functional form F of the error distribution is known. Then the following steps are adequate to compute the p -value for \bar{F} .

1. Generate independent observations $\{y_{ij} : j = 1, \dots, n_i, i = 1, \dots, k\}$ from $F\{(t - \mu_0)/\sigma_0\}$ where (μ_0, σ_0) can have any arbitrary values, but must be held fixed for different (i, j) .

2. Compute the \bar{F} statistic in (2.13) with $RSS(H_0)$ and $RSS(H_1)$ as in (2.11) and (2.12), respectively.
3. Repeat the previous two steps N times (say $N = 10000$), and estimate the p -value by M/N where M is the number of times the \bar{F} statistic in the second step exceeded its sample value. ■

Note that in the first step of the simulation, the observations may be generated from a distribution with any value for the common location and scale parameters because, in view of Theorem 2.1.2, the null distribution of \bar{F} does not depend on the common value of μ_1, \dots, μ_k or σ . For example, if the errors are assumed to be $N(0, \sigma^2)$ where σ is unknown, then it suffices to generate the observations from $N(0, 1)$ in the first step of the simulation for computing the p -value. Similarly, if the error distribution is a logistic with unknown variance, then it suffices to generate the observations from a logistic distribution with any values for its mean and variance.

The simulation approach can also be used to implement a bootstrap-type procedure in which the observations may be generated from the empirical distribution of the within treatment residuals. To improve robustness of validity against violation of the assumption of normal error distribution, we may estimate the p -value corresponding to several forms of error distribution and then choose their maximum as the p -value. The second part of Theorem 2.1.2 says that if the error distribution is $F(t/\sigma)$ for some (F, σ) , which may be unknown, then the asymptotic distribution of \bar{F} does not depend on (F, σ) . Therefore, if n_i is large for $i = 1, \dots, k$ then the foregoing simulation for normal error would estimate the asymptotic p -value, irrespective of whether or not the true error distribution is normal. Similarly, the simulation method with any error distribution would also estimate the asymptotic p -value if n_i is large for $i = 1, \dots, k$, irrespective of the true error distribution. It is reasonable to conjecture that closer the $F\{(\cdot - \mu)/\sigma\}$ to the true error distribution, for some (μ, σ) , the better the approximation.

Suppose that the precise form of F is unknown but it is known that F is a member of the class of distributions \mathcal{F} for some \mathcal{F} . Let p_F denote the p -value corresponding to F in Table 2.3. Then we have

$$p\text{-value} = \sup_{F \in \mathcal{F}} p_F \quad (2.14)$$

As an example, suppose that the error distribution is known to be normal, logistic, or a t -distribution with four degrees of freedom (T_4). Then, the foregoing p -value in (2.14) for \bar{F} can be computed by applying the foregoing simulation procedure to compute the p -values corresponding to normal, logistic, and T_4 and then taking their maximum as the p -value in (2.14). In most practical applications, we would not know \mathcal{F} , although it would not be difficult to specify a small class of distributional forms that would be well spread out in \mathcal{F} . In this case, a reasonable procedure is to compute the p -values corresponding to a range of choices of F (such as normal, logistic, T_r , and χ_r^2) that are well spread out in \mathcal{F} and take their maximum as a reasonable approximation to the p -value. ■

Example 2.12

The experiment in Example 1.2.1 also included a fourth treatment; it was not considered in the previous worked example for simplicity. The group of children receiving the fourth treatment were checked weakly for progress but they did not receive any special exercises.

Table 2.4 Effect of exercise on the age at which a child first walks

Treatment	Age	n	\bar{y}
Group 1:	9.00	9.50	9.75
Group 2:	11.00	10.00	10.00
Group 3:	13.25	11.50	12.00
Group 4:	11.50	12.00	9.00

Reprinted from: *Science*, Volume 176, P. R. Zelazo, N. A. Zelazo, and S. Kolb, Walking in the newborn, Pages 314–315, Copyright (1972), with permission from AAAS.

Since Groups 1 and 2 were also checked weekly, it may be argued that Group 4 is more appropriate as a control than Group 3; we shall not pursue this issue any further here. The data for all four groups are given in Table 2.4.

Suppose that we wish to test $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$ against an alternative that captures our *a priori* information. There is no unique way of formulating the alternative hypothesis. For illustrative purposes, let us suppose that we wish to test

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 \text{ against } H_1 : \mu_1 \leq \mu_3, \mu_2 \leq \mu_3, \mu_1 \leq \mu_4, \mu_2 \leq \mu_4.$$

This says that Treatments 1 and 2 are at least as good as Treatments 3 and 4, but no ordering is suggested between Treatments 1 and 2, or between Treatments 3 and 4. In this case, there is no convenient formula similar to (2.8) for computing the p-value for \bar{F} ; this is a common feature of most of the problems that involve test against inequality constraints. However, in view of Theorem 2.1.2, the simulation method on page 32 is still applicable.

Since $(\bar{y}_1, \bar{y}_2, \bar{y}_3, \bar{y}_4) = (10.1, 11.4, 12.4, 11.7)$, the sample means satisfy the restrictions in H_1 . Therefore, the constrained and unconstrained estimators of $(\mu_1, \mu_2, \mu_3, \mu_4)$ are the same; in particular, $(\tilde{\mu}_1, \tilde{\mu}_2, \tilde{\mu}_3, \tilde{\mu}_4) = (\bar{y}_1, \bar{y}_2, \bar{y}_3, \bar{y}_4) = (10.1, 11.4, 12.4, 11.7)$. Now substituting directly into the formula for \bar{F} in (2.13), the sample value of \bar{F} is 6.43. The p-values corresponding to a range of error distributions were computed using the simulation approach. The computed values are given in the first row of Table 2.5; the figures in the second row are also p-values for a different statistic, which will be discussed later.

It is also of interest to compute the p-value by resampling; to this end, we set the error distribution equal to the empirical distribution of the residuals about the treatment means. This can be implemented by resampling with replacement from the set of within treatment residuals. The computed p-value is also given in the last column of Table 2.5. This example illustrates that simulation is a very convenient

Table 2.5 The p-values for the \bar{F} - and \bar{E}^2 -tests for different error distributions

Test	$N(0, \sigma^2)$	T_4	T_{10}	Error Distribution	χ^2_1	χ^2_2	χ^2_4	χ^2_7	RSS
\bar{F}	0.052	0.051	0.058	0.050	0.048	0.049	0.051	0.047	0.048
\bar{E}^2	0.046	0.047	0.053	0.046	0.044	0.044	0.047	0.047	0.043

The symbol RSS refers to resampling from the unrestricted within treatment residuals. To compute the p-value corresponding to $N(0, \sigma^2)$, we used $\sigma = 1$.

way of implementing tests against any order restriction in the one-way classification model when the errors are iid even if the common error distribution is not normal. It is well known that the critical values from the F-tables are reasonably robust even if the error distribution is not normal for the unrestricted F-test of $H_0 : \mu_1 = \dots = \mu_k$ against $H_2 : \mu_1, \dots, \mu_k$ are unequal. It will be argued later that a similar comment is likely to be applicable for the \bar{F} -test as well. Since the p-values in Table 2.5 are close for different error distributions, the results of this example are consistent with this conjecture.

The value of the standard unconstrained F-statistic for testing $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$ against $H_2 : \{\mu_1, \mu_2, \mu_3, \mu_4\}$ are not all equal is 2.14 and its p-value is 0.129 based on F-statistic $\sim F_{3,19}$. The numerical results for this example again highlight the fact that if the sample means satisfy the constraints that correspond to those in the alternative hypothesis, then the estimate of μ under H_1 and H_2 are the same and the p-value for the constrained test would be smaller than that for the unconstrained F-test. ■

2.1.2.1 Exact finite sample results when the error distribution is $N(0, \sigma^2)$

The results for normal theory likelihood ratio test of H_0 against H_1 are very similar to those for the \bar{F} -test. Assume that the errors are independent and distributed as $N(0, \sigma^2)$. Let us define

$$\bar{E}^2 = \{RSS(H_0) - RSS(H_1)\} / RSS(H_0).$$

Then the \bar{E}^2 -test rejects H_0 for large values of \bar{E}^2 . If the error distribution is normal then, it may be verified that

$$\bar{E}^2 = \{1 - \exp(-LRT/n)\} \quad (2.15)$$

where LRT denotes the likelihood ratio statistic ($= -2 \log \Lambda$) for testing H_0 against H_1 . It follows from (2.15) that LRT is an increasing function of \bar{E}^2 , and, therefore, two tests are equivalent.

It is worth noting that \bar{F} and \bar{E}^2 have the same numerator, namely $\{RSS(H_0) - RSS(H_1)\}$, which measures the discrepancy between H_0 and H_1 . The denominators of these statistics act as scaling factors so that the null distribution of the test statistics

do not depend on the unknown scale parameter of the error distribution; they do not make any contribution to quantifying the discrepancy between H_0 and H_1 . Now, Theorem 2.1.2 holds with \bar{F} replaced by \bar{E}^2 as well; simply replace F by \bar{E}^2 in the proof. If the error distribution is not normal then \bar{E}^2 is not necessarily a monotonic function of the likelihood ratio statistic, but it is still a suitable test statistic, and the three-step simulation for computing the p -value for \bar{F} (see page 32) can also be applied for \bar{E}^2 , with the latter replacing the former in the second step of the simulation.

For the Walking Exercise example discussed earlier in this section, the sample value of \bar{E}^2 is 0.253. The corresponding p -values for several error distributions are given in Table 2.5; they range from 0.043 for the resampling method to 0.053 for T_{10} . Note that the results for \bar{F} and \bar{E}^2 are close; we conjecture that the difference between \bar{F} -test and \bar{E}^2 -test in terms of p -value for a given set of data and in terms of power, would be small in most practical situations.

2.1.2.2 Computing $\min_{\mu \in H_1} \sum \sum (y_{ij} - \mu_i)^2$, \bar{F} , and \bar{E}^2 .

The implementation of the \bar{F} - and \bar{E}^2 -tests requires the computation of $RSS(H_1)$ that is equal to

$$\min_{\mu \in H_1} \sum \sum (y_{ij} - \mu_i)^2 \quad \text{where } H_1 : \mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B.$$

Let

$$q(\mu) = (\bar{y} - \mu)^T W (\bar{y} - \mu) \quad (2.16)$$

where $\mu = (\mu_1, \dots, \mu_k)^T$, $\bar{y} = (\bar{y}_1, \dots, \bar{y}_k)^T$ and $W = \text{diag}\{n_1, \dots, n_k\}$, and let A be a matrix such that

$$\{\mu : \mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B\} = \{\mu : A\mu \geq 0\}. \quad (2.17)$$

In this particular case, each row of A is a permutation of the k -vector $(1, -1, 0, \dots, 0)$. Since

$$\sum \sum (y_{ij} - \mu_i)^2 = q(\mu) + C(y),$$

where $C(y) = \sum \sum (y_{ij} - \bar{y}_i)^2$, which does not depend on μ , we have that

$$\bar{F} = \left\{ \min_{H_1} q(\mu) - \min_{H_1} q(\mu) \right\} / S^2,$$

$$\text{and } \bar{E}^2 = \left\{ \min_{H_0} q(\mu) - \min_{H_1} q(\mu) \right\} / \left\{ \min_{H_0} q(\mu) + C(y) \right\}.$$

We can use the following general method for computing the minimum of $q(\mu)$ subject to $\mu \in H_i$, ($i = 1, 2$). Let $x \in \mathbb{R}^k$ be given, B be a given $k \times k$ positive definite matrix and

$$g(\mu) = (x - \mu)^T B (x - \mu), \quad \mu \in \mathbb{R}^k. \quad (2.18)$$

Suppose that we wish to solve

$$\min g(\mu) \quad \text{subject to } A_1\mu \geq 0 \text{ and } A_2\mu = 0 \quad (2.19)$$

for some matrices A_1 and A_2 that do not depend on μ . This constrained minimization problem in which the objective function is a quadratic in μ and the constraints are linear equality and inequality constraints in μ is called a *quadratic program*. There are efficient computer algorithms and software for this optimization problem. For example, we found that the subroutine QPROG in IMSL worked well; similar procedures in packages such as MATLAB, GAUSS, and SPLUS are likely to work just as well.

The foregoing quadratic programming problem is sometimes expressed in the following slightly different, but equivalent, form. Note that

$$g(\mu) = 2f(\mu) + \text{constant}, \quad \text{where } f(\mu) = a^T \mu + (1/2)\mu^T B\mu$$

and $a = -B^T x$. Therefore, the minimization of $g(\mu)$ subject to some constraints on μ is equivalent to the minimization of $f(\mu)$ subject to the same constraints on μ . Therefore, $\hat{\mu}$, the solution to (2.19), is also the solution to

$$\min\{a^T \mu + (1/2)\mu^T B\mu\} \text{ subject to } A_1\mu \geq 0 \text{ and } A_2\mu = 0. \quad (2.20)$$

2.1.3 General Remarks

It will be shown in the next chapter that, if the errors are independent and distributed as $N(0, \sigma^2)$, then the p -value for \bar{F} is

$$\sum_{i=0}^k w_i(H_0, H_1) \text{pr}(iF_{i,\nu} \geq \bar{f}_{obs}), \quad (2.21)$$

where $\{w_i(H_0, H_1)\}$ are some nonnegative weights, which depend on the null hypothesis H_0 and the alternative hypothesis H_1 , and \bar{f}_{obs} is the sample value of \bar{F} . The quantities $\{w_i(H_0, H_1)\}$ are known as *chi-bar-square weights* and also as *level probabilities*; the reasons for these terms will become clear in the next chapter. They appear in the null distribution of several test statistics when there are inequality constraints on parameters.

Often, when there are order restrictions, the null distributions of various test statistics, including \bar{F} and \bar{E}^2 , take the form $w_0 \text{pr}(X_0 \leq c) + \dots + w_k \text{pr}(X_k \leq c)$ where $w_i = w_i(H_0, H_1)$ and X_0, \dots, X_k are some random variables. In this expression, computation of $\text{pr}(X_i \leq c)$ is usually not difficult ($i = 0, \dots, k$). However, simple closed-forms are not available for computing w_0, \dots, w_k exactly, in general. Consequently, there have been various attempts to obtain approximations/bounds for these quantities, so that they can be used to obtain approximations/bounds for the p -value. The details of these are not easy and/or tedious. It appears that the complicated nature of these technical details and unavailability of closed-forms for computing the p -value have led to the opinion among some authors that tests of H_0 against an order restriction are difficult to apply. Fortunately, such an opinion is ill-founded, because there is no need to compute the weights $\{w_i\}$ exactly. In view of the fact that the p -value can be computed sufficiently precisely by simulation, the aforementioned

approximations/bounds on $\{w_1, \dots, w_k\}$, and the related bounds on the exact null distributions of \bar{F} and \bar{E}^2 are not that important for implementing tests of hypotheses involving order restrictions on μ_1, \dots, μ_k . Even if w_0, \dots, w_k can be computed exactly, the simulation approach can be used to compute a p -value that would be more robust against violation of distributional assumptions about the error term, although for a one-way classification problem we believe that the p -values corresponding to normal and other error distributions are likely to be close. However, the bounds on $\{w_1, \dots, w_k\}$ may be useful in nonlinear models; this will be discussed in Chapter 4.

2.2 ORDERED NULL HYPOTHESES

Now we consider testing problems in which there are inequality constraints in the null hypothesis. To introduce the general form of the testing problems, let μ_i denote the mean of population i , and let y_{ij} denote the j th observation from population i ($i = 1, \dots, k, j = 1, \dots, n_i$). Assume that the observations are independent; thus the setting is the same as that in Table 2.3 for one-way classification. Suppose that we wish to test

$$H_1 : \mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B \quad \text{vs} \quad H_2 : \text{No restriction on } \mu_1, \dots, \mu_k \quad (2.22)$$

where B is a given subset of $\{(i, j) : i, j = 1, \dots, k\}$. Thus the hypothesis H_1 in (2.22) is the same as that in (2.9), although in the present context it is the null hypothesis; further, let us also note that in view of the convention introduced earlier (see page 28), the alternative parameter space is the set of all values $\{\mu_1, \dots, \mu_k\}$ that do not satisfy the constraints H_1 in (2.22). This section provides an introduction and the essentials for carrying out a test of H_1 vs H_2 . The distribution theory will be presented in the next chapter.

First, let us consider a general testing problem of the form H_a vs H_b , where H_a and H_b are given. By motivations similar to those leading to (2.13), we define

$$\bar{F} = \{RSS(H_a) - RSS(H_a \cup H_b)\}/S^2, \quad (2.23)$$

for testing H_a against H_b , where RSS is defined as in (2.10). This is a natural generalization of the usual unrestricted F -statistic for H_0 vs H_2 to a test of H_a vs H_b ; clearly, the \bar{F} in (2.13) is a special case of (2.23). Thus, for testing H_1 vs H_2 , we have

$$\bar{F} = \{RSS(H_1) - RSS(H_1 \cup H_2)\}/S^2, \quad (2.24)$$

where, by the definition of $RSS(H)$ in (2.10), we have that

$$\begin{aligned} RSS(H_1) &= \min_{\mu} \sum \sum (y_{ij} - \mu_i)^2, \text{ subject to } (\mu_i - \mu_j) \geq 0 \text{ for } (i, j) \in B \\ RSS(H_1 \cup H_2) &= \min_{\mu} \sum \sum (y_{ij} - \mu_i)^2, \text{ subject to no restriction on } \mu. \end{aligned} \quad (2.25)$$

This leads to

$$RSS(H_1) = \sum \sum (y_{ij} - \tilde{\mu}_i)^2, \quad RSS(H_1 \cup H_2) = \sum \sum (y_{ij} - \bar{y}_i)^2. \quad (2.26)$$

Now, with $q(\mu)$ as in (2.16), the \bar{F} in (2.24) for the testing problem in (2.22) takes the form,

$$\bar{F} = \min_{\mu \in H_1} q(\mu)/S^2, \quad (2.27)$$

and therefore the computational methods of subsection 2.1.2.2 are also applicable and adequate.

Let \bar{f}_{obs} denote the sample value of \bar{F} . We need to take some care in defining the p -value and critical value for this \bar{F} -test of H_1 vs H_2 because $\text{pr}(\bar{F} \geq \bar{f}_{obs} \mid \text{the null hypothesis}, H_1)$ is a function of the particular value of (μ_1, \dots, μ_k) in the null parameter space $\{(\mu_1, \dots, \mu_k) : \mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B\}$, for any $\bar{f}_{obs} > 0$. This is in sharp contrast to the testing problem studied in the previous subsection where $\text{pr}(\bar{F} \geq \bar{f}_{obs} \mid \text{the null hypothesis}, H_0)$ does not depend on the particular value of (μ_1, \dots, μ_k) in the null parameter space $\{(\mu_1, \dots, \mu_k) : \mu_1 = \dots = \mu_k\}$. For testing H_1 vs H_2 we define

$$p\text{-value} = \sup_{\mu \in H_1} \text{pr}_{\mu}(\bar{F} \geq \bar{f}_{obs}), \quad (2.28)$$

and a level- α test of H_1 vs H_2 rejects H_1 if this p -value $\leq \alpha$. An interpretation of this standard approach is that if the data are consistent with at least one point in the null parameter space then we do not reject the null hypothesis. Let μ^* denote the value of μ at which the supremum in (2.28) is reached. Then, μ^* is the point in the null parameter space with which the data are most consistent, and the right-hand side of (2.28) is a measure of the degree of this consistency. If this is large, then the data are consistent with the value μ^* in the null parameter space and hence we would not reject the null hypothesis. If, on the other hand, the right-hand side of (2.28) is small then the data are not consistent with any point in the null parameter space and hence we would reject the null hypothesis, H_1 .

At first, it may seem that computing the supremum in (2.28) is a formidable task; fortunately, it is not. It can be shown that for the particular null hypothesis (i.e., H_1) in (2.22), the supremum in (2.28) is reached where $\mu_1 = \dots = \mu_k$; further, the right-hand side of (2.28) does not depend on the common value of μ_1, \dots, μ_k (see the theorem below). Hence, (2.28) can be written as

$$p\text{-value} = \text{pr}(\bar{F} \geq \bar{f}_{obs} \mid \mu_1 = \dots = \mu_k). \quad (2.29)$$

The main result concerning the \bar{F} -test of H_1 vs H_2 is given in the next result.

Theorem 2.2.1 Suppose that the setting is as in Table 2.3, the testing problem is (2.22), and \bar{F} is as in (2.24). Let μ denote the common value of μ_1, \dots, μ_k when they are all equal. Then we have the following:

1. $\sup_{\mu \in H_1} p_{\mu}(\bar{F} \geq c) = \text{pr}(\bar{F} \geq c \mid \mu_1 = \dots = \mu_k), \text{ for any } c > 0.$

2. $\text{pr}(\bar{F} \geq c \mid \mu_1 = \dots = \mu_k = \mu)$, depends on the functional form F of the error distribution but not on (μ, σ) .

3. The asymptotic distribution of \bar{F} at $\mu_1 = \dots = \mu_k (= \mu, \text{say})$, and hence the value of $\lim \text{pr}(\bar{F} \geq c \mid \mu_1 = \dots = \mu_k = \mu)$ as $n \rightarrow \infty$, do not depend on (F, μ, σ) where $c > 0$, $n = (n_1, \dots, n_k)$ and $0 < \lim(n_i/n) < 1$. ■

The proof of the first part will be discussed in the next chapter; see Theorem 3.8.1 and Section 3.9. The proof of the second part follows by arguments very similar to those for the corresponding result in the previous subsection, and hence is omitted; the proof of the last part will be provided in Chapter 4. It follows from the foregoing result that the p -value and critical value for testing H_1 vs H_2 can be estimated very precisely by simulation; the main steps are given below.

Computing the exact p -value for the \bar{F} -test of H_1 vs H_2

Suppose that the setting is as in Table 2.3, the testing problem is (2.22) and the functional form F of the error distribution is known. Thus, $y_{ij} \sim F\{(t - \mu_i)/\sigma\}$ where F is known but σ may be unknown.

1. Generate independent observations $\{y_{ij} : j = 1, \dots, n_i, i = 1, \dots, k\}$ from $F\{(t - \mu_0)/\sigma_0\}$ where we may choose (μ_0, σ_0) to have any convenient value but must be the same for different values of (i, j) .
2. Compute the \bar{F} statistic; to this end we may either use (2.24) with RSS as in (2.26), or use (2.27).
3. Repeat the previous two steps N times (say $N=10000$), and estimate the p -value by M/N where M is the number of times the simulated value of \bar{F} in the second step was not less than its sample value. ■

Note that in the first step of the simulation, it suffices to generate the observations from a distribution with any values for the mean and variance because, in view of Theorem 2.1.2, the null distribution of \bar{F} does not depend on the common location or scale of the error distribution.

Now, suppose that the functional form F of the error distribution is unknown, but is known to be member of the class \mathcal{F} . Let p_F denote the p -value when the error distribution is $F(t/\sigma)$ for some σ as in Table 2.3; by Theorem 2.2.1, p_F does not depend on σ . Therefore, we have

$$p\text{-value} = \sup_{F \in \mathcal{F}} p_F. \quad (2.30)$$

As in the arguments that followed (2.14), a procedure to compute a reasonable approximation to the foregoing p -value is to compute p_F corresponding to a range of choices of F (such as normal, logistic, t_r and χ^2_r) that are well spread out in \mathcal{F} and take their maximum as an approximation to the p -value in (2.30).

The results for the normal theory likelihood ratio test of H_1 against H_2 are closely related to those for the \bar{F} test. For testing H_1 vs H_2 , let us define

$$\bar{E}^2 = \{RSS(H_1) - RSS(H_1 \cup H_2)\}/RSS(H_1).$$

Then, it may be verified that

$$\bar{E}^2 = \{1 - \exp(-LRT/n)\}. \quad (2.31)$$

where LRT is the normal theory likelihood ratio statistic for testing H_1 vs H_2 . Therefore, \bar{E}^2 is a monotonically increasing function of the LRT, and hence the LRT and \bar{E}^2 -tests of H_1 vs H_2 are equivalent. Further, $\bar{E}^2 [= \bar{F}/(1 + \bar{F})]$ is also a monotonic function of \bar{F} for testing H_1 against H_2 . Therefore, for this testing problem \bar{E}^2 -test and the \bar{F} -test are equivalent.

The exact finite sample null distributions of \bar{F} and \bar{E}^2 when the errors are normally distributed will be obtained in the next chapter. It is worth emphasizing that these exact results are for the specific case when the errors are normal; the exact null distributions of \bar{F} and \bar{E}^2 are unknown for other error distributions. Hence the p -value computed using these exact results correspond to normal error only. By contrast the simulation approach to computing the p -value does not require the errors to be normally distributed.

In general, hypothesis testing problems with order restrictions in the null hypothesis arise in contexts that are more complicated than the simple one-way layout studied in this section. Nevertheless, the ideas developed in this section provide a valuable introduction and will be helpful in later sections where these are explored further. The following numerical example illustrates the main results of this subsection on test of hypothesis when there are order restrictions in the null hypothesis.

Example 2.2.1 An example with order restrictions in the null hypothesis

The sample mean \bar{y}_i and the sample size n_i for population i are given in Table 2.6 for $i = 1, \dots, 5$; the error mean square, $s^2 = 0.817$.

Table 2.6 Summary statistics for comparing μ_1, \dots, μ_5

i	1	2	3	4	5
n_i	6	5	6	6	7
\bar{y}_i	11.825	10.811	11.751	11.837	11.231

Let the null and alternative hypothesis be

$H_1 : \mu_1 \leq \mu_2, \mu_1 \leq \mu_3, \mu_4 \leq \mu_2, \mu_4 \leq \mu_3, \mu_1 \leq \mu_5$, and $H_2 : \text{not } H_1$, respectively. Assume that the setting in Table 2.3 holds; in particular, the observations are assumed to be *iid*. Let us apply the \bar{F} -test. It follows from the definition of \bar{F} ,

Table 2.7 Estimated p -values of \bar{F} for different error distributions

Dist:	$N(0, \sigma^2)$	T_4	T_{10}	χ_1^2	χ_4^2	χ_8^2	RS
p -value:	0.091	0.092	0.092	0.074	0.086	0.085	0.091

RS : resample with replacement from the empirical distribution of the within treatment residuals, which is the same as bootstrap.

(2.27), (2.24), and (2.16) that

$$\begin{aligned}\bar{F} &= \{RSS(H_1) - RSS(H_1 \cup H_2)\}/S^2 = \min_{\mu \in H_1} q(\mu)/S^2 \\ &= [\min_{A\mu \geq 0} (\bar{y} - \mu)^T W(\bar{y} - \mu)]/S^2\end{aligned}$$

where A is the 5×5 matrix $(-1, 1, 0, 0 | -1, 0, 1, 0, 0 | 0, 1, 0, -1, 0 | 0, 0, 1, -1, 0 | -1, 0, 0, 1)$; the five rows correspond to the five order restrictions in H_1 . We used the subroutine QPROG in IMSL for solving this constrained minimization problem. The sample value of \bar{F} is 5.04. If the error distribution is $F(t/\sigma)$, then p -value = $p_F(\bar{F} \geq 5.04 | \mu_1 = \dots = \mu_5)$, where the probability is computed when the error distribution is $F(t)$. In view of Theorem 2.2.1, the unknown σ can be set equal to one or any other convenient value when computing the p -value. Because the distribution function F is unknown, we approximate the p -value by the maximum of the p -values corresponding to several choices of F (see Table 2.7). In view of the results in Table 2.7, we conclude that the p -value ≥ 0.074 . The results in Table 2.7 show that the p -values are not that sensitive to the form of the error distribution in this particular example; we conjecture that this is likely to be the case more generally. ■

The foregoing results and simulation approaches to computing p -values extend to tests involving order restrictions on the regression parameter β in the linear model, $\mathbf{Y} = \mathbf{X}\beta + \mathbf{E}$ where \mathbf{E} has mean zero and covariance $\sigma^2 U$, σ is unknown and U is known; the inequality constraints may be in the null or alternative hypothesis. This includes (1) test of $H_0 : R\beta = 0$ against $H_1 : R_1\beta \geq 0$, where R_1 is a submatrix of R , and (2) test of $H_1 : R_1\beta \geq 0, R_2\beta = 0$ against $H_2 : \beta$ is unrestricted. These will be discussed in the next chapter.

Consider k populations with μ_i denoting a scalar parameter of interest for group i , $i = 1, \dots, k$. Suppose that the μ_i 's are known to satisfy the pairwise constraints,

$$\begin{aligned}\mu_i - \mu_j &\geq 0 \text{ for } (i, j) \in B \\ &\mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B\end{aligned}\tag{2.32}$$

where B is a given subset of $\{(i, j) : i, j = 1, \dots, k\}$. The objective in isotonic regression is to address statistical inference on μ_1, \dots, μ_k when they are constrained by (2.32).

Suppose that the random variable Y_i is an estimate of μ_i for $i = 1, \dots, k$; for example, Y_i could be an observation from a population with mean μ_i or it could be the mean of several *iid* observations from a population with mean μ_i . A natural method of estimating $\{\mu_1, \dots, \mu_k\}$ when they are known to satisfy (2.32) is to

$$\text{minimize } \sum_{i=1}^k (Y_i - \mu_i)^2 w_i \quad \text{subject to (2.32)}\tag{2.33}$$

where $\{w_1, \dots, w_k\}$ are given or suitably chosen weights. For example, if $Y_i \sim N(\mu_i, \sigma^2/w_i)$ for $i = 1, \dots, k$ and Y_1, \dots, Y_k are independent then (2.33) leads to the maximum likelihood estimate of (μ_1, \dots, μ_k) under (2.32). On the other hand, if Y_1, \dots, Y_k are not necessarily normal but the variance of Y_i is σ^2/w_i , where σ is unknown and w_1, \dots, w_k are known positive numbers, then (2.33) is the method of weighted least squares with weights, w_i ($i = 1, \dots, k$). For the time being, we shall not assume that Y_1, \dots, Y_k are normal.

Let $(\tilde{\mu}_1, \dots, \tilde{\mu}_k)$ denote the value of (μ_1, \dots, μ_k) , which solves the constrained minimization problem (2.33). Then we say that $(\tilde{\mu}_1, \dots, \tilde{\mu}_k)$ is the *isotonic regression* of (Y_1, \dots, Y_k) with respect to the weights (w_1, \dots, w_k) and the order restriction (2.32). Thus, an *isotonic regression* is a weighted least squares estimate subject to a set of pairwise constraints; different weights and/or pairwise constraints correspond to different isotonic regressions.

As an example, consider the one-way ANOVA setting with k treatments: $E(X_{ij}) = \mu_i$, $\text{var}(X_{ij}) = \sigma^2$ for $j = 1, \dots, n_i$ and $i = 1, \dots, k$, and the X_{ij} 's are independent. Let \bar{Y}_i denote the sample mean \bar{X}_i for $i = 1, \dots, k$; then $E(\bar{Y}_i) = \mu_i$ and $\text{var}(\bar{Y}_i) = \sigma^2/n_i$. Suppose that μ_1, \dots, μ_k are known to satisfy

$$\text{simple order : } \mu_1 \leq \dots \leq \mu_k.$$

Clearly this order restriction is of the form (2.32). Then the isotonic regression of (Y_1, \dots, Y_k) with respect to the weights (n_1, \dots, n_k) and the simple order is the *value* of (μ_1, \dots, μ_k) , which minimizes $\sum(Y_i - \mu_i)^2 n_i$ subject to $\mu_1 \leq \dots \leq \mu_k$. Two other important examples of pairwise constraints are:

1. *Tree order*: $\mu_1 \leq \mu_2, \dots, \mu_1 \leq \mu_k$, also written as $\mu_1 \leq [\mu_2, \dots, \mu_k]$; and
2. *Umbrella order*: $\mu_1 \leq \dots \leq \mu_m \geq \dots \geq \mu_k$, where m is known.

Usually, the tree order arises when a control (= group 1) is compared with several treatments. It is also worth noting that $\mu_1 \leq \mu_2 = \mu_3$ is also of the type in (2.32), because the constraint $\mu_2 = \mu_3$ is equivalent to $\mu_2 \leq \mu_3$ and $\mu_3 \leq \mu_2$. As was noted in (2.17), the set of constraints in (2.32) can also be expressed as

$$A\mu \geq 0\tag{2.34}$$

for some matrix A in which each row is a permutation of the k -vector, $(-1, 1, 0, \dots, 0)$, and $\mu = (\mu_1, \dots, \mu_k)^T$. For example, for the simple and tree orders, the matrix A

in (2.34) of order $(k - 1) \times k$ takes the forms

$$\begin{bmatrix} -1 & 1 & 0 & \cdots & 0 \\ 0 & -1 & 1 & \cdots & 0 \\ \cdots & 0 & -1 & \cdots & 0 \\ 0 & \cdots & 0 & -1 & 1 \end{bmatrix} \text{ and } \begin{bmatrix} -1 & 1 & 0 & \cdots & 0 \\ -1 & 0 & 1 & \cdots & 0 \\ \cdots & 1 & 0 & \cdots & 0 \\ -1 & 0 & \cdots & 0 & 1 \end{bmatrix}, \quad (2.35)$$

respectively. Since

$$\sum (Y_i - \mu_i)^2 w_i = (\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{W}(\mathbf{Y} - \boldsymbol{\mu}),$$

where $\mathbf{Y} = (Y_1, \dots, Y_k)^T$, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_k)^T$, and $\mathbf{W} = \text{diag}\{w_1, \dots, w_k\}$, the minimization problem (2.33) can be expressed as

$$\text{minimize } (\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{W}(\mathbf{Y} - \boldsymbol{\mu}) \quad \text{subject to } A\boldsymbol{\mu} \geq 0.$$

This is a standard quadratic program; see subsection 2.1.2.2 for a general method of solving it.

There are two types of inference problems that arise in isotonic regression: (i) k is fixed, and (ii) k increases with the sample size. In most parts of this text, only the first case is considered. The second case arises when a real function is to be estimated subject to shape constraints, for example, the nonparametric estimation of (i) a monotone density function, (ii) unimodal density function, (iii) monotone regression function, and (iv) convex regression function; the statistical inference problems arising in such cases are more challenging and they will be discussed later.

2.3.1 Quasi-order and Isotonic Regression

Let the labels x_1, \dots, x_k denote the k groups and let $X = \{x_1, \dots, x_k\}$. Let x be a qualitative explanatory variable taking values in X ; for group i , the explanatory variable x takes the value x_i . Note that x_1, \dots, x_k are not necessarily real numbers. In this section, the following notation is adopted to identify a real function on X with a point in \mathbb{R}^k . To any given real function g on X , we associate the point \mathbf{g} in \mathbb{R}^k defined by $\mathbf{g} = (g_1, \dots, g_k)^T$ where $g_i = g(x_i)$ for $i = 1, \dots, k$. Conversely, to any given point $\mathbf{g} = (g_1, \dots, g_k)^T$ in \mathbb{R}^k we associate the real function g on X defined by $g(x_i) = g_i$ for $i = 1, \dots, k$. Therefore, we shall freely interchange the roles of functions like g on X and the corresponding vector, $\mathbf{g} = (g_1, \dots, g_k)^T$.

Let us now consider our original problem of estimating μ_1, \dots, μ_k subject to (2.32). Let y_i denote an estimator of μ_i , $i = 1, \dots, k$, and let

$$y(x_i) = y_i, \quad \mathbf{y} = (y_1, \dots, y_k)^T, \quad \mu(x_i) = \mu_i, \quad \text{and } \boldsymbol{\mu} = (\mu_1, \dots, \mu_k)^T.$$

Consider the regression model

$$y(x_i) = \mu(x_i) + \epsilon_i, \quad (2.36)$$

where ϵ_i is the error term, $i = 1, \dots, k$. Thus, we have a regression model in which the domain of the regression function $\mu(x)$ is the set of labels X , and the function

$\mu(x)$ satisfies the restrictions (2.32). An example of the foregoing setting is when y_i is the sample mean response of several observations from group i and μ_i is the population mean. Another example is where μ_i is the probability of success and y_i is the sample proportion of successes in group i , $i = 1, \dots, k$.

First, let us recall some definitions on binary relations. A binary relation \preceq on the set $X = \{x_1, \dots, x_k\}$ is said to be (a) *reflexive* if $x \preceq x$ for every $x \in X$; (b) *transitive*, if $x \preceq y$ and $y \preceq z$ imply $x \preceq z$; (c) *antisymmetric* if $x \preceq y$ and $y \preceq x$ imply $x = y$; (d) a *partial order* if it is reflexive, transitive, and antisymmetric; and (e) a *quasi-order* if it is reflexive and transitive. If \preceq is a quasi-order and if $x_i \preceq x_j$, then we may interpret it as “group i is not higher than group j .”

Let \preceq be a given quasi-order on the set of labels $X = \{x_1, \dots, x_k\}$. A real function, μ , on X is said to be *isotonic* with respect to \preceq if

$$x_i \preceq x_j \text{ implies that } \mu(x_i) \leq \mu(x_j), \text{ for every } x_i \text{ and } x_j \text{ in } X.$$

Thus, an isotonic function has the *same tone* as the quasi-order, and hence the term *iso tone*.

Suppose that $\{\mu_1, \dots, \mu_k\}$ satisfy the constraints (2.32). Then, these constraints induce the order \preceq on X , defined by

$$x_i \preceq x_j \quad \text{if and only if} \quad (i, j) \in B.$$

It is easily seen that the induced relation \preceq is a quasi-order. Conversely, suppose that a quasi-order \preceq on X is given. Then, define the corresponding set B in (2.32) as $\{(i, j) : x_i \preceq x_j\}$. Therefore, the pairwise constraints (2.32) may be referred to as a quasi-order.

It follows that estimation of $\boldsymbol{\mu}$ subject to (2.32), and estimation of the function $\mu(x)$ subject to it being isotonic with respect to the quasi-order induced by (2.32) are equivalent problems. A formal statement is given in the next proposition.

Proposition 2.3.1 *For a given a quasi-order \preceq on X there exists a matrix A such that each row of A is a permutation of the k -vector $(1, -1, 0, \dots, 0)$, and*

$$A\boldsymbol{\mu} \geq 0 \text{ is equivalent to } \mu \text{ is isotonic with respect to } \preceq. \quad (2.37)$$

Conversely, given a matrix A in which each row is a permutation of the k -vector $(1, -1, 0, \dots, 0)$, there exists a quasi-order \preceq on X such that (2.37) holds. ■

Let y be a given real function on X and w be a given nonnegative function on X . Let

$$\mathcal{F} = \{\mu : \mu \text{ is an isotonic function on } X\} \quad \text{and} \quad \mathcal{Q} = \{\mu : A\boldsymbol{\mu} \geq 0\}$$

where A is as in Proposition 2.3.1. It follows from the foregoing proposition that \mathcal{F} and \mathcal{Q} are equivalent, in the sense that there is a one-to-one correspondence between \mathcal{F} and \mathcal{Q} .

An isotonic function y^* on X is said to be *the least squares isotonic regression of y with respect to the weight w* if

$$\sum_{x \in X} \{y(x) - \mu(x)\}^2 w(x) \quad (2.38)$$

reaches its minimum over $\mu \in \mathcal{F}$ at $\mu = y^*$. It is also possible to consider L^P isotonic regression with respect to w by minimizing $\sum |y(x) - \mu(x)|^p w(x)$; we shall not consider them for the time being. If there is no possibility of ambiguity, then we shall refer to y^* simply as the *isotonic regression of y without referring to least squares explicitly.*

Let

$$W = \text{diag}\{w_1, \dots, w_k\}.$$

Then

$$\sum_{x \in X} \{y(x) - \mu(x)\}^2 w(x) = (\mathbf{y} - \boldsymbol{\mu})^T W(\mathbf{y} - \boldsymbol{\mu}).$$

Now, since \mathcal{F} and Q are equivalent sets, the following two minimization problems are also equivalent:

$$(A) \quad \underset{x \in X}{\text{minimize}} \quad \sum_{x \in X} \{y(x) - \mu(x)\}^2 w(x) \quad \text{subject to } \boldsymbol{\mu} \in \mathcal{F}, \quad (2.39)$$

$$(B) \quad \underset{\boldsymbol{\mu}}{\text{minimize}} \quad (\mathbf{y} - \boldsymbol{\mu})^T W(\mathbf{y} - \boldsymbol{\mu}) \quad \text{subject to } A\boldsymbol{\mu} \geq \mathbf{0}. \quad (2.40)$$

Let y^* denote the isotonic regression of y in the foregoing minimization problem (A). Then y^* is the *least square projection* of y onto Q with respect to the inner product $(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T W \mathbf{y}$. Such projections will be studied in more detail in a later chapter; see Section 2.1.2.2 for some comments on computations relating to the problem (B).

2.4 ISOTONIC REGRESSION: RESULTS RELATED TO COMPUTATIONAL FORMULAS

In this section, several results relating to the computation of the isotonic regression are discussed. They are included here because the details are instructive and/or the results may be useful for studying the properties of the estimator. No attempt is made to provide a comprehensive discussion of the numerical algorithms that may be used for computing the isotonic regression; a discussion of such computational issues is outside the scope of this book. Since the computation of the isotonic regression is a problem in quadratic programming, interested readers may wish to consult the relevant literature in optimization. In this chapter, we will provide only an introduction.

2.4.1 Isotonic Regression Under Simple Order

Let y_i and μ_i denote an observation and a parameter of interest, respectively, for group i , $i = 1, \dots, k$. Assume that $y_i = \mu_i + \epsilon_i$ and that the parameters, μ_1, \dots, μ_k are unknown but known to satisfy the simpler order $\mu_1 \leq \dots \leq \mu_k$. Suppose that it is desired to estimate these parameters by the method of weighted least squares:

$$\min_{\mu_1 \leq \dots \leq \mu_k} \sum_{i=1}^k (y_i - \mu_i)^2 w_i \quad (2.41)$$

where w_1, \dots, w_k are given positive numbers. This problem has been studied extensively in the early literature, and detailed discussions may be found in Barlow et al. (1972) and Robertson et al. (1988).

Let y^* denote the constrained least squares estimate of μ resulting from (2.41); y^* is also the isotonic regression of y with respect to weights $\{w_i\}$ and the simple order. It turns out that y^* has simple and instructive explicit forms. It can be computed by a very simple algorithm known as the *Pool Adjacent Violators Algorithm*; this is described below.

Pool Adjacent Violators Algorithm (PAVA)

A set of consecutive elements of X will be called a *block*. For example, $\{x_1, x_2\}$ and $\{x_3, x_4, x_5, x_6\}$ are blocks, but $\{x_1, x_2, x_6\}$ is not a block because 1, 2, and 6 are not consecutive integers. Initially, X is partitioned into the k blocks $\{x_1\}, \dots, \{x_k\}$ with one label in each block. Now, PAVA pools consecutive blocks as follows:

- (I) If $y_1 \leq \dots \leq y_k$ then stop the iteration and the solution is $y^* = y_j$ for $j = 1, \dots, k$.
- (II) Otherwise, let i be the smallest index such that $y_i > y_{i+1}$. Now, pool the categories x^i and x_{i+1} to form a single block. Let the response variable for this block be the weighted response, $(w_i y_i + w_{i+1} y_{i+1}) / (w_i + w_{i+1})$ with weight $(w_i + w_{i+1})$.
- (III) Now repeat the process for the new blocks with their weighted responses and corresponding weights until the responses are nondecreasing.
- (IV) Once this process is completed, y^* is equal to the weighted response of the block which contains x_i .

The following example illustrates the main steps of PAVA.

Example: Let $k = 7$ and the values of $\{y_i\}$ and $\{w_i\}$ be as in the top three lines of Table 2.8. The values of y_i^* are given in the last line. The calculations in the table show that the estimator of μ subject to the simple order $\mu_1 \leq \dots \leq \mu_7$ is $(-5/2, -5/2, 2/3, 2/3, 3, 5, 6)$.

Although the foregoing discussions are limited to the case when the objective function is a weighted sum of squares, the result has been extended to include the case when it is of the form $(\mathbf{y} - \boldsymbol{\mu})^T G(\mathbf{y} - \boldsymbol{\mu})$ where G is a positive definite matrix (Duz and Salvador (1988)) and the simple order set is replaced by an *acute cone*. The PAVA has also been extended to include more general quadratic programming problems and concave regression problems (Tang and Lin (1991), Qian (1994b)).

The isotonic regression under simple order has another form that is instructive. Although the variable x is treated as qualitative in isotonic regression, let us consider

Table 2.8 Illustration of PAVA

i	w_i	y_i	y_i^*
First pool:	1	1	1
Blocks	-1	-4	-4
Pooled w	{1, 2}	{3}	{2}
Pooled y	5/2	2	1
Second pool:	{1, 2}	{3, 4}	{3}
Blocks	-5/2	1	0
Pooled w	{1, 2}	{3}	2/3
Pooled y	-5/2	-5/2	2/3

the special case when it also has a quantitative interpretation. Let $\mu(x)$ be a nondecreasing function of x and let $F(x) = \int^x \mu(t)dt$. Then $F'(x)$ is a convex function and $\mu(x)$ is the left slope of $F(x)$. This suggests that an alternative approach to estimating $\mu(x)$ subject to the constraint that μ is nondecreasing, is to estimate F subject to the constraint that it is convex and then estimate $\mu(x_i)$ as the left slope of the estimate of F at x_i . In fact a version of this approach and PAVA are equivalent. Such an equivalence holds, with suitable modifications, even when x_1, \dots, x_k are not quantitative but k categories. These are discussed below.

Let P_0 denote the origin in two-dimensions, and let the coordinates of P_i be defined by

$$P_i \equiv (w_1 + \dots + w_i, w_1 y_1 + \dots + w_i y_i), \quad i = 1, \dots, k.$$

Let the *Cumulative Sum Diagram* (CSD) be the curve which joins the consecutive points, P_0, P_1, \dots, P_k by straight line segments. The *Greatest Convex Minorant* (GCM) of CSD is defined as

$$GCM(t) = \max\{g(t) : g \text{ is convex and } g(t) \leq CSD(t) \text{ for every } t\}. \quad (2.42)$$

Therefore, GCM of CSD is the maximum of the convex functions that do not lie above $CSD(t)$ at any t . Alternatively, $GCM(t)$ can be seen as the path of a taut string that joins P_0 and P_k such that P_0, \dots, P_k are on or above the string (see Figures 2.1 and 2.2). Now we have the following; see pages 9–10 in Barlow et al. (1972) or Theorem 1.2.1 in Robertson et al. (1988) for a proof.

Proposition 2.4.1 *For the simple increasing order, $\mu_1 \leq \dots \leq \mu_k$, we have that*

$$y_i^* = \text{left slope of } GCM(t) \text{ at } t = (w_1 + \dots + w_i), \quad i = 1, \dots, k. \quad \blacksquare$$

The solution based on GCM can be given an explicit form using min-max and max-min functions. To this end, let

$$Av(s, t) = (w_s y_s + \dots + w_t y_t) / (w_s + \dots + w_t) \quad (2.43)$$

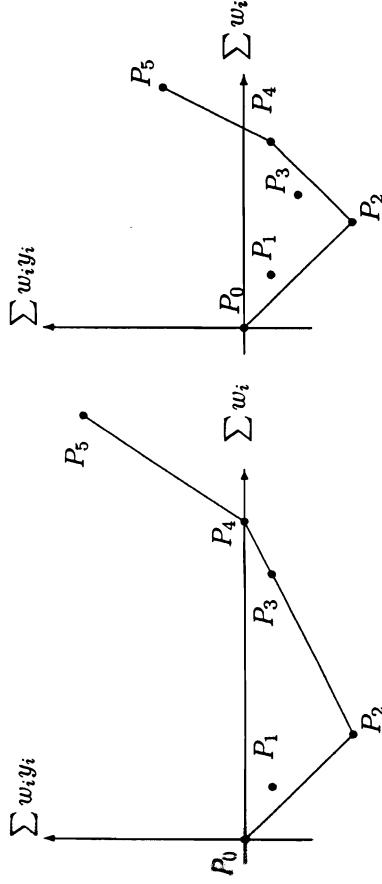


Fig. 2.1 Greatest convex minorant of $P_1 P_2 P_3 P_4 P_5$

for $t \geq s$. Thus, $Av(s, t)$ is the weighted mean of the observations y_s, \dots, y_t and it is also the left slope of the line segment $P_{s-1}P_t$. Then equivalent closed-forms of y_i^* given in the next result, and it is followed by an illustrative example; for more details see Barlow et al. (1972), Robertson et al. (1988), and Qian (1992).

Proposition 2.4.2 *For the simple increasing order, $\mu_1 \leq \dots \leq \mu_k$, we have that*

$$\begin{aligned} y_i^* &= \min_{t \geq i} \max_{s \leq i} Av(s, t) &= \min_{t \geq i} \max_{s \leq t} Av(s, t) \\ &= \max_{s \leq i} \min_{t \geq i} Av(s, t) &= \max_{s \leq i} \min_{t \geq s} Av(s, t) \end{aligned} \quad \blacksquare$$

Example: Let $k = 6$. The values of (w_i, y_i) and the coordinates of the points P_i are given in Table 2.9, and Fig. 2.1 shows the points. Let us denote the slope of $P_i P_j$ by $[P_i P_j]$. Now, by the foregoing Proposition, we have that

$$\begin{aligned} y_2^* &= \max_{s \leq 2} \min_{t \geq 2} Av(s, t) &= \max_{s \leq 2} \min_{t \geq 2} [P_{s-1}P_t] \\ &= \max_{s \leq 2} \{ \min_{t \geq 2} [P_0P_t], \min_{t \geq 2} [P_1P_t] \} \\ &= \max \{ \min \{ [P_0P_2], [P_1P_2], [P_0P_3], [P_0P_4], [P_0P_5], \\ &\quad \min \{ [P_1P_3], [P_1P_4], [P_1P_5] \} \} \\ &= \max \{ [P_0P_2], [P_1P_2] \} = [P_0P_2]. \end{aligned}$$

This simple example illustrates how the max-min formula is related to the left slope of the $GCM(t)$ at $t = w_1 + w_2$. Similarly, it would be instructive to verify that y_3^* is equal to the left slope of the $GCM(t)$ at $t = (w_1 + \dots + w_3)$. Fig. 2.2 shows another example of a greatest convex minorant of $P_0 P_1 P_2 P_3 P_4 P_5$. There are more general versions of the *min max* and *max min* formulas; see Section 1.4 in Robertson et al. (1988).

Table 2.9 Greatest convex minorant

	Greatest convex minorant					
i	1	2	3	4	5	6
w_i	1	1	1	2	2	2
y_i	-1	-3	1	1	3	4
Pooled blocks	{1,2}	{3}	{4}	{5}	{6}	
Pooled w	-2	1	1	3	4	
Pooled y						
Final Estimate	-2	-2	1	1	3	4
y_i^*						
Greatest convex minorant :						
$\sum_{j=1}^i w_j$	1	2	5	6	8	10
$\sum_{j=1}^i w_j y_j$	-1	-4	-1	0	6	14
P_i	(1,1)	(2,4)	(5,-1)	(6,0)	(8,6)	(10,14)

2.4.2 Ordered Means of Exponential Family

This section contains more advanced topics and it can be ignored at first reading. Further, most of the results in this section will not be required until Chapter 7.

There are different forms of what is known as the exponential family. First, we shall consider the form that is used in Generalized Linear Models where there may be an unknown scale parameter (for example, see McCullagh and Nelder (1989)); another form in which there are no unknown scale parameters will be considered at the end of this subsection. Suppose that the variable y has probability function,

$$f(y; \theta, \phi) = \exp\left\{\frac{y\theta - b(\theta)}{a(\phi)} + c(y; \phi)\right\}, \quad (2.44)$$

where θ and ϕ are scalars and $a(\cdot) > 0$, $b(\cdot)$ and $c(\cdot)$ are some functions. The function f may be the probability density of a continuous random variable such as normal or gamma, or it may the probability function of a discrete random variable such as poisson or binomial. Let $\ell(y; \theta, \phi)$ denote $\log f(y; \theta, \phi)$. Then, it is well known that (for example, see McCullagh and Nelder (1989), p 22)

$$\begin{aligned} (\partial/\partial\theta)\ell(\theta, \phi) &= \{y - b(\theta)\}/a(\phi), & (\partial^2/\partial\theta^2)\ell(\theta, \phi) &= -\dot{b}(\theta)/a(\phi), \\ E(y) &= \mu = \dot{b}(\theta) \text{ and } \text{var}(y) = \ddot{b}(\theta)a(\phi). \end{aligned}$$

Let $\hat{\theta}$ and $\hat{\mu}$ denote the unconstrained *mles* of θ and μ based on a single observation of y . Then, we have $\hat{\mu} = \dot{b}(\hat{\theta}) = y$.

Now, suppose that we have k groups and assume that the variable y_i in group i has probability function,

$$f(y_i; \theta_i, \phi_i) = \exp\left\{\frac{y_i\theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i; \phi_i)\right\}, \quad (2.45)$$

where θ_i and ϕ_i are scalars and $a(\cdot) > 0$, $b(\cdot)$ and $c(\cdot)$ are some functions. Let μ_i denote the population mean of y_i . The function $a(\phi_i)$ usually takes the form ϕ/m_i ;

ϕ and m_i are known as the dispersion parameter and prior weight, respectively. As an example, if y_i is the mean of m_i independent observations from $N(\mu_i, \sigma^2)$ then $a(\phi_i) = \sigma^2/m_i$. Let y_{i1}, \dots, y_{in_i} denote n_i independent observations from group i , $i = 1, \dots, k$. Let $\theta = (\theta_1, \dots, \theta_k)^T$, $\phi = (\phi_1, \dots, \phi_k)^T$, $\bar{y}_i = (y_{i1} + \dots + y_{in_i})/n_i$, $\bar{y} = (\bar{y}_1, \dots, \bar{y}_k)^T$. We are interested to obtain the *mle* of μ subject to

$$\mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B \quad (2.46)$$

where B is a given subset of $\{(i, j) : i, j = 1, \dots, k\}$. Let A be defined by

$$\{\mu : A\mu \geq 0\} = \{\mu : \mu_i - \mu_j \geq 0 \text{ for } (i, j) \in B\}.$$

It follows from (2.45) that the loglikelihood $L(\theta, \phi)$ can be expressed as

$$L(\theta, \phi) = \sum_{i=1}^k \frac{n_i \{\bar{y}_i\theta_i - b(\theta_i)\}}{a(\phi_i)} + H(\bar{y}; \phi) \quad (2.47)$$

where H is some function that does not depend on θ . Let $(\hat{\theta}, \hat{\phi})$ denote the unconstrained *mle* of (θ, ϕ) ; thus $(\hat{\theta}, \hat{\phi})$ is the unconstrained maximizer of $L(\theta, \phi)$. Let the corresponding estimator of μ be denoted by $\hat{\mu}$. If $\mu(\theta, \phi)$ expresses μ as a function of (θ, ϕ) then by the invariance of *mle*, we have that $\hat{\mu} = \mu(\hat{\theta}, \hat{\phi})$. Since the exponential family was expressed in the canonical form, $\hat{\mu} = \bar{y}$. Similarly, let $(\tilde{\phi}, \tilde{\mu})$ denote the *mle* of (ϕ, μ) subject to $\mu \in Q$ where $Q = \{\mu : A\mu \geq 0\}$. Since the *mle* is invariant under transformations, the constrained *mle* of θ is $\tilde{\theta} = \theta(\tilde{\mu})$. It turns out that, under some conditions concerning the dispersion parameter ϕ , $\tilde{\mu}$ is equal to the isotonic regression of μ with respect to some weights, and, therefore, $\tilde{\mu}$ can be obtained by minimizing a quadratic function of \bar{y} . This is stated in the next result; the proof is given in the Appendix.

Proposition 2.4.3 Suppose that the parameter θ is to be estimated by maximizing the function $L(\theta, \phi)$ in (2.47) subject to (2.46), which we write as $A\mu \geq 0$ where $\mu_i = E(y_i)$. Let $(\hat{\theta}, \hat{\mu})$ denote the resulting estimator of (θ, μ) .

(a) Assume that ϕ_1, \dots, ϕ_k are known constants. Then $\hat{\mu}$ is also the value of μ at which $\sum_i (\bar{y}_i - \mu_i)^2 \{n_i/a(\phi_i)\}$ reaches its minimum subject to $A\mu \geq 0$.

(b) Assume that $\phi_1 = \dots = \phi_k (= \phi, \text{say})$, and ϕ may be an unknown nuisance parameter. Then $\hat{\mu}$ is also the value of μ at which $\sum_i (\bar{y}_i - \mu_i)^2 n_i$ reaches its minimum subject to $A\mu \geq 0$. ■

To consider important special cases of this result, let y_{i1}, \dots, y_{in_i} denote n_i independent observations from any one of the following (here we use the same notation as in McCullagh and Nelder (1989) page 30):

1. Normal with mean μ_i and unknown variance σ^2 .
2. Poisson with mean μ_i .
3. Binary with $\text{pr}(y_{ij} = 1) = \mu_i$ and $\text{pr}(y_{ij} = 0) = 1 - \mu_i$.

4. Gamma with parameters (μ_i, ν) , where $\text{var}(y_i) = \mu_i^2/\nu$ and ν is unknown.

5. Inverse gaussian with parameter (μ_i, σ^2) where σ is unknown.

All of these cases belong to the canonical exponential family in (2.44). Therefore, for each of the five cases just listed, it follows from Proposition 2.4.3 that

$$\hat{\boldsymbol{\mu}} = \arg \min_{A\boldsymbol{\mu} \geq 0} \sum_{i \leq s} (\bar{y}_i - \mu_i)^2 n_i.$$

If the constraint $A\boldsymbol{\mu} \geq 0$ is the same as $\mu_1 \leq \dots \leq \mu_k$ then it follows from Proposition 2.4.2 that

$$\hat{\mu}_i = \min_{t \geq i} \max_{s \leq i} (n_s \bar{y}_s + \dots + n_t \bar{y}_t) / (n_s + \dots + n_t); \quad (2.48)$$

further, $\hat{\mu}_i$ can also be expressed as the left slope of a greatest convex minorant of a particular piecewise linear function as shown in Proposition 2.4.1.

Now, let us consider the case when ϕ_1, \dots, ϕ_k are unknown nuisance parameters; this was not considered in Proposition 2.4.3. Let $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\phi}})$ denote the solution of

$$\min L(\boldsymbol{\theta}, \boldsymbol{\phi}) \text{ subject to } A\boldsymbol{\mu} \geq 0.$$

It can be shown using Lemmas 2.5.3 and 2.5.2 in the appendix to this chapter that $\hat{\boldsymbol{\mu}}$ is still the isotonic regression of \bar{y} with respect to the weights $\tilde{w}_i = n_i/a(\tilde{\boldsymbol{\phi}}_i)$ as in part (a) of Proposition 2.4.3. However, since the weights depend on $\tilde{\boldsymbol{\phi}}$, which in turn may depend on $\hat{\boldsymbol{\mu}}$, it follows that $\hat{\boldsymbol{\mu}}$ is not just a function of \bar{y} alone, as is the case in parts (a) and (b) of Proposition 2.4.3. Therefore, the foregoing results do not provide a simple closed-form for computing $\hat{\boldsymbol{\mu}}$ when ϕ_1, \dots, ϕ_k are unknown.

If the conditions of Proposition 2.4.3 are satisfied, then it provides a convenient way of computing $\hat{\boldsymbol{\mu}}$ starting from the unconstrained estimator $\hat{\boldsymbol{\mu}}$. Further, it provides analytically simpler forms for the solution. This may be helpful in the study of the properties of the solutions. For example, to study the asymptotic properties of $\hat{\boldsymbol{\mu}}$, it may be convenient to use the fact that $\hat{\boldsymbol{\mu}}$ is obtained by minimizing a quadratic function of $\hat{\boldsymbol{\mu}}$, since properties of $\hat{\boldsymbol{\mu}}$ are easier to obtain, and minimizers of sum of squares are more tractable, this approach may turn out to be easier.

If there are no unknown scale parameters, then a result similar to the foregoing proposition holds even when the exponential family is not in the canonical form (2.44). Suppose that the density function of y_{ij} is

$$f(y; \boldsymbol{\theta}, \boldsymbol{\phi}_i) = \exp\{p_1(\boldsymbol{\theta}_i)p_2(\boldsymbol{\phi}_i)K(y; \boldsymbol{\theta}_i) + S(y; \boldsymbol{\phi}_i) + q(\boldsymbol{\theta}_i, \boldsymbol{\phi}_i)\},$$

where $(\partial/\partial\boldsymbol{\theta}_i)q(\boldsymbol{\theta}_i, \boldsymbol{\phi}_i) = -\boldsymbol{\theta}_i(\partial/\partial\boldsymbol{\theta}_i)p_1(\boldsymbol{\theta}_i)p_2(\boldsymbol{\phi}_i)$ and $\boldsymbol{\phi}_i$ is known ($i = 1, \dots, k$). Then, we have $E\{K(y_i; \boldsymbol{\phi}_i)\} = \boldsymbol{\theta}_i$ and $\text{var}\{K(Y_i; \boldsymbol{\phi}_i)\} = \{(\partial/\partial\boldsymbol{\theta}_i)p_1(\boldsymbol{\theta}_i)p_2(\boldsymbol{\phi}_i)\}^{-1}$. Further, the unconstrained mle $\hat{\boldsymbol{\theta}}_i$ of $\boldsymbol{\theta}_i$ is $(1/n_i)\{K(y_1; \boldsymbol{\phi}_i) + \dots + K(y_{n_i}; \boldsymbol{\phi}_i)\}$; let $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_k)^T$. Now, the mle of $\boldsymbol{\theta}$ subject to $A\boldsymbol{\theta} \geq 0$ is the isotonic regression of $\hat{\boldsymbol{\theta}}$ with weights $w_i = n_i p_2(\boldsymbol{\phi}_i)$ (see Theorem 1.5.2 in Robertson et al. (1988); they attribute this result to the student, Zehua Chen).

The next result (proof is given in the appendix to this chapter) appears on page 45. In Barlow et al. (1972); also on page 38 in Robertson et al. (1988). This result will be used for obtaining simpler formulas for the constrained mle of a multinomial parameter, and also the constrained nonparametric mle of a monotonic density function in a later chapter.

Corollary 2.4.4 (Maximum of a product). Let c_i and y_i be given positive numbers, s be a given real number,

$$L(\boldsymbol{\theta}) = \theta_1^{y_1} \dots \theta_k^{y_k},$$

A *be a given matrix in which every row is a permutation of the k -vector $(-1, 1, 0, \dots, 0)$, and $\mathbf{z}_i = (y_i/\sum y_i)(s/c_i)$ ($i = 1, \dots, k$). Then*

$$\begin{aligned} &\text{the maximum of } L(\boldsymbol{\theta}) \text{ subject to } A\boldsymbol{\theta} \geq 0 \text{ and } \sum c_i \theta_i = s \\ &\text{and} \end{aligned} \quad \begin{aligned} &\text{the minimum of } \sum (z_i - \theta_i)^2 c_i \text{ subject to } A\boldsymbol{\theta} \geq 0 \\ &\text{are reached at the same value of } \boldsymbol{\theta}. \end{aligned}$$

Example 2.4.1 Multinomial Parameters

Consider a multinomial distribution with k categories, and probability function

$$n!(y_1! \dots y_k!)^{-1} \theta_1^{y_1} \dots \theta_k^{y_k},$$

where $n = (y_1 + \dots + y_k)$ is the total number of observations. Let $\hat{\boldsymbol{\theta}}$ denote the mle of $\boldsymbol{\theta}$ subject to $A\boldsymbol{\theta} \geq 0$ and $\sum \theta_i = 1$. Then, we have that

$$\hat{\boldsymbol{\theta}} = \arg \min_{A\boldsymbol{\theta} \geq 0} \sum_{i=1}^k (\hat{\theta}_i - \theta_i)^2 \text{ where } \hat{\theta}_i = y_i/n \quad \text{for } i = 1, \dots, k.$$

This can be deduced from the previous corollary as follows.

Since $n!(y_1! \dots y_k!)^{-1}$ does not depend on $\boldsymbol{\theta}$ and $\sum \theta_i = 1$ is equivalent to $\sum c_i \theta_i = s$ where $c_i = 1 = s$, the results of the previous corollary on maximization of a product are applicable. Note that $z_i = \{y_i/(\sum y_i)\}(s/c_i) = (y_i/n) = \hat{\theta}_i$, which is the unrestricted mle of θ_i ($i = 1, \dots, k$). Further, $w_i = (y_i/z_i) = n$. Now, the desired result follows from the previous corollary. ■

2.5 APPENDIX: PROOFS

The proof of Proposition 2.4.3 involves several intermediate results that are of independent interest. These are stated below as lemmas. The proofs of these lemmas, given below, use various mathematical results on projections onto polyhedrals. A detailed account of the relevant results and definitions are given in the next chapter. Therefore, these proofs are meant to be read only if the reader is familiar with the

basic results on projections on to polyhedrals discussed in the Appendix to the next chapter. Proofs of the following lemmas are also given from first principles in Barlow et al. (1972); see also Section 1.7 in Robertson et al. (1988).

Lemma 2.5.1 *Let A be a matrix in which every row is a permutation of $(-1, 1, 0, 0, \dots, 0)$, $\mathcal{Q} = \{\mu : A\mu \geq 0\}$, and let \mathbf{y}^* denote the value of μ at which $\sum (y_i - \mu_i)^2 w_i$ reaches its minimum over $A\mu \geq 0$. For any real function ϕ on \mathbb{R} , define $\phi^* = (\phi(y_1^*), \dots, \phi(y_k^*))^T$. Then we have the following:*

1. *Let F be the unique face of \mathcal{Q} such that \mathbf{y}^* is in the relative interior of F . Then, ϕ^* is in the linear space spanned by F .*
2. *$(\mathbf{y} - \mathbf{y}^*)^T W \phi^* = 0$ where $W = \text{diag}\{w_1, \dots, w_k\}$.*

Proof: Let a_i denote a typical row of A . Let I and J be such that

$$F = \{\mu : a_i^T \mu = 0 \text{ for } i \in I, \text{ and } a_j^T \mu \geq 0 \text{ for } j \in J\}.$$

Then

$$ri(F) = \{\mu : a_i^T \mu = 0 \text{ for } i \in I, \text{ and } a_j^T \mu > 0 \text{ for } j \in J\},$$

where $ri(\cdot)$ denotes the relative interior. Let $\langle F \rangle = \{\mu : a_i^T \mu = 0 \text{ for } i \in I\}$, the linear space spanned by F . Then $(\mathbf{y} - \mathbf{y}^*) \perp_{W^{-1}} \langle F \rangle$. Suppose that $a_i^T \mu = 0$ for some $i \in I$ and that the nonzero elements of a_i occur at positions p and q . Then $\mu_p = \mu_q$. Consequently, $\phi(\mu_p) = \phi(\mu_q)$ and $a_i^T \phi_0 = 0$ where $\phi_0 = (\phi(\mu_1), \dots, \phi(\mu_k))^T$. Since $\mathbf{y}^* \in F \subset \langle F \rangle$, it follows that $a_i^T \mathbf{y}^* = 0$ and hence $a_i^T \phi^* = 0$ for every $i \in I$; hence, ϕ^* also belongs to $\langle F \rangle$. This completes the proof of the first part. The proof of the second part follows from $(\mathbf{y} - \mathbf{y}^*) \perp_{W^{-1}} \langle F \rangle$. ■

The first part of the lemma provides a geometric picture and the second part follows from the first. The next lemma says that the minimum of a weighted sum of squares and that of a more general function, both subject to the same constraints, are reached at the same point.

Lemma 2.5.2 *Let y_1, \dots, y_k be given real numbers, and let I be an interval containing $\{y_1, \dots, y_k\}$. Let Φ be a convex function that is finite on I and ∞ elsewhere. Let ϕ be the derivative of Φ (at a corner, choose a value between the left and right derivatives); assume that ϕ is finite on I . Let*

$$\Delta(u, v) = \begin{cases} \Phi(u) - \Phi(v) - (u - v)\phi(v), & \text{if } u, v \in I \\ \infty, & \text{if } u \notin I \text{ or } v \notin I. \end{cases}$$

For any two points, y and μ in \mathbb{R}^k , let

$$\Delta[\mathbf{y}, \mu] = \sum_{i=1}^k \Delta(y_i, \mu_i) w_i.$$

Let A be a matrix in which every row is a permutation of $(-1, 1, 0, 0, \dots, 0)$, $\mathcal{Q} = \{\mu : A\mu \geq 0\}$, and let \mathbf{y}^* denote the value of μ at which $\sum (y_i - \mu_i)^2 w_i$ reaches its minimum over $A\mu \geq 0$. Then,

$$(i) \quad \Delta[\mathbf{y}, \mu] \geq \Delta[\mathbf{y}, \mathbf{y}^*] + \Delta[\mathbf{y}^*, \mu], \text{ for any } \mu \in \mathcal{Q}.$$

$$(ii) \quad \mathbf{y}^* \text{ solves } \min_{\mu \in \mathcal{Q}} \Delta[\mathbf{y}, \mu].$$

$$(iii) \quad \mathbf{y}^* \text{ also solves } \max_{\mu \in \mathcal{Q}} \sum \{\Phi(\mu_i) + (y_i - \mu_i)\phi(\mu_i)\} w_i.$$

$$(iv) \quad \text{If } \Phi \text{ is strictly convex then } \Delta[\mathbf{y}, \mu] \text{ has a unique minimum over } \mu \in \mathcal{Q}.$$

Proof:

[This proof is essentially the same as that in Barlow et al. (1972), page 42]. It follows from the definition of Δ that, for real numbers r, s and t ,

$$\Delta(r, t) = \Delta(r, s) + \Delta(s, t) - (r - s)\{\phi(s) - \phi(t)\}.$$

Let μ be a point in \mathcal{Q} , $\phi_0 = (\phi(\mu_1), \dots, \phi(\mu_k))^T$, and $\phi^* = (\phi(y_1^*), \dots, \phi(y_k^*))^T$. It follows from (2.49) that

$$\Delta[\mathbf{y}, \mu] = \Delta[\mathbf{y}, t^*] + \Delta[t^*, \mu] - (\mathbf{y} - \mathbf{y}^*)^T W(\phi^* - \phi_0). \quad (2.49)$$

Since $A\mu \geq 0$ and ϕ is nondecreasing, it follows that $A\phi_0 \geq 0$ and hence $\phi_0 \in \mathcal{Q}$. Therefore, it follows from Lemma 3.12.3 (see page 114) that

$$(\mathbf{y} - \mathbf{y}^*)^T W \phi_0 \leq 0. \quad (2.50)$$

Further, by Lemma 2.5.1, we have that

$$(\mathbf{y} - \mathbf{y}^*)^T W \phi^* = 0. \quad (2.51)$$

Now, part (i) of the theorem follows from (2.49), (2.50), and (2.51). Part (ii) follows from part (i) because $\Delta[\cdot, \cdot] \geq 0$. Part (iii) follows from

$$\Delta[\mathbf{y}, \mu] = \sum_{i=1}^k \Phi(y_i) w_i - \sum_{i=1}^k \{\Phi(\mu_i) + (y_i - \mu_i)\phi(\mu_i)\} w_i]$$

and the fact that the first term, $\sum \Phi(y_i) w_i$, does not depend on μ . Part (iv) follows from $\Delta[\mathbf{y}^*, \mu] > 0$ if $\mathbf{y}^* \neq \mu$. ■

Lemma 2.5.3 *Let $L(\theta, \phi)$ denote the loss function in (2.47). Then,*

$$\begin{aligned} L(\theta, \phi) &= \sum_{i=1}^k \{y_i \theta_i - b(\theta_i)\} \{n_i/a(\phi_i)\} + \sum_{i,j} c(y_{ij}; \phi_i) \\ &= \sum_{i=1}^k \{\Phi(\mu_i) + \theta(\mu_i)(\bar{y}_i - \mu_i)\} \{n_i/a(\phi_i)\} + H(\mathbf{y}; \phi, \mu_0), \end{aligned}$$

for some function H , $\phi = (\phi_1, \dots, \phi_k)^T$, and $\Phi(t) = \int_{\mu_0}^t \theta(\mu) d\mu$.

Proof:

Let θ, μ and ϕ denote scalar parameters. Let $\dot{b}(\theta)$ and $\ddot{b}(\theta)$ denote the first and second derivatives of $b(\theta)$. For the exponential family in (2.47), we have that

$$\mu = \dot{b}(\theta) \text{ and } \text{var}(Y) = \ddot{b}(\theta)a(\phi);$$

for example, see McCullagh and Nelder (1989), page 29. Since $\ddot{b}(\theta) > 0$, it follows that $\dot{b}(\theta)$ is a strictly increasing function. Therefore, $\mu = \dot{b}(\theta)$ can be inverted, and hence we have that

$$\theta(\mu) = \dot{b}^{-1}(\mu), \mu = \dot{b}\{\theta(\mu)\}, \text{ and } (d/d\mu)b\{\theta(\mu)\} = \dot{b}(\theta)\dot{\theta}(\mu) = \mu\dot{\theta}(\mu).$$

Now, to obtain the unconstrained *mle*, $\hat{\theta}$, we solve $(\partial/\partial\theta)L(\theta, \phi) = 0$. This leads to

$$\hat{\mu}_i = \bar{y}_i = \dot{b}(\hat{\theta}_i) = \dot{b}\{\theta(\hat{\mu}_i)\}. \quad (2.52)$$

Let μ_0 be a value in the allowable range of μ . Then

$$\begin{aligned} b\{\theta(\mu)\} &= \int_{\mu_0}^{\mu} (d/dt)b\{\theta(t)\}dt + C(\mu_0) \\ &= \int_{\mu_0}^{\mu} t\dot{\theta}(t)dt + C(\mu_0) \\ &= \mu\dot{\theta}(\mu) - \int_{\mu_0}^{\mu} \theta(t)dt + C(\mu_0). \end{aligned} \quad (2.53)$$

Substituting (2.52) and (2.53) in (2.47), we have

$$\text{Likelihood} = \sum \{n_i/a(\phi_i)\} \{\Phi(\mu_i) + \theta(\mu_i)(\hat{\mu}_i - \mu_i)\} + H(\mathbf{y}, \phi)$$

where $\Phi(x) = \int_{\mu_0}^x \theta(t)dt$.

Proof of Proposition 2.4.3:

Follows from the previous two lemmas.

Proof of Corollary 2.4.4:

Let $w_i = (y_i/z_i)$ for $i = 1, \dots, k$, and $\ell(\theta) = \log L(\theta)$. Then

$$\ell(\theta) = \sum y_i \log \theta_i = \sum (z_i \log \theta_i)w_i. \quad (2.54)$$

Since $\sum (z_i - \theta_i)w_i = s^{-1}(\sum y_i)(s - \sum c_i\theta_i)$, it follows that

$$\sum c_i\theta_i = s \text{ if and only if } \sum (z_i - \theta_i)w_i = 0. \quad (2.55)$$

Let $\Phi(t) = t \log(t)$ and let Δ be defined as in Lemma 2.5.2. Then

$$\begin{aligned} \Delta[a, b] &= \Phi(a) - \Phi(b) - (a - b)\phi(b) \\ &= a \log a - a \log b - (a - b) \end{aligned}$$

and hence

$$\Delta[z_i, \theta_i]w_i = \sum (z_i \log z_i)w_i - \{\ell(\theta) + r(\theta)\} \quad (2.56)$$

where $r(\theta) = \sum (z_i - \theta_i)w_i$. Since $\sum (z_i \log z_i)w_i$ does not depend on θ , we can restrict our attention to $\{\ell(\theta) + r(\theta)\}$ for the purposes of finding the minimum of $\sum \Delta[z_i, \theta_i]w_i$. Let $\mathcal{A} = \{\theta : A\theta \geq 0\}$ and $\mathcal{R} = \{\theta : r(\theta) = 0\}$. Let

$$\tilde{\theta} = \arg \min_{\theta \in \mathcal{A}} \sum (z_i - \theta_i)^2 w_i.$$

Let $\phi(t) = 1$ for every $t \in \mathbb{R}$. Then, by Lemma 2.5.1, we have that $\phi^* = (\phi(\tilde{\theta}_1), \dots, \phi(\tilde{\theta}_k))^T = 1^T$, and $(z - \tilde{\theta})^T W 1 = 0$; hence $r(\tilde{\theta}) = \sum (z_i - \tilde{\theta}_i)w_i = 0$. Further,

$$\begin{aligned} \tilde{\theta} &= \arg \min_{\theta \in \mathcal{A}} \sum \Delta[z_i, \theta_i]w_i && \text{by Lemma 2.5.2, part(ii)} \\ &= \arg \max_{\theta \in \mathcal{A}} \{\ell(\theta) + r(\theta)\} && \text{by (2.56)} \end{aligned} \quad (2.57)$$

Now,

$$\begin{aligned} \max_{\theta \in \mathcal{A} \cap \mathcal{R}} \ell(\theta) &= \max_{\theta \in \mathcal{A} \cap \mathcal{R}} \{\ell(\theta) + r(\theta)\}, \text{ because } r(\theta) = 0 \text{ for } \theta \in \mathcal{R}. \\ &\leq \max_{\theta \in \mathcal{A}} \{\ell(\theta) + r(\theta)\}, \text{ because } \mathcal{A} \cap \mathcal{R} \subset \mathcal{A} \\ &= \ell(\tilde{\theta}) + r(\tilde{\theta}) && \text{by (2.57) and (2.58)} \\ &= \ell(\tilde{\theta}) && \text{because } r(\tilde{\theta}) = 0. \end{aligned} \quad (2.58)$$

Since $r(\tilde{\theta}) = 0$, we have that $\tilde{\theta} \in \mathcal{R}$, and hence $\max_{\theta \in \mathcal{A} \cap \mathcal{R}} \ell(\theta) = \ell(\tilde{\theta})$. Since Φ is strictly convex, $\tilde{\theta}$ is uniquely determined. ■

Problems**2.1**

In Example 1.2.2, assume that the error distribution is normal. Now, test the following hypotheses.

1. $H_0 : \mu_1 = \mu_2 = \mu_3 \quad \text{vs} \quad H_1 : \mu_1 \geq \mu_2 \geq \mu_3$.

2. $H_0 : \mu_1 \geq \mu_2 \geq \mu_3 \quad \text{vs} \quad H_1 : \mu_1 \geq \mu_2 \geq \mu_3 \text{ does not hold.}$

[Hint: Apply the \bar{F} -test and follow the steps as in the worked example on walking exercise.]

2.2

For the Example 2.1.2 on page 33, compute the *p*-values when the error distribution is (a) Normal, (b) Logistic, and (c) χ_5^2 for each of the following:

1. $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 \text{ vs } H_1 : \mu_1 \leq \mu_2 \leq \mu_3, \mu_2 \leq \mu_4$.
2. $H_1 : \mu_1 \leq \mu_2 \leq \mu_3, \mu_2 \leq \mu_4 \text{ vs } H_2 : \text{No restriction on } (\mu_1, \dots, \mu_4)$.

2.3 Consider the normal theory one-way ANOVA context in Table 2.1. Let H_1 be as in (2.9), and $(\tilde{\mu}_1, \dots, \tilde{\mu}_k)$ be the solution of

$$\min \sum \sum (y_{ij} - \mu_i)^2 \text{ subject to } H_1.$$

Show that $\tilde{\mu}_1, \dots, \tilde{\mu}_k$ is the constrained *mle* of (μ_1, \dots, μ_k) under H_1 .

[Hint: Let $L(\mu_1, \dots, \mu_k, \sigma)$ denote the log likelihood corresponding to $y_{ij} \sim N(\mu_i, \sigma^2)$. Then $L(\mu_1, \dots, \mu_k, \sigma) = -(2\sigma^2)^{-1} \sum \sum (y_{ij} - \mu_i)^2 - N \log \sigma^2$. For any given, (μ_1, \dots, μ_k) , L reaches its minimum at $\hat{\sigma}^2 = N^{-1} \sum \sum (y_{ij} - \mu_i)^2$. Therefore, the concentrated loglikelihood obtained by eliminating σ is

$$L^*(\mu_1, \dots, \mu_k) = L(\mu_1, \dots, \mu_k, \hat{\sigma}) = - \sum \sum (y_{ij} - \mu_i)^2 - 2N^{-1}.$$

The *mle* of (μ_1, \dots, μ_k) under H_1 is its value at which $\sum \sum (y_{ij} - \mu_i)^2$ reaches its minimum.]

2.4 Let y be a binary random variable taking the values 1 and 0 with probability μ and $(1 - \mu)$, respectively, ($0 < \mu < 1$). Let $f(y)$ denote the probability function of y . Verify that $f(y) = \mu^y(1 - \mu)^{1-y} = \exp[y \log\{\mu/(1 - \mu)\} + \log(1 - \mu)]$, and deduce that the probability function of y is of the form (2.45) with $\theta(\mu) = \log[\mu/(1 - \mu)]$, $a(\phi) = 1$ and $b(\theta) = -\log(1 - \mu)$.

Suppose that there are k groups and for group i , let y_{i1}, \dots, y_{in_i} denote n_i independent binary random variables with $\text{pr}(y_{ij} = 1) = \mu_i$ ($i = 1, \dots, k$). Verify that

$$\text{likelihood} = \prod_{i=1}^k \prod_{j=1}^{n_i} \mu_i^{y_{ij}} (1 - \mu_i)^{1-y_{ij}}.$$

Let A be a matrix with k columns in which every row is a permutation of $(1, -1, 0, \dots, 0)$. Use Proposition 2.4.3 and show that the *mle* $\tilde{\mu}$ of μ subject to $A\mu \geq 0$ is also the value of μ which

$$\text{maximizes } \sum (\bar{y}_i - \mu_i)^2 n_i \text{ subject to } A\mu \geq 0. \quad (2.59)$$

Now, suppose that the constraint $A\mu \geq 0$ is the same as $\mu_1 \leq \dots \leq \mu_k$. Show that

$$\tilde{\mu}_i = \min_{t \geq i} \max_{s \leq i} (n_s \bar{y}_s + \dots + n_t \bar{y}_t) / (n_s + \dots + n_t).$$

Suppose that $k = 5$, $(n_1, \dots, n_5) = (10, 15, 25, 12, 16)$ and $\bar{y} = (3/10, 5/15, 8/25, 8/12, 14/16)$. Use a diagram to illustrate that $\tilde{\mu}_i$ can also be expressed as the left slope of a greatest convex minorant of a particular piece-wise linear function. [Hint: see Proposition 2.4.1].

3

Tests on Multivariate Normal Mean

3.1 INTRODUCTION

In this chapter we consider two broad classes of hypothesis testing problems when the population distribution is normal. For example, the observations may be independently and identically distributed from the multivariate normal distribution, $N(\boldsymbol{\theta}, \mathbf{V})$, or they may satisfy a linear model of the form $\mathbf{Y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{E}$ with normal errors. For the standard problem of testing $H_0 : \mathbf{R}\boldsymbol{\theta} = \mathbf{0}$ against $H_2 : \mathbf{R}\boldsymbol{\theta} \neq \mathbf{0}$, where \mathbf{R} is a given fixed matrix, it is easy to apply the likelihood ratio test because the (likelihood ratio) test statistic can be computed easily, and the statistical tables for its null distribution, χ_q^2 where $q = \text{rank}(\mathbf{R})$, are easily available. In what follows, we shall use the acronym LRT for the likelihood ratio statistic and the test as well.

If the hypotheses involve inequalities in $\boldsymbol{\theta}$, then the theory becomes more complicated, and several difficulties in applying the results are also encountered. One such difficulty is caused by the fact that the null distribution of LRT depends on the particular inequalities. For example, when \mathbf{V} is known, the null distribution of LRT for testing

$$H_0 : \mathbf{R}\boldsymbol{\theta} = \mathbf{0} \text{ against } H_1 : \mathbf{R}\boldsymbol{\theta} \geq \mathbf{0}, \mathbf{R}\boldsymbol{\theta} \neq \mathbf{0} \quad (3.1)$$

depends on the matrix \mathbf{R} through $\mathbf{R}\mathbf{V}\mathbf{R}^T$, not just on $\text{rank}(\mathbf{R})$, as is the case for testing $H_0 : \mathbf{R}\boldsymbol{\theta} = \mathbf{0}$ against $H_2 : \mathbf{R}\boldsymbol{\theta} \neq \mathbf{0}$. Consequently, it is not practicable to have statistical tables of critical values because we would need one table for each $\mathbf{R}\mathbf{V}\mathbf{R}^T$ that may arise in practice, and there are infinite number of them. Another difficulty is that, except in some very special cases, it is difficult to compute the critical values exactly.

Fortunately, the p -values and critical values of the tests discussed in this chapter can be computed easily by simulation. A simple example illustrated this in the previous chapter. Its extension to the setting of this chapter is easy. The required computer programs can be developed so that they are applicable to most practical situations in which the constraints on θ are linear. The main requirement for this is a computer program for solving a quadratic program; such programs are available in IMSL, NAG and several other software packages. See Bazarra et al. (1993) for relevant theory and references.

Sometimes, constrained inference problems may lead to unexpected results; the simple example in Silvapulle (1997a) illustrated this. For example, the sample means of a bivariate response may appear to suggest that a treatment performed substantially worse than the control, but a likelihood ratio test may conclude to accept that the treatment is better than the control. The crux of these issues are also discussed. We adopt the general approach of computing the p -values approximately by simulation. For most practical purposes, the simulation approach is simple, adequate, and easy to implement.

The next section provides a discussion of the two types of testing problems that are discussed in this chapter, which we call Type A and Type B problems. Subsection 3.2.1 shows that for a large class of problems considered in this chapter, considerations may be restricted to a single observation from a normal distribution; the reasoning is based on sufficiency of the sample mean for a random sample from a multivariate normal. Some examples in two dimensions are discussed in detail in Section 3.3 to introduce the essential basics. A distribution, called chi-bar-square, arises naturally in Type A and Type B testing problems. This class of distributions is defined in Section 3.4. Computing the critical values of a chi-bar-square distribution is not as easy as reading a number from a table of critical values. In fact, except in some very special cases, it is virtually impossible to compute the tail probability of a chi-bar-square distribution using a hand calculator. Various results relating to the tail probability of a chi-bar-square distribution are discussed in Section 3.5. Then, in Sections 3.7 and 3.8 we consider the two types of testing problems in multivariate analysis with known covariance matrices. These results are extended in the following section to deal with the same type of testing problems in the linear model when the covariance matrix is $\sigma^2 U$ where σ^2 is unknown but U is known. Then, we consider the same two types of testing problems when the observations are from a multivariate normal and the covariance matrix is completely unknown. The proofs of several results are demanding, and therefore they are given in an Appendix to ensure that the main ideas and results are presented without the interruption of lengthy technical details.

3.2 STATEMENT OF TWO GENERAL TESTING PROBLEMS

Let us first define some terms. A set $\mathcal{A} \subset \mathbb{R}^p$ is said to be *convex* if $\{\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}\} \in \mathcal{A}$ whenever $\mathbf{x}, \mathbf{y} \in \mathcal{A}$ and $0 < \lambda < 1$. Therefore \mathcal{A} is a convex set if the line segment joining \mathbf{x} and \mathbf{y} is in \mathcal{A} whenever the points \mathbf{x} and \mathbf{y} are in \mathcal{A} . A set \mathcal{A} is said to be a *cone* with vertex \mathbf{x}_0 if $\mathbf{x}_0 + k(\mathbf{x} - \mathbf{x}_0) \in \mathcal{A}$ for every $\mathbf{x} \in \mathcal{A}$ and $k \geq 0$; if the

vertex is the origin, then we shall simply refer to it as a *cone*. Therefore, a cone is a set that consists of infinite straight lines starting from the origin. A *polyhedral* \mathcal{P} in \mathbb{R}^p is a set of the form $\{\theta \in \mathbb{R}^p : a_1^T \theta \geq 0, \dots, a_k^T \theta \geq 0\}$ where a_1, \dots, a_k are given elements of \mathbb{R}^p . Thus, \mathcal{P} is the intersection of the *half-spaces*, $\{a_1^T \theta \geq 0\}, \dots, \{a_k^T \theta \geq 0\}$. It is easily seen that a polyhedral is a closed convex cone. A more detailed discussion of polyhedrals is given in the first Appendix to this chapter.

Notation: Throughout this chapter, we shall adopt the following notation without any further comment:

- (a) C_a, C_b, C , and \mathcal{M} are subsets of an Euclidean space, for example \mathbb{R}^p .
- (b) C is a closed convex cone, \mathcal{M} is a linear space, $\mathcal{M} \subset C$, and $C_a \subset C_b$.

We shall consider two types of inequality constrained testing problems. These are stated below in their most general forms, where $\theta \in \mathbb{R}^p$.

- Type A:** Test $H_0 : \theta \in \mathcal{M}$ against $H_1 : \theta \in C$,
Type B: Test $H_1 : \theta \in C$ against $H_2 : \theta \notin C$.

Clearly, the null hypothesis in Type A is of the form $R\theta = \mathbf{0}$, for some matrix R . Strictly speaking, the alternative hypothesis in the foregoing Type A problem should be stated as $H_1 : \theta \in C, \theta \notin \mathcal{M}$. However, we shall continue to adopt the convention introduced in the previous chapter (see page 28): “*Test of H_0 against H_1 is to be interpreted as “test of H_0 against $H_1 - H_0$.”*” Therefore, the foregoing Type A and Type B testing problems may also be stated in the following equivalent forms:

- Type A:** Test $H_0 : \theta \in \mathcal{M}$ against $H_1 : \theta \in C$ and $\theta \notin \mathcal{M}$.
Type B: Test $H_1 : \theta \in C$ against $H_2 : \theta \in \mathbb{R}^p$.

Some examples of Type A problems:

1. $H_0 : \theta = \mathbf{0}$ against $H_1 : \theta \geq \mathbf{0}$.
2. $H_0 : (\theta_a^T, \theta_b^T)^T = \mathbf{0}$ against $H_1 : \theta_a \geq \mathbf{0}$.
3. $H_0 : R\theta = \mathbf{0}$ against $H_1 : R_1\theta \geq \mathbf{0}$, where R_1 is a submatrix of R .

4. $H_0 : A(\theta) = \mathbf{0}$ against $H_1 : A(\theta)$ is negative semi-definite, where $\theta = (\theta_1, \theta_2, \theta_3)^T$ and $A(\theta)$ is the 2×2 matrix, $(\theta_1, \theta_3 \mid \theta_3, \theta_2)$; this is equivalent to $H_0 : (\theta_1, \theta_2, \theta_3) = \mathbf{0}$ against $H_1 : \theta_1 \leq 0, \theta_2 \leq 0, \theta_1\theta_2 - \theta_3^2 \geq 0$.

In the last example, although the parameter space $\{\theta : \theta_1 \leq 0, \theta_2 \leq 0, \theta_1\theta_2 - \theta_3^2 \geq 0\}$ involves nonlinear inequalities, it turns out to be a closed convex cone. In view of the convention stated above, $\theta \neq \mathbf{0}, (\theta_a^T, \theta_b^T)^T \neq \mathbf{0}, R\theta \neq \mathbf{0}$, and $A(\theta) \neq \mathbf{0}$ are assumed to be implicit in the alternative hypothesis of the foregoing examples (1), (2), (3), and (4), respectively.

Some examples of Type B problems:

1. $H_1 : \theta \geq \mathbf{0}$ against $H_2 : \theta \not\geq \mathbf{0}$.
2. $H_1 : \theta_a \geq \mathbf{0}$ and $\theta_b = \mathbf{0}$ against $H_2 : \theta \in \mathbb{R}^p$, where $\theta = (\theta_a^T, \theta_b^T)^T$.

3. $H_1 : \mathbf{R}_1\boldsymbol{\theta} \geq 0$ and $\mathbf{R}_2\boldsymbol{\theta} = \mathbf{0}$ against $H_2 : H_1$ is not true, where \mathbf{R}_1 and \mathbf{R}_2 are some known fixed matrices.

4. $H_1 : \mathbf{A}(\boldsymbol{\theta})$ is negative semi-definite against $H_2 : H_1$ is not true, where $\boldsymbol{\theta} = (\theta_1, \theta_2, \theta_3)$ and $\mathbf{A}(\boldsymbol{\theta})$ is the 2×2 matrix $(\theta_1, \theta_3 \mid \theta_3, \theta_2)$; this is equivalent to $H_1 : \theta_1 \leq 0, \theta_2 \leq 0, \theta_1\theta_2 - \theta_3^2 \geq 0$, and $H_2 : \text{no restriction on } \boldsymbol{\theta}$.

These examples provide some idea of the type of testing problems studied in this chapter. The main difference between Type A and Type B problems is that inequalities may be present in the alternative hypothesis of Type A and the null hypothesis of Type B problems.

Here are some examples that are not Type A or Type B testing problems:

1. $H_0 : \theta_1 \leq 0$ or $\theta_2 \leq 0$ against $H_1 : \theta_1 > 0$ and $\theta_2 > 0$.
2. $H_0 : \boldsymbol{\theta} \geq \mathbf{0}$ or $\boldsymbol{\theta} \leq \mathbf{0}$ against $H_1 : \boldsymbol{\theta} \not\geq \mathbf{0}$ and $\boldsymbol{\theta} \not\leq \mathbf{0}$.
3. $H_0 : \theta_1 = \theta_2 = 0$ against $H_1 : \theta_2^2 \geq |\theta_1|$.

The first of these arises in clinical trials involving *combination drugs*; this testing problem is known as the *sign testing* problem. The second one arises in tests against *cross-over interactions* between treatments and groups. For the first two examples, some exact finite sample results are available. These will be studied in Chapter 9. The third is an artificial example; this example resembles a Type A problem but the alternative parameter space is not a closed convex cone. For such problems, large sample results will be studied in Chapter 4.

When the population distribution is $N(\boldsymbol{\theta}, \mathbf{V})$, the nature of the solutions to inequality constrained testing problems also depend on what is known about the covariance matrix \mathbf{V} in addition to the structure of the null and alternative parameter spaces. We shall consider three cases: (i) \mathbf{V} is a known positive definite matrix, (ii) $\mathbf{V} = \sigma^2 \mathbf{U}$, where \mathbf{U} is a known positive definite matrix and σ is unknown, and (iii) \mathbf{V} is unknown. In the next chapter, we shall study asymptotic results when \mathbf{V} is a function of some nuisance parameters.

Almost all of the results in this chapter are for Type A or Type B problems. Occasionally we shall consider test of $H_a : \boldsymbol{\theta} \in C_a$ against $H_b : \boldsymbol{\theta} \in C_b$ because it incorporates Type A and Type B problems and it is instructive to introduce some results within this general context.

3.2.1 Reduction by Sufficiency

Often, we shall consider test of hypotheses based on a single observation from a multivariate normal distribution. Although this case is considered for simplicity, it can also be justified using the sufficiency of sample mean. To illustrate this, let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be n independent and identically distributed observations from $N(\boldsymbol{\theta}, \mathbf{V})$, where \mathbf{V} is a known positive definite matrix and $\boldsymbol{\theta} \in \mathbb{R}^p$. Suppose that we wish to test

$$H_a : \boldsymbol{\theta} \in C_a \text{ against } H_b : \boldsymbol{\theta} \in C_b,$$

where, as was indicated earlier, $C_a \subset C_b$. In this case, the likelihood ratio test can be formulated in terms of a single observation of $\bar{\mathbf{X}}$, where $\bar{\mathbf{X}}$ is the sample mean. As is shown below, this is a consequence of the fact that $\bar{\mathbf{X}}$ is sufficient for $\boldsymbol{\theta}$.

Let $L(\boldsymbol{\theta})$ denote the loglikelihood for the n observations $\mathbf{X}_1, \dots, \mathbf{X}_n$. Since

$$L(\boldsymbol{\theta}) = -(n/2)^{-1} \log |\mathbf{V}| - (np/2) \log(2\pi) - 2^{-1} \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\theta})^T \mathbf{V}^{-1} (\mathbf{X}_i - \boldsymbol{\theta}),$$

we have that

$$L(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta}) + g(\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{V}),$$

where

$$\ell(\boldsymbol{\theta}) = -2^{-1} (\bar{\mathbf{X}} - \boldsymbol{\theta})^T (n^{-1} \mathbf{V})^{-1} (\bar{\mathbf{X}} - \boldsymbol{\theta})$$

and $g(\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{V})$ does not depend on $\boldsymbol{\theta}$. Therefore, $\ell(\boldsymbol{\theta})$ is the *kernel* of the loglikelihood for $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$. Now, the LRT for testing $H_a : \boldsymbol{\theta} \in C_a$ against $H_b : \boldsymbol{\theta} \in C_b$, based on $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ is

$$2[\max\{L(\boldsymbol{\theta}) : \boldsymbol{\theta} \in C_b\} - \max\{L(\boldsymbol{\theta}) : \boldsymbol{\theta} \in C_a\}].$$

Therefore, it follows that

$$LRT = 2[\max\{\ell(\boldsymbol{\theta}) : \boldsymbol{\theta} \in C_b\} - \max\{\ell(\boldsymbol{\theta}) : \boldsymbol{\theta} \in C_a\}]. \quad (3.2)$$

On the other hand, since $\bar{\mathbf{X}} \sim N(\boldsymbol{\theta}, n^{-1} \mathbf{V})$, the kernel of the loglikelihood for the single observation $\bar{\mathbf{X}}$ is $\ell(\boldsymbol{\theta})$. Therefore, the maximum likelihood estimator (*mle*) and the LRT based on $\bar{\mathbf{X}}$ and those based on $\mathbf{X}_1, \dots, \mathbf{X}_n$ are identical. Hence, in much of the developments in this chapter where the covariance matrix \mathbf{V} is known, we will consider statistical inference based on a single observation that may be thought of as $\bar{\mathbf{X}}$.

3.3 THEORY: THE BASICS IN TWO DIMENSIONS

Before we study general testing problems, it would be instructive to consider some simple special cases. We shall illustrate the derivation of the null distribution of the LRT for some special cases and then use the ideas developed to introduce the general results.

Example 3.3.1 Maximum Likelihood Estimation

Let $\mathbf{X} = (X_1, X_2)^T \sim N(\boldsymbol{\theta}, I)$, where I is the 2×2 identity matrix and $\boldsymbol{\theta} = (\theta_1, \theta_2)^T$. Consider the maximum likelihood estimation of $\boldsymbol{\theta}$ based on one observation on \mathbf{X} and subject to the constraint $\boldsymbol{\theta} \in C$ where C is the nonnegative orthant $\{(\theta_1, \theta_2)^T : \theta_1 \geq 0, \theta_2 \geq 0\}$. For the single observation \mathbf{X} , the kernel $\ell(\boldsymbol{\theta})$ of the loglikelihood is given by

$$-2\ell(\boldsymbol{\theta}) = \{(X_1 - \theta_1)^2 + (X_2 - \theta_2)^2\} = \|\mathbf{X} - \boldsymbol{\theta}\|^2$$

where $\|\cdot\|$ is the Euclidean norm defined by $\|\mathbf{x}\|^2 = (x_1^2 + x_2^2)$. Let $\tilde{\theta}$ be the *mle* of θ subject to $\theta \geq 0$. Since $-2\ell(\theta)$ is equal to the squared distance between \mathbf{X} and θ , $\tilde{\theta}$ is the point in C that is closest to \mathbf{X} ; in other words, $\tilde{\theta}$ is the *projection* of \mathbf{X} onto C (see Fig. 3.1). With Π denoting the projection function, we can write

$$\tilde{\theta} = \Pi(\mathbf{X} | \mathbb{R}^{+2}).$$

Let Q_1, Q_2, Q_3, Q_4 denote the four quadrants in the two-dimensional plane as in Fig. 3.1. Then the *mle* $\tilde{\theta}$ is given by

$$(\tilde{\theta}_1, \tilde{\theta}_2) = \begin{cases} (X_1, X_2) & \text{if } \mathbf{X} \in Q_1 \\ (0, X_2) & \text{if } \mathbf{X} \in Q_2 \\ (0, 0) & \text{if } \mathbf{X} \in Q_3 \\ (X_1, 0) & \text{if } \mathbf{X} \in Q_4 \end{cases} \quad (3.3)$$

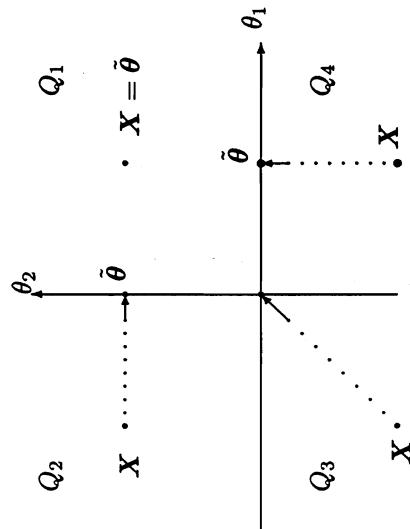


Fig. 3.1 The maximum likelihood estimator $\tilde{\theta}$ of θ subject to $\theta \geq 0$, based on a single observation of \mathbf{X} , where $\mathbf{X} \sim N(\theta, I)$.

Since the distribution function of \mathbf{X} is known and $\tilde{\theta}$ is a function of \mathbf{X} only, we can write down explicit expressions for the distribution of $\tilde{\theta}$. In contrast to the familiar situation where the parameter space for θ is \mathbb{R}^2 , here we observe that the distribution of the constrained *mle*, $\tilde{\theta}$, is not normal. Further, the distribution of $(\tilde{\theta} - \theta)$ and that of any scaled form of it depend on θ . Therefore, we cannot construct a confidence region for θ in the usual way based on the distribution of $(\tilde{\theta} - \theta)$. Consequently, the explicit form of the distribution of $(\tilde{\theta} - \theta)$ has not found much use in statistical inference.

Example 3.3.2 Type A testing problem (Example 3.3.1 continued).

Let $\mathbf{X} = (X_1, X_2)^T \sim N(\theta, I)$, where I is the 2×2 identity matrix and $\theta = (\theta_1, \theta_2)^T$. Consider the likelihood ratio test of $H_0 : \theta_1 = \theta_2 = 0$ against $H_1 : \theta_1 \geq 0, \theta_2 \geq 0$,

based on a single observation of \mathbf{X} . Since $-2\ell(\theta) = \|\mathbf{X} - \theta\|^2$ and $LRT = 2[\max\{\ell(\theta) : \theta \in H_1\} - \max\{\ell(\theta) : \theta \in H_0\}]$ we have that

$$LRT = \|\mathbf{X}\|^2 - \|\mathbf{X} - \tilde{\theta}\|^2 = \|\tilde{\theta}\|^2.$$

The last step follows since $(\mathbf{X} - \tilde{\theta})^T \tilde{\theta} = 0$ as is clear from Figure 3.1; it may also be verified by considering each of the four cases in (3.3) (see Fig. 3.1). Since we are interested in the distribution of LRT under the null hypothesis, let us suppose that the null hypothesis is true for the rest of these derivations. To obtain an expression for $\Pr(LRT \leq c)$, note that

$$\begin{aligned} \Pr(LRT \leq c) &= \sum_{i=1}^4 \Pr(LRT \leq c \mid \mathbf{X} \in Q_i) \Pr(\mathbf{X} \in Q_i) \\ &= \sum_{i=1}^4 \Pr(LRT \leq c \mid \mathbf{X} \in Q_i) \Pr(\mathbf{X} \in Q_i). \end{aligned} \quad (3.4)$$

Let us evaluate each of the four summands in the last expression.

Using the circular symmetry of $N(0, I)$, it can be shown that the direction ($= \mathbf{X}/\|\mathbf{X}\|$) and the length ($= \|\mathbf{X}\|$) of \mathbf{X} are statistically independent (see Exercise 3.1). Consequently, the conditional distribution of $(X_1^2 + X_2^2)$, given that the direction of \mathbf{X} is in Q_1 , is the same as that of its unconditional distribution. Therefore,

$$\begin{aligned} \Pr(LRT \leq c \mid \mathbf{X} \in Q_1) &= \Pr(X_1^2 + X_2^2 \leq c \mid \mathbf{X} \in Q_1) \\ &= \Pr(X_1^2 + X_2^2 \leq c) = \Pr(\chi_2^2 \leq c). \end{aligned} \quad (3.5)$$

Now, consider the second summand in (3.4) :

$$\begin{aligned} \Pr(LRT \leq c \mid \mathbf{X} \in Q_2) &= \Pr(X_2^2 \leq c \mid X_2 \geq 0, X_1 \leq 0) \\ &= \Pr(X_2^2 \leq c \mid X_2 \geq 0) \quad \text{since } X_1 \text{ and } X_2^2 \text{ are independent} \\ &= \Pr(X_2^2 \leq c) \quad \text{since } X_2 \text{ is symmetric} \\ &= \Pr(\chi_1^2 \leq c) \quad \text{since } X_2 \sim N(0, 1). \end{aligned} \quad (3.6)$$

Similarly, $\Pr(LRT \leq c \mid \mathbf{X} \in Q_4) = \Pr(\chi_1^2 \leq c)$. Therefore, we have that

$$LRT = \begin{cases} X_1^2 + X_2^2 & \text{given } \mathbf{X} \in Q_1 \\ X_2^2 & \text{given } \mathbf{X} \in Q_2 \\ 0 & \text{given } \mathbf{X} \in Q_3 \\ X_1^2 & \text{given } \mathbf{X} \in Q_4 \end{cases} \sim \begin{cases} \chi_2^2 & \text{given } \mathbf{X} \in Q_1 \\ \chi_1^2 & \text{given } \mathbf{X} \in Q_2 \\ 0 & \text{given } \mathbf{X} \in Q_3 \\ \chi_1^2 & \text{given } \mathbf{X} \in Q_4 \end{cases} \quad (3.7)$$

Now, it follows from (3.4)-(3.7) that the null distribution of the *LRT* is the following weighted sum of chi-square distributions:

$$\Pr(LRT \leq c) = (1/4) + (1/2)\Pr(\chi_1^2 \leq c) + (1/4)\Pr(\chi_2^2 \leq c) \quad (3.8)$$

for $c > 0$. For notational convenience, we define χ_0^2 , the *chi-square distribution with zero degrees of freedom*, to be the distribution that takes the value zero with probability

one. With this notation, the foregoing result is usually written as

$$\Pr(LRT \leq c \mid H_0) = \sum_{i=0}^2 w_i \Pr(\chi_i^2 \leq c),$$

where $(w_0, w_1, w_2) = (0.25, 0.5, 0.25)$. For later use, let us note that w_0, w_1 , and w_3 are the probabilities that \mathbf{X} falls in the cones $Q_3, Q_2 \cup Q_4$, and Q_1 respectively. It follows from (3.7) that the critical region, $\{LRT \geq c\}$, is the region to the upper-right of the curve $PQRS$ in Fig. 3.2; PQ is parallel to the x_1 -axis, RS is parallel to x_2 -axis and QR is a circular arc of radius \sqrt{c} . ■

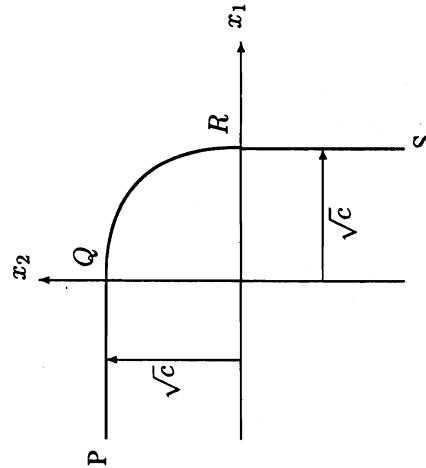


Fig. 3.2 The likelihood ratio test of $\theta = 0$ against $\theta \geq 0$ based on a single observation of $\mathbf{X} \sim N(\theta, I)$; the critical region, $\{\mathbf{x} : LRT \geq c\}$ is the upper-right region bounded by the curve $PQRS$.

The foregoing two examples show, in a very simple situation, how the constrained MLE over C behaves like a projection of a normal random variate onto C (see Fig. 3.1), and how a mixture of chi-square distributions arise when the alternative hypothesis involves inequality constraints. In fact, more general mixtures of chi-square distributions arise when the hypotheses involve other forms of inequalities. The ideas underlying such general results are introduced below using simple bivariate examples.

Example 3.3.3 Likelihood ratio test of $H_0 : \theta = 0$ against $H_1 : R\theta \geq 0$.

Let $\mathbf{X} = (X_1, X_2)^T \sim N(\theta, I)$. Consider the LRT of

$$H_0 : \theta = 0 \quad \text{against} \quad H_1 : R\theta \geq 0$$

where R is the 2×2 nonsingular matrix $(1, 4 \mid 1, -2)$ (i.e., the two rows of R are $(1, 4)$ and $(1, -2)$ respectively). Let $C = \{\theta : R\theta \geq 0\}$. Then, C is a closed convex cone (see Fig. 3.3).

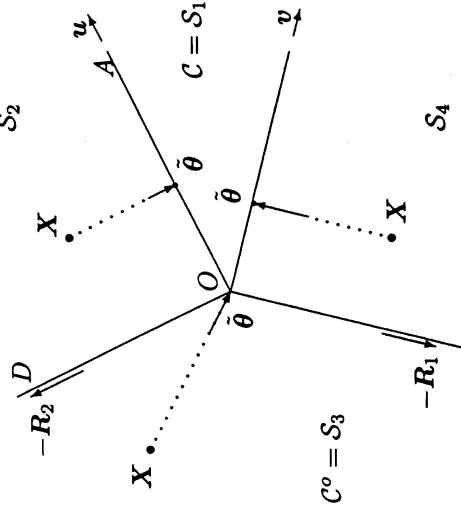


Fig. 3.3 The mle $\tilde{\theta}$ of θ based on a single observation of \mathbf{X} when $\mathbf{X} \sim N(\theta, I)$ and θ is constrained to lie in the cone C .

Let $C^o = \{\alpha : \alpha^T \theta \leq 0 \text{ for every } \theta \in C\}$; C^o is called the *negative dual or polar cone* of C with respect to the inner product, $\langle \alpha, \theta \rangle = \alpha^T \theta = \alpha_1 \theta_1 + \alpha_2 \theta_2$. Thus, C^o is the collection of vectors which do not form an acute angle with any vector in C . It may be verified that the boundaries of C^o are the perpendiculars to the boundaries of C , and that C^o is the closed convex cone formed by these perpendiculars (see Fig. 3.3). Let R_1^T and R_2^T denote the two rows of R . It may also be verified that $-R_1$ and $-R_2$ are parallel to the boundaries of C^o as shown in Fig. 3.3. Clearly, C and C^o partition the plane into 4 cones; let us denote these by S_1, S_2, S_3, S_4 with $S_1 = C$ and $S_3 = C^o$ (see Fig. 3.3). The partition $\{S_1, S_2, S_3, S_4\}$ in this example corresponds to $\{Q_1, Q_2, Q_3, Q_4\}$ in Example 3.3.2.

Let u and v denote unit vectors parallel to the upper and lower boundaries of C . Since the kernel of the loglikelihood is equal to $-0.5\|\mathbf{X} - \theta\|^2$, the mle of θ subject to $\theta \in C$ is the point in C closest to \mathbf{X} . In other words, the mle is the projection of \mathbf{X} onto C , which we write as $\Pi(\mathbf{X} \mid C)$. As in Examples 3.3.1 and 3.3.2, we obtain the following for the mle:

$$\tilde{\theta} = \begin{cases} \mathbf{X} & \text{given } \mathbf{X} \in S_1 \\ (u^T \mathbf{X})u & \text{given } \mathbf{X} \in S_2 \\ 0 & \text{given } \mathbf{X} \in S_3 \\ (v^T \mathbf{X})v & \text{given } \mathbf{X} \in S_4. \end{cases} \quad (3.9)$$

By considering each of the cases, $\mathbf{X} \in S_i, i = 1, 2, 3, 4$ separately, it may be verified that $(\mathbf{X} - \tilde{\theta})$ and $\tilde{\theta}$ are orthogonal when $\tilde{\theta}$ is not zero; Fig. 3.3 shows this clearly. Therefore, it follows from (3.2) that

$$LRT = \|\mathbf{X}\|^2 - \|\mathbf{X} - \tilde{\theta}\|^2 = \|\tilde{\theta}\|^2.$$

To derive the null distribution of LRT , first note that

$$\Pr(LRT \leq c) = \sum_{i=1}^4 \Pr(LRT \leq c \mid \mathbf{X} \in \mathcal{S}_i) \Pr(\mathbf{X} \in \mathcal{S}_i). \quad (3.10)$$

Suppose that H_0 holds. Then $\mathbf{X} \sim N(\mathbf{0}, I)$. Therefore, the length and direction of \mathbf{X} are independent (see Exercise 3.1). Hence $\|\mathbf{X}\|^2 \leq c$ and $\mathbf{X} \in \mathcal{S}_i$ are independent, and

$$\begin{aligned} \Pr(LRT \leq c \mid \mathbf{X} \in \mathcal{S}_i) &= \Pr(X_1^2 + X_2^2 \leq c \mid \mathbf{X} \in \mathcal{S}_i) \\ &= \Pr(X_1^2 + X_2^2 \leq c) = \Pr(\chi_2^2 \leq c). \end{aligned}$$

To obtain an expression for $\Pr(LRT \leq c \mid \mathbf{X} \in \mathcal{S}_2)$, note that we can think of $O\mathbf{A}$ and OD as the first and second axes. In this new orthogonal coordinate system, the LRT given $\mathbf{X} \in \mathcal{S}_2$ is the squared length of the first coordinate and hence is distributed as χ_1^2 ; essentially the same argument was used in the previous example leading to (3.6). Similarly, the LRT given $\mathbf{X} \in \mathcal{S}_4$ is distributed as χ_1^2 . Thus, as in Examples 3.3.1 and 3.3.2, we obtain the following for the conditional null distribution of LRT :

$$LRT = \begin{cases} X_1^2 + X_2^2 & \text{given } \mathbf{X} \in \mathcal{S}_1 \\ (\mathbf{u}^T \mathbf{X})^2 & \sim \begin{cases} \chi_2^2 \\ \chi_1^2 \end{cases} \\ 0 & \text{given } \mathbf{X} \in \mathcal{S}_3 \\ (\mathbf{v}^T \mathbf{X})^2 & \sim \begin{cases} \chi_2^2 \\ \chi_1^2 \end{cases} \\ \{\mathbf{X}, C\} & \text{given } \mathbf{X} \in \mathcal{S}_4. \end{cases} \quad (3.11)$$

Collecting these results and substituting in (3.10), we have that

$$\Pr(LRT \leq c \mid H_0) = q\Pr(\chi_0^2 \leq c) + 0.5\Pr(\chi_1^2 \leq c) + (0.5 - q)\Pr(\chi_2^2 \leq c), \quad (3.12)$$

where $q = \Pr(\mathbf{X} \in \mathcal{S}_3 \mid H_0) = \Pr(\mathbf{X} \in C^o \mid H_0) = (2\pi)^{-1}\gamma$ and γ is the angle (in radians) of the cone C^o at its vertex; therefore,

$$q = (2\pi)^{-1} \cos^{-1} [\mathbf{R}_1^T \mathbf{R}_2 / \{\mathbf{R}_1^T \mathbf{R}_1\} \{\mathbf{R}_2^T \mathbf{R}_2\}]^{1/2}], \quad (3.13)$$

where \mathbf{R}_1^T and \mathbf{R}_2^T are the two rows of \mathbf{R} . The critical region $\{LRT \geq c\}$ is the region to the right of the curve $PQRS$ in Figure 3.4; PQ is orthogonal to the upper boundary of C , RS is orthogonal to the lower boundary of C and QR is a circular arc of radius \sqrt{c} . ■

Example 3.3.4 Likelihood ratio test of $H_0 : \boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \mathbf{R}\boldsymbol{\theta} \geq \mathbf{0}$ when $\mathbf{X} \sim N(\boldsymbol{\theta}, V)$ and V is arbitrary.

Let $\mathbf{X} = (X_1, X_2)^T \sim N(\boldsymbol{\theta}, V)$ where V is an arbitrary positive definite matrix. Consider the LRT of $H_0 : \boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \boldsymbol{\theta} \in C$ where $C = \{\boldsymbol{\theta} : \mathbf{R}\boldsymbol{\theta} \geq \mathbf{0}\}$ is a closed convex cone and \mathbf{R} is a 2×2 nonsingular matrix. We can transform this problem to the set up in Example 3.3.3. Let V^{-1} be factorized as $V^{-1} = \mathbf{A}^T \mathbf{A}$, for some \mathbf{A} (for example, use Cholesky decomposition; see Anderson (1984) page 586);

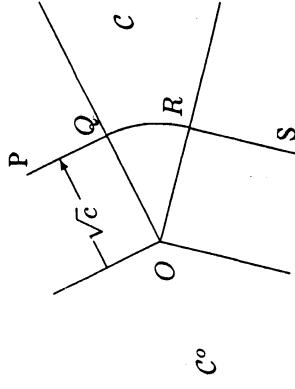


Fig. 3.4 The region to the right of $PQRS$ is the critical region $\{\mathbf{x} : \|\tilde{\boldsymbol{\theta}}(\mathbf{x})\|^2 \geq c\}$ for the LRT of $H_0 : \boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \boldsymbol{\theta} \in C$ based on a single observation of \mathbf{X} where $\mathbf{X} \sim N(\boldsymbol{\theta}, I)$; $\tilde{\boldsymbol{\theta}}(\mathbf{x})$ is the mle of $\boldsymbol{\theta}$ subject to $\boldsymbol{\theta} \in C$ when $\mathbf{X} = \mathbf{x}$, as shown in Figure 3.3. if $\mathbf{V} = [v_{ij}]$ then $\mathbf{A} = (v_{11}v_{22} - v_{12}^2)^{-1/2} [\sqrt{v_{22}}, -v_{12}/\sqrt{v_{22}} \mid 0, v_{11} - (v_{12}^2/v_{22})]$. Now, transform \mathbf{X} and the associated quantities, $\boldsymbol{\theta}$ and C , by \mathbf{A} :

$$\mathbf{Y} = \mathbf{AX}, \boldsymbol{\alpha} = \mathbf{A}\boldsymbol{\theta}, \mathcal{P} = \mathbf{AC}$$

where $\mathbf{AC} = \{\mathbf{A}\boldsymbol{\theta} : \boldsymbol{\theta} \in C\}$. Then $\mathcal{P} = \{\boldsymbol{\alpha} : \mathbf{RA}^{-1}\boldsymbol{\alpha} \geq \mathbf{0}\}$, $\mathbf{Y} \sim N(\boldsymbol{\alpha}, I)$, and the testing problem is equivalent to the test of $H_0 : \boldsymbol{\alpha} = \mathbf{0}$ against $H_1 : \boldsymbol{\alpha} \in \mathcal{P}$ based on \mathbf{Y} . This is precisely the set-up in the last example with $\{\mathbf{Y}, \mathcal{P}\}$ in place of $\{\mathbf{X}, C\}$. Since the mle and LRT are invariant under nonsingular transformations, we have that

$$\tilde{\boldsymbol{\theta}} = \mathbf{A}^{-1}\tilde{\boldsymbol{\alpha}} \quad \text{and} \quad LRT = \|\tilde{\boldsymbol{\alpha}}\|^2 = \tilde{\boldsymbol{\theta}}^T \mathbf{V}^{-1} \tilde{\boldsymbol{\theta}}.$$

Now, by (3.13) and (3.12), we have that

$$\Pr(LRT \leq c \mid H_0) = q\Pr(\chi_0^2 \leq c) + 0.5\Pr(\chi_1^2 \leq c) + (0.5 - q)\Pr(\chi_2^2 \leq c), \quad (3.14)$$

where

$$q = (2\pi)^{-1} \cos^{-1} [(\mathbf{R}_1^T \mathbf{V} \mathbf{R}_2) / \{(\mathbf{R}_1^T \mathbf{V} \mathbf{R}_1)\} \{(\mathbf{R}_2^T \mathbf{V} \mathbf{R}_2)\}]^{1/2}], \quad (3.15)$$

and \mathbf{R}_1^T and \mathbf{R}_2^T are the two rows of \mathbf{R} .

This example shows that we do not have to consider the case for a general V separately because the inference problem can be restated in terms of the identity covariance matrix and therefore can be handled by the method for the previous example. However, it is instructive to review the details in terms of the original variables and parameters without using such a transformation because some important features tend to be masked by the transformation.

Geometry of mle and LRT when $\mathbf{X} \sim N(\boldsymbol{\theta}, V)$, $V \neq I$.

Since V^{-1} is positive definite, it induces the inner product defined by $\langle \mathbf{x}, \mathbf{y} \rangle_V = \mathbf{x}^T V^{-1} \mathbf{y}$. This further induces the distance function defined by $\|\mathbf{x} - \mathbf{y}\|_V = \{(\mathbf{x} - \mathbf{y})^T V^{-1} (\mathbf{x} - \mathbf{y})\}^{1/2}$.

$\langle \cdot \rangle_V$ if $\langle \mathbf{x}, \mathbf{y} \rangle_V = 0$ and denote it by $\mathbf{x} \perp_V \mathbf{y}$. If necessary, we shall use the abbreviation V -orthogonal, V -distance, etc. to indicate that the relevant inner product is $\langle \cdot \rangle_V$.

Let \mathcal{B} be a given subset of \mathbb{R}^2 and let $\tilde{\theta}$ denote the mle of θ over $\theta \in \mathcal{B}$ based on a single observation of \mathbf{X} where $\mathbf{X} \sim N(\theta, V)$. Then, since $\ell(\theta) = (-1/2)(\mathbf{X} - \theta)^T V^{-1}(\mathbf{X} - \theta) = (-1/2)\|\mathbf{X} - \theta\|_V^2$, it follows that $\tilde{\theta} = \arg \min_{\theta \in \mathcal{B}} \|\mathbf{X} - \theta\|_V^2$ and hence $\tilde{\theta}$ is the point in \mathcal{B} that is V -closest to \mathbf{X} . In this case, we say that $\tilde{\theta}$ is the V -projection of \mathbf{X} onto \mathcal{B} and denote it by $\Pi_V(\mathbf{X} \mid \mathcal{B})$ or equivalently by $\Pi_V(\mathbf{X}, \mathcal{B})$. More details about such projections are given in the appendix to this chapter. The following result relating to the contours of $N(0, V)$ is useful.

Proposition 3.3.1 *The ellipse $\mathbf{x}^T V^{-1} \mathbf{x} = \text{Constant}$ and any line through the center of the ellipse are V -orthogonal at their points of intersection.*

Proof: Let us write

$$V = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}, \text{ and hence } V^{-1} = \Delta^{-1} \begin{pmatrix} \sigma_2^2 & -\rho\sigma_1\sigma_2 \\ -\rho\sigma_1\sigma_2 & \sigma_1^2 \end{pmatrix} \quad (3.16)$$

where $\Delta = \{\sigma_1^2\sigma_2^2(1 - \rho^2)\}$. Now, the ellipse $\mathbf{x}^T V^{-1} \mathbf{x} = \text{Constant}$ can be expressed as

$$\sigma_2^2 x_1^2 - 2\rho\sigma_1\sigma_2 x_1 x_2 + \sigma_1^2 x_2^2 = \text{Constant}. \quad (3.17)$$

By implicit differentiation, we have that

$$dx_2/dx_1 = (\sigma_2/\sigma_1)(\rho\sigma_1 x_2 - \sigma_2 x_1)/(\sigma_1 x_2 - \rho\sigma_2 x_1); \quad (3.18)$$

this is the slope of the tangent to the contour (3.16) at (x_1, x_2) . Now consider the line ℓ defined by $x_2 = mx_1$ ($m \neq 0$) (see Figure 3.5). For the time being, assume that $\sigma_1 m \neq \sigma_2 \rho$ to avoid division by zero; this can be relaxed by taking the limit $\sigma_1 m \rightarrow \sigma_2 \rho$.

Note that the slope, a , of the tangent to the contour (3.16) at the point of intersection with the line ℓ is given by

$$a = (\sigma_2/\sigma_1)(\rho\sigma_1 m - \sigma_2)/(\sigma_1 m - \rho\sigma_2). \quad (3.19)$$

Further, $(1, m)$, a vector parallel to the line ℓ , and $(1, a)$, a vector parallel to the tangent of the ellipse at the point of intersection with ℓ , are V -orthogonal. Therefore, the line ℓ and any contour (3.16) are V -orthogonal at their points of intersections. ■

To illustrate this further, let C_1 and C_2 be two contours defined by $\mathbf{x}^T V^{-1} \mathbf{x} = c_1^2$ and $\mathbf{x}^T V^{-1} \mathbf{x} = c_2^2$ respectively, for some c_1 and c_2 (see Fig. 3.5). Let O denote the origin and, as shown in Fig. 3.5, let F and Q be the points at which ℓ intersects C_1 and C_2 , respectively. Let T_Q be the tangent to C_2 at Q ; similarly, let T_F be the tangent to C_1 at F . Let P and E be points on T_Q and T_F respectively. Now, we have the following:

- For a given i , ($i = 1, 2$), every point on C_i is at V -equidistant from O .

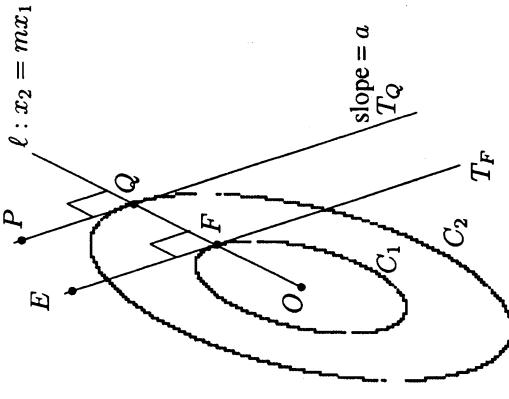


Fig. 3.5 *V*-projection onto a radial line ℓ through the center O of the ellipses C_1 and C_2 ; the point on ℓ that is V -closest to P is Q ; the point on T_F that is V -closest to O is F .

2. T_Q and T_F are parallel and their common slope a is given by (3.18).

3. T_Q and T_F are V -orthogonal to ℓ .

By Pythagoras theorem, the point on a line that is V -closest to a point, say P , is obtained by dropping a V -perpendicular to the line from P . Therefore, in Fig. 3.5 the point on ℓ that is V -closest to P is Q ; the point on T_Q that is V -closest to O is also Q . If $\mathbf{X} \sim N(\theta, V)$, the observed value of \mathbf{X} is at O , and θ is restricted to T_Q then the mle of θ is at Q ; if the observation is at P and θ is restricted to ℓ then the mle of θ is at Q .

The foregoing are also helpful to illustrate the form of the mle of θ when it is constrained to lie in a cone. Let $\mathbf{X} \sim N(\theta, V)$ and let us consider the mle of θ based on a single observation of \mathbf{X} when θ is constrained to lie on a convex cone. Let \mathcal{C} be the convex cone AOB in Fig. 3.6 and suppose that θ is constrained to lie in it.

Fig. 3.6 shows the mle when θ is restricted to \mathcal{C} and Fig. 3.7 shows the critical region for testing $H_0 : \theta = 0$ against $H_1 : \theta \in \mathcal{C}$. Let $OC \perp_V OA$ and $OD \perp_V OB$. Then, COD is the polar cone of \mathcal{C} with respect to $\langle \cdot \rangle_V$. Let Q and R be the points of intersection of an arbitrary contour $(\mathbf{x}^T V^{-1} \mathbf{x} = \text{constant})$ with OA and OB , respectively. Let PQ and SR be the tangents to the contour at Q and R , respectively.

Thus, $PQRS$ is a smooth curve with continuous slope everywhere. If $\mathbf{X} \in AOB$ then $\tilde{\theta} = \mathbf{X}$; if $\mathbf{X} \in COA$, say $\mathbf{X} = \overline{OP}$, then $\theta = \overline{OQ}$; if $\mathbf{X} \in DOC$ then $\theta = 0$; and if $\mathbf{X} \in DOB$, say $\mathbf{X} = \overline{OS}$, then $\theta = \overline{OR}$. Now, the boundary of a typical critical region is $PQRS$, where the QR is the segment of the contour C (ie. $\mathbf{x}^T V^{-1} \mathbf{x} = \text{constant}$) that lies in AOB .

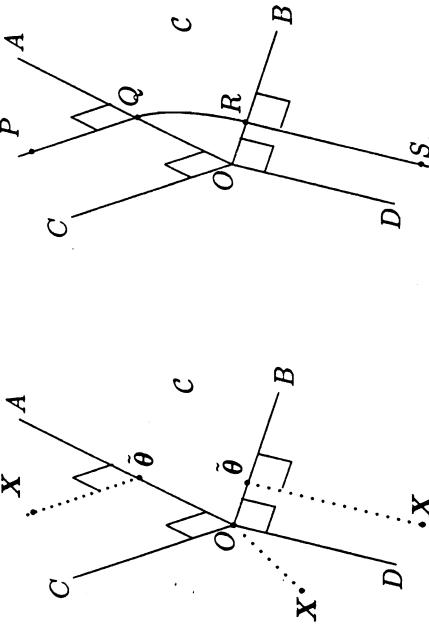


Fig. 3.6 The constrained mle of θ subject to $\theta \in C$ when $\mathbf{X} \sim N(\theta, \mathbf{V})$. The figure shows two plots. The top plot shows a coordinate system with axes X and \bar{X} . A point O is on the X -axis, and a point O' is on the \bar{X} -axis. A vector $\hat{\theta}$ starts from O' . A dotted line passes through O' and a point on the \bar{X} -axis. A solid line passes through O' and a point on the \bar{X} -axis. A shaded region C is bounded by the \bar{X} -axis and the solid line. The bottom plot shows a similar setup but with a different shaded region C , which is bounded by the \bar{X} -axis and a curve.

Example 3.3.5 A numerical example: likelihood ratio test against inequality constraints

Let \mathbf{X} be a bivariate random variable distributed as $N(\theta, \mathbf{V})$ where $\mathbf{V} = (1, \rho; \rho, 1)$ with $\rho = 0.9$. Now, consider the estimation of θ and test of $H_0 : \theta = 0$ against $H_1 : \theta \geq 0$ based on the following five observations on \mathbf{X} : $(-4, -3), (-4, -3), (-3, -2), (-2, -1)$, and $(-2, -1)$. For these data, the mean \bar{x} is $(-3, -2)$; as was pointed out in Section 3.2.1, the estimation and tests depend on the data only through this mean because \mathbf{V} is known. The kernel, $f(\theta)$, of the loglikelihood is $\{-0.5n(\bar{\mathbf{X}} - \theta)^T \mathbf{V}^{-1} (\bar{\mathbf{X}} - \theta)\}$. To estimate θ , we minimize $(\bar{\mathbf{x}} - \theta)^T \mathbf{V}^{-1} (\bar{\mathbf{x}} - \theta)$ over $\{\theta_1 \geq 0, \theta_2 \geq 0\}$; this achieves its minimum at $(0, 0.7)$. This may be verified by using the geometric method in the previous example. Alternatively, since the dimension of θ is only two, the minimum can be computed as explained in the next paragraph; while we would not use this method in high dimensions, it is instructive in two dimensions.

Note that $(\bar{\mathbf{x}} - \theta)^T \mathbf{V}^{-1} (\bar{\mathbf{x}} - \theta)$ is proportional to $f(\theta_1, \theta_2)$ defined by

$$f(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2 - 1.8\theta_1\theta_2 + 2.4\theta_1 - 1.4\theta_2 + 2.2.$$

This function is convex because its second derivative is proportional to \mathbf{V}^{-1} , which is positive definite. The minimum of $f(\theta)$, without any restrictions on (θ_1, θ_2) , is achieved at $\theta = \bar{\mathbf{x}} = (-3, -2)$, which is outside the region $\{\theta \geq 0\}$. Because $f(\theta)$ is convex, the minimum of $f(\theta)$ subject to $\{\theta \geq 0\}$ is achieved on the boundary of $\{\theta \geq 0\}$, which consists of the nonnegative θ_1 - and θ_2 -axes. Now, to find the absolute minimum of $f(\theta)$ over $\theta \geq 0$, it suffices to find the minima of $f(\theta)$ on the nonnegative segments of the θ_1 - and θ_2 -axes separately; the lowest of the two minima is the absolute minimum of $f(\theta)$ subject to $\theta \geq 0$.

The null distribution of LRT is given by (see (3.14) and (3.15)):

$$\text{pr}(LRT \geq c | H_0) = 0.5\text{pr}(\chi_1^2 \geq c) + 0.5(1 - \pi^{-1} \cos^{-1} \rho)\text{pr}(\chi_2^2 \geq c).$$

Substituting $\rho = 0.9$, the sample value of the likelihood ratio statistic, $n\hat{\theta}^T \mathbf{V}^{-1} \hat{\theta}$, is 12.9 for which the p -value is

$$0.5\text{pr}(\chi_1^2 \geq 12.9) + 0.5(1 - \pi^{-1} \cos^{-1} 0.9)\text{pr}(\chi_2^2 \geq 12.9),$$

which is less than 0.05. ■

Now, let us consider a curious phenomenon of the LRT based on the foregoing example.

Example 3.3.6 A curious example involving the LRT; the data are inconsistent with the the null and the alternative hypotheses (Silvapulle (1997a))

Consider an experiment to establish that a new treatment for pain relief is better than a placebo. For each patient, let X_1 and X_2 denote the reductions in pain at two locations (say neck and back) due to the treatment. For illustrative purposes, suppose that the distribution of \mathbf{X} and the data are the same as in the previous example. Suppose also that the new treatment is expected to be at least as good as the placebo in controlling pain. Thus, it is assumed *a priori* that $\theta \geq 0$. Suppose that the inference problem for establishing that the new treatment is better than the placebo be formulated as test of $H_0 : \theta = 0$ against $H_1 : \theta \geq 0$. Since every observed value is negative, the pain has increased under the new treatment for each patient in the sample; clearly, the new treatment has performed worse than the placebo. For these data, the mle of (θ_1, θ_2) under $H_1 : \theta \geq 0$ is $(0, 0.7)$ and the LRT of H_0 vs H_1 would reject H_0 in favor of H_1 at 5% level (LRT=12.9, p -value < 0.05).

At first, it would appear that there is something wrong because there is no way that we could claim that the new treatment is better than the placebo when, in fact, the pain under the new treatment turned out to be worse for each patient. In fact, there is nothing wrong as far as the statistical calculations are concerned. However, several issues arise in the interpretation of the results.

To provide an insight into why the mle of θ_2 is positive and the LRT rejects H_0 , first recall that the maximum likelihood estimator and the likelihood ratio statistic are based on the distribution of $\bar{\mathbf{X}}$ being $N(\theta, n^{-1}\mathbf{V})$. The contours of $N(\theta, n^{-1}\mathbf{V})$ passing through the observed value $\bar{\mathbf{x}} = (-3, -2)$ for $\theta = (0, 0)$ and $\theta = (0, 0.7)$ are shown in Fig. 3.8; these are indicated by C_2 and C_1 , respectively. Note that C_1 is smaller than C_2 . Therefore, if $g(\mathbf{x}|\theta)$ denotes the density function of $\bar{\mathbf{X}}$ then $g((-3, -2)|(0, 0.7)) > g((-3, -2)|(0, 0))$. Thus, the likelihood at $\theta = (0, 0.7)$ is larger than that at $\theta = (0, 0)$. Therefore, the data are more likely to have arisen from $N(\theta, \mathbf{V})$ with $\theta = (0, 0.7)$ than $\theta = 0$.

A small p -value for $H_0 : \theta = 0$ vs $H_1 : \theta \geq 0$ says that the data are considerably more inconsistent with $\theta = 0$ than with $\theta \geq 0$. If we accept the formulation of the problem as given above, then there are no contradictions; the LRT does exactly what we expect it to do.

property of LRT translates to it being not *concordant monotone*; for a discussion of this see Cohen et al. (2000). For further discussions and different points, see Cohen and Sackrowitz (2004), Chaudhuri and Perlman (2004), and Perlman and Wu (1999, 2002b).

It is possible to construct tests of $\theta = 0$ vs $\theta \geq 0$ in which points such as $(-3, -2)$ do not lie in the critical region. Cohen and Sackrowitz (1998) introduced a class of tests that they called the *Directed Tests* for testing this type of hypotheses. They developed their ideas on precise mathematical formulations of the problem; this will be introduced in Chapter 9. However, at the time of writing this book, there is a continuing debate concerning the advantages and disadvantages of various tests for constrained inference problems, in particular, as to whether or not LRT is deficient, and as to whether or not the directed tests are themselves deficient (see the references in the previous paragraph).



Fig. 3.8 The mle of θ subject to $\theta \geq 0$ when $\mathbf{X} \sim N(\theta, \mathbf{V})$, $\mathbf{V} = (1, 0.9 | 0.9, 1)$ and the sample mean is $(-3, -2)$.

Up to this point there are no difficulties in the statistical calculations. However, recently, there have been debates about the interpretation of these results. If the assumption $\mathbf{X} \sim N(\theta, \mathbf{V})$ where $\mathbf{V} = (1, 0.9 | 0.9, 1)$ is correct and $\bar{\mathbf{x}} = (-3, -2)^T$, should we accept that the new treatment is better? Some authors argue that LRT is deficient, while others argue there is nothing wrong with LRT. It is not possible to discuss all the relevant issues concerning this. Our objective here has been to draw attention to the phenomenon. We refer the readers to Silvapulle (1997a), Cohen et al. (2000), Cohen and Sackrowitz (2004), Perlman and Wu (2002b), and Chaudhuri and Perlman (2004). Although these authors discuss different examples, the statistical issues therein are the same as those in the foregoing example.

If we wish to establish that the new treatment is better than the placebo, then the testing problem could have been formulated as $H_0 : \theta \notin A$ vs $H_1 : \theta \in A$ where A is the set of all values of θ that correspond to the treatment being better than the placebo. For example, one possible choice is $A = \{\theta_1 > 0, \theta_2 > 0\}$; for this particular choice of A , if the test rejects H_0 , then clearly there is statistical evidence not only to reject H_0 but also to support H_1 . This particular testing problems will be discussed in later chapters under the topic, *sign testing*. More generally, in a test of H_a vs H_b for some H_a and H_b , if we wish to interpret the rejection of H_a as statistical evidence in support of H_b , it may be necessary to ensure that $H_a \cup H_b$ is the full parameter space. ■

The phenomenon in the simple example just discussed manifests in other settings as well; see Problem 3.9 for an example. In the analysis of ordinal response data, this

3.4 CHI-BAR-SQUARE DISTRIBUTION

The family of chi-bar-square distributions, to be introduced in this section, plays an important role when the null and/or alternative hypothesis involve parameter inequalities. The null distributions of likelihood ratio statistics for Type A and Type B problems with multivariate normal data turn out to be chi-bar-square. The role of chi-bar-square distributions in Type A and Type B testing problems is similar to that of chi-square distributions in the application of likelihood ratio for testing $\theta = 0$ against $\theta \neq 0$ in small and in large samples. In this section, we introduce the general form of chi-bar-square distribution and show how it relates to an important fundamental result on likelihood ratio test.

Let $\mathcal{C} \subset \mathbb{R}^p$ (\mathcal{C} is a closed convex cone) and let $\mathbf{Z}_{p \times 1} \sim N(0, \mathbf{V})$, where \mathbf{V} is a positive definite matrix. We define $\tilde{\chi}^2(\mathbf{V}, \mathcal{C})$ to be the random variable, which has the same distribution as $[\mathbf{Z}^T \mathbf{V}^{-1} \mathbf{Z} - \min_{\theta \in \mathcal{C}} (\mathbf{Z} - \theta)^T \mathbf{V}^{-1} (\mathbf{Z} - \theta)]$. If there is no possibility of any ambiguity, we shall also write

$$\tilde{\chi}^2(\mathbf{V}, \mathcal{C}) = \mathbf{Z}^T \mathbf{V}^{-1} \mathbf{Z} - \min_{\theta \in \mathcal{C}} (\mathbf{Z} - \theta)^T \mathbf{V}^{-1} (\mathbf{Z} - \theta). \quad (3.19)$$

It would be instructive to provide a geometric interpretation of $\tilde{\chi}^2$ and its related quantities. Fig. 3.9 provides a schematic diagram of \mathcal{C} and a typical value of \mathbf{Z} represented by OA . Let the inner product be defined as $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{V}} = \mathbf{x}^T \mathbf{V}^{-1} \mathbf{y}$, for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$; this induces the norm $\|\mathbf{x}\|_{\mathbf{V}} = \{\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x}\}^{1/2}$ and the distance $\|\mathbf{x} - \mathbf{y}\|_{\mathbf{V}} = \{(x - y)^T \mathbf{V}^{-1} (x - y)\}^{1/2}$ between \mathbf{x} and \mathbf{y} . Let B be the point in \mathcal{C} that is \mathbf{V} -closest to A , and let $\tilde{\mathbf{Z}}$ denote the vector OB . Therefore, $\tilde{\mathbf{Z}}$ is the value of \mathbf{x} at which

$$(\mathbf{Z} - \mathbf{x})^T \mathbf{V}^{-1} (\mathbf{Z} - \mathbf{x}), \quad \mathbf{x} \in \mathcal{C}$$

is the minimum. In other words, $\tilde{\mathbf{Z}}$ is the \mathbf{V} -projection of \mathbf{Z} onto \mathcal{C} ; let us denote it by $\Pi_{\mathbf{V}}(\mathbf{Z} | \mathcal{C})$ or equivalently $\Pi_{\mathbf{V}}(\mathbf{Z}, \mathcal{C})$. It is shown in the Appendix that $\mathbf{Z} - \tilde{\mathbf{Z}}$ is \mathbf{V} -orthogonal to \mathcal{C} ; in other words, OB is \mathbf{V} -orthogonal to AB . Therefore, by

It would also be useful to note that

$$\bar{\chi}^2(\mathbf{V}, \mathcal{C}^\circ) = \|OC\|_{\mathbf{V}}^2 = \|\mathbf{Z} - \mathcal{C}\|^2.$$

These results are summarized below.

Proposition 3.4.1 Let \mathbf{V} be a $p \times p$ positive definite matrix. Then

1. $\|\mathbf{Z}\|_{\mathbf{V}}^2 = \|\Pi_{\mathbf{V}}(\mathbf{Z} | \mathcal{C})\|_{\mathbf{V}}^2 + \|\mathbf{Z} - \Pi_{\mathbf{V}}(\mathbf{Z} | \mathcal{C})\|_{\mathbf{V}}^2.$
2. $\Pi_{\mathbf{V}}(\mathbf{Z} | \mathcal{C}^\circ) = \mathbf{Z} - \Pi_{\mathbf{V}}(\mathbf{Z} | \mathcal{C}).$
3. If $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{V})$ then

Fig. 3.9 \mathbf{V} -projection: OB and OC are the \mathbf{V} -projections of OA onto \mathcal{C} and \mathcal{C}° respectively.

applying the Pythagoras theorem to OAB , we have that

$$\|OA\|_{\mathbf{V}}^2 = \|OB\|_{\mathbf{V}}^2 + \|BA\|_{\mathbf{V}}^2,$$

and hence

$$\mathbf{Z}^T \mathbf{V}^{-1} \mathbf{Z} = \tilde{\mathbf{Z}}^T \mathbf{V}^{-1} \tilde{\mathbf{Z}} + \min_{\mathbf{x} \in \mathcal{C}} (\mathbf{Z} - \mathbf{x})^T \mathbf{V}^{-1} (\mathbf{Z} - \mathbf{x}).$$

Therefore,

$$\bar{\chi}^2(\mathbf{V}, \mathcal{C}) = \|OB\|_{\mathbf{V}}^2.$$

We denote the \mathbf{V} -distance between the point \mathbf{Z} and the set \mathcal{C} by $\|\mathbf{Z} - \mathcal{C}\|_{\mathbf{V}}$ and it is defined by

$$\|\mathbf{Z} - \mathcal{C}\|_{\mathbf{V}}^2 = \inf\{(\mathbf{Z} - \mathbf{x})^T \mathbf{V}^{-1} (\mathbf{Z} - \mathbf{x}) : \mathbf{x} \in \mathcal{C}\};$$

in Fig. 3.9, $\|\mathbf{Z} - \mathcal{C}\|_{\mathbf{V}}$ is equal to $\|AB\|_{\mathbf{V}}$.

The polar cone \mathcal{C}° of \mathcal{C} with respect to the inner product $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{V}} = \mathbf{x}^T \mathbf{V}^{-1} \mathbf{y}$ is defined by

$$\mathcal{C}^\circ = \{\mathbf{x} : \mathbf{x}^T \mathbf{V}^{-1} \mathbf{y} \leq 0 \text{ for every } \mathbf{y} \in \mathcal{C}\}.$$

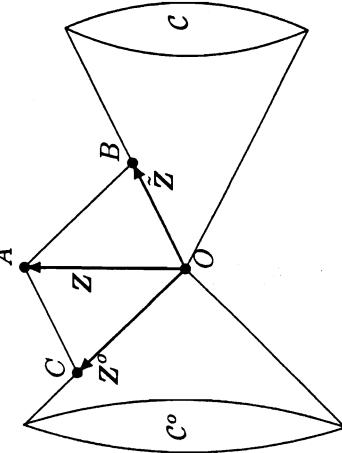
Because \mathcal{C} is a closed convex cone, it may be verified easily that \mathcal{C}° is also a closed convex cone. It consists of the vectors that form obtuse angles with every vector in \mathcal{C} ; here obtuse means that $\mathbf{x}^T \mathbf{V}^{-1} \mathbf{y} \leq 0$.

Let \mathbf{Z}° denote the point in \mathcal{C}° that is closest to \mathbf{Z} . In other words, \mathbf{Z}° is the projection of \mathbf{Z} onto \mathcal{C}° and is defined as the value of \mathbf{x} at which

$$(\mathbf{Z} - \mathbf{x})^T \mathbf{V}^{-1} (\mathbf{Z} - \mathbf{x}), \quad \mathbf{x} \in \mathcal{C}^\circ$$

is a minimum. Let C in \mathcal{C}° be the point such that $\mathbf{Z}^\circ = OC$. Again, the Pythagoras theorem is applicable to OAC . In fact, $OBAC$ is a rectangle, and hence

$$OC = BA, CA = OB, \|OA\|_{\mathbf{V}}^2 = \|OC\|_{\mathbf{V}}^2 + \|CA\|_{\mathbf{V}}^2 = \|OB\|_{\mathbf{V}}^2 + \|BA\|_{\mathbf{V}}^2.$$



Now we have the following fundamental result; proof is given in the Appendix.

Theorem 3.4.2 Let \mathcal{C} be a close convex cone in \mathbb{R}^p and \mathbf{V} be a $p \times p$ positive definite matrix. Then the distribution of $\bar{\chi}^2(\mathbf{V}, \mathcal{C})$ is given by

$$pr\{\bar{\chi}^2(\mathbf{V}, \mathcal{C}) \leq c\} = \sum_{i=0}^p w_i(p, \mathbf{V}, \mathcal{C}) pr(\chi_i^2 \leq c) \quad (3.20)$$

where $w_i(p, \mathbf{V}, \mathcal{C})$, $i = 0, \dots, p$ are some nonnegative numbers and $\sum_{i=0}^p w_i(p, \mathbf{V}, \mathcal{C}) = 1$. ■

More details about the quantities $w_i(p, \mathbf{V}, \mathcal{C})$ and their computation are discussed in the next section. Several proofs of Theorem 3.4.2 have appeared in the literature for progressively more general cases. The proofs in Gourieroux et al. (1982) and Shapiro (1985) are noteworthy. The proof in Gourieroux et al. (1982) is for the special case when $\mathcal{C} = \{\theta : R\theta \geq 0\}$ and R is a full row-rank matrix. This proof is instructive because it is a direct extension of the ideas presented earlier. The proof in Shapiro (1985) is for general closed convex cones and is the most general known; further, Shapiro's proof is based on results relating to projections on to polyhedrals.

The proof given in the Appendix to this chapter is based on the proof due to Shapiro (1985). The special case of the theorem for $p = 2$ was discussed earlier in Section 1.3. The main idea underlying the proof for the general case is similar although the technical details are more involved.

In the special case of two dimensions, we saw that the sum of the three weights is 1. Theorem 3.4.2 shows that this generalizes to the case in (3.20). Thus the right-hand side of (3.20) is a weighted mean of several tail probabilities of χ^2 -distributions and hence is known as a chi-bar-square distribution. We shall refer to $\{w_i(p, \mathbf{V}, \mathcal{C})\}$ as chi-bar-square weights or simply as weights. Another term used for these weights is level probabilities.

Now, suppose that we wish to test $H_0 : \theta = 0$ against $H_1 : \theta \in \mathcal{C}$ based on a single observation of the p -dimensional vector \mathbf{X} where $\mathbf{X} \sim N(\theta, \mathbf{V})$ and \mathbf{V} is a

known positive definite matrix. It follows from the definition, that the *LRT* takes the form

$$\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} - \min\{(\mathbf{X} - \boldsymbol{\theta})^T \mathbf{V}^{-1} (\mathbf{X} - \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathcal{C}\}, \quad (3.21)$$

which is the expression on the right-hand side of (3.19) with \mathbf{X} in place of \mathbf{Z} . Therefore, the distribution of *LRT* under $H_0 : \boldsymbol{\theta} = \mathbf{0}$ is $\bar{\chi}^2(V, C)$. In view of this, (3.21) is sometimes called a $\bar{\chi}^2$ -statistic. It is worth noting that the $\bar{\chi}^2$ -statistic is based on principles of generalized least squares, and therefore it is a reasonable test statistic even if the distribution of \mathbf{X} is not normal. The null distribution of the $\bar{\chi}^2$ -statistic in (3.21) is a $\bar{\chi}^2$ -distribution if $\mathbf{X} \sim N(\mathbf{0}, \mathbf{V})$. ■

3.5 COMPUTING THE TAIL PROBABILITIES OF CHI-BAR-SQUARE DISTRIBUTIONS

As was indicated earlier, the null distribution of several test statistics for or against inequality constraints turn out to be $\bar{\chi}^2$. Therefore, we need to be able to compute its tail probability to obtain the p-value and/or critical value. This would be easy if the chi-bar-square weights $\{w_i\}$ are known. Unfortunately, the exact computation of $\{w_i\}$ is quite difficult in general. The difficulties in computing $\{w_i\}$ may give the false impression that likelihood ratio tests of hypotheses for Type A and Type B testing problems are too complicated to implement. However, we can compute the tail probability of a chi-bar-square distribution by simulation; this is easy to use and fast. Therefore, for most practical needs, the fact that w_i has a complicated form does not cause any serious difficulties. Later, we will see situations where the chi-bar-square weights depend on unknown nuisance parameters; in those cases our inability to compute the weights accurately and fast may cause difficulties.

Simulation 1: To compute $\text{pr}\{\bar{\chi}^2(\mathbf{V}, C) \geq c\}$

(1) Generate \mathbf{Z} from $N(\mathbf{0}, \mathbf{V})$. (2) Compute $\bar{\chi}^2(\mathbf{V}, C)$ in (3.19). (3) Repeat the first two steps N times (say, $N = 10000$). (4) Estimate $\text{pr}\{\bar{\chi}^2(\mathbf{V}, C) \geq c\}$ by (n/N) where n is the number of times $\bar{\chi}^2(\mathbf{V}, C)$ in the second step turned to be greater than or equal to c . ■

If C involves only linear constraints, then a quadratic program can be used for computing $\bar{\chi}^2$ (see Section 2.1.2.2). We can also try to compute the tail probability of a $\bar{\chi}^2$ distribution by computing the weights first. However, the weights of a chi-bar-square distribution are likely to be required only for computing the tail probability; if this is the case, then it is likely to be just as easy to compute the tail probability directly by the foregoing simulation method.

Now, let us consider the case when

$$C = \{\boldsymbol{\theta} \in \mathbb{R}^p : \mathbf{a}_i^T \boldsymbol{\theta} \geq 0, i = 1, \dots, k\}. \quad (3.22)$$

If the weights $w_i(p, V, C)$ are of interest, then formula (3.29) given in the next section can be used for a simulation-based method for computing them. The simulation steps are easy to implement. It requires only a quadratic program. Such programs are

available in several software packages; a simple and easy to implement algorithm for solving quadratic programs is also given in Wollan and Dykstra (1987).

Simulation 2: To compute $w_i(p, V, C)$ when C is the polyhedral in (3.22)

(1) Generate \mathbf{Z} from $N(\mathbf{0}, \mathbf{V})$. (2) Compute $\tilde{\mathbf{Z}}$, the point at which $(\mathbf{Z} - \boldsymbol{\theta})^T \mathbf{V}^{-1} (\mathbf{Z} - \boldsymbol{\theta})$ is a minimum over $\boldsymbol{\theta} \in C$. (3) Let $J = \{j : \mathbf{a}_j^T \tilde{\mathbf{Z}} = 0\}$ and $\phi = \{\boldsymbol{\theta} : \mathbf{a}_j^T \boldsymbol{\theta} = 0 \text{ for every } j \in J\}$ and s be the dimension of ϕ ; compute s . (4) Repeat the previous steps N times (say $N = 10000$). (5) Estimate $w_i(p, V, C)$ by the proportion of times $\#$ turned out to be equal to i , ($i = 1, \dots, p$). ■

For more details and proofs relating to the foregoing simulation method see Silvapulle (1996) and Shapiro (1985).

The case when C is the nonnegative orthant \mathbb{R}^{+p} arises frequently in inequality constrained statistical inference. When the \mathbf{a}_j 's defining the polyhedral $C = \{x : \mathbf{a}_j^T x \geq 0\}$ are linearly independent, a result in the next section shows that $\{w_i(p, V, C)\}$ can be obtained from $\{w_i(k, W, \mathbb{R}^{+k})\}$ for some $k \times k$ positive definite matrix W . Thus, the case when the constrained set is the nonnegative orthant is of particular interest.

If the weights take the form $w_i(k, W, \mathbb{R}^{+k})$, then the foregoing Simulation 2 reduces to the following simpler form:

Simulation 3: To compute $w_i(p, W, \mathbb{R}^{+p})$, $i = 1, \dots, p$

(1) Generate \mathbf{Z} from $N(\mathbf{0}, W)$. (2) Compute $\tilde{\mathbf{Z}}$, the point at which $(\mathbf{Z} - \boldsymbol{\theta})^T \mathbf{W}^{-1} (\mathbf{Z} - \boldsymbol{\theta})$ is the minimum over $\boldsymbol{\theta} \geq 0$. (3) Count the number of positive components of $\tilde{\mathbf{Z}}$ (this is equal to s in Step 3 of Simulation 2). (4) Repeat the previous three steps N times (say $N = 10000$). (5) Estimate $w_i(p, W, \mathbb{R}^{+p})$ by the proportion of times $\#$ turned out to have exactly i positive components, $i = 1, \dots, p$. ■

This approach for the case when the polyhedral is \mathbb{R}^{+k} appears in Wolak (1987). The general idea of using simulation to compute the weights has been around for sometime, for example, see Perlman (1969) and Gourieroux et al. (1982). The authors have a FORTRAN program that implements the foregoing procedure using subroutines from IMSL; this works quite fast.

If $\text{pr}\{\bar{\chi}^2(V, C) \geq c\}$ is very small, then Simulation 1 would require a large number of simulations. In this case, we can overcome the problem by computing the weights of the chi-bar-square distribution separately, for example, by simulation. Once this is done, the tail probabilities of the chi-square distributions in $\sum w_i(p, V, C) \text{pr}(\chi_i^2 \geq c)$ may be computed using methods that are already available. The advantage of this method is that the value of c has no bearing on $w_i(p, V, C)$ and standard statistical software packages have programs for computing $\text{pr}(\chi_i^2 \geq c)$ very accurately for large values of c . ■

In most practical applications, C is a polyhedral. However, for completeness, let us also consider the case when C is not a polyhedral. In this general case, closed-form expressions for $w_i(p, V, C)$ are quite complicated. They can be expressed as some integrals (see Kuriki (1993), and Lin and Lindsay (1997)). At this stage, it is unclear whether or not such integral representations can be used for computing them. However, a simulation approach can be used for computing the weights even in this general case. ■

Simulation 4: To compute $w_i(p, V, C)$ when C is not necessarily a polyhedral

- (1) Let $0 < c_1 < \dots < c_{p+1} < \infty$ and $a_{ij} = \text{pr}(\chi_i^2 \leq c_{j+1})$ for $i, j = 0, \dots, p$.
- (2) Generate Z from $N(\mathbf{0}, V)$, and compute $\bar{\chi}^2(V, C)$ in (3.19). (3) Repeat the previous step a large number of times and let p_j be the proportion of times $\bar{\chi}^2(V, C)$ computed in the previous step fell in the interval $(0, c_{j+1}), j = 0, \dots, p$. (4) Estimate $w_j(p, V, C)$ by solving the system of $(p+1)$ equations $p_j = a_{0j}w_0 + \dots + a_{pj}w_p, (j = 0, \dots, p)$.

For most of the Type A and Type B testing problems, the tail probability P of the likelihood ratio statistic takes the form

$$P = \sum_{i=0}^p w_i(p, V, C) \text{pr}(\chi_{p+k}^2 \geq c),$$

for some k . Since $\text{pr}(\chi_i^2 \geq c) \leq \text{pr}(\chi_j^2 \geq c)$ for $i < j$, $\sum_0^p w_i(p, V, C) = 1$ and $0 \leq w_i(p, V, C) \leq (1/2)$ by part 4 in Proposition 3.6.1, we have the following lower and upper bounds for P :

$$(1/2)\{\text{pr}(\chi_k^2 \geq c) + \text{pr}(\chi_{k+1}^2 \geq c)\} \leq P \leq (1/2)\{\text{pr}(\chi_{p+k-1}^2 \geq c) + \text{pr}(\chi_p^2 \geq c)\}. \quad (3.23)$$

These bounds are easy to compute and may turn out to be adequate for practical purposes in some cases; for example, if the upper bound is small then there may not be any need to compute the exact p -value to reject the null hypothesis. It is worth pointing out that these may be crude bounds if p is large.

The proof of Theorem 3.4.2 shows that $\text{pr}\{\bar{\chi}^2(V, C) \geq c\}$ and $w_i(p, V, C)$ are continuous in (V, C) . Therefore, if C is close to \mathbb{R}^{+p} and V is close to I , then $w_i(p, V, C)$ is close to $w_i(p, I, \mathbb{R}^{+p})$, and an approximation to $\text{pr}(\bar{\chi}^2(V, C) \geq c)$ is $\sum_{i=0}^p 2^{-p}p!/i!(p-i)! \text{pr}(\chi_i^2 \geq c)$ because $w_i(p, I, \mathbb{R}^{+p}) = 2^{-p}p!/i!(p-i)!$ by part 2 of Proposition 3.6.1. It may be verified that this approximation lies between the lower and the upper bounds given in (3.23). This provides a quick and rough estimate of the tail probability $\text{pr}(\bar{\chi}^2(V, C) \geq c)$. The closer the (V, C) to (I, \mathbb{R}^{+p}) , the better the approximation; for some comparisons on this, see Piegorsch (1990).

The following closed-form expressions for $w_i(p, V, \mathbb{R}^{+p})$ are useful.

1. Let $p = 1$. Then

$$w_0(1, V) = w_1(1, V) = 0.5. \quad (3.24)$$

Therefore, $\text{pr}(\bar{\chi}^2 \geq t) = 0.5\text{pr}(\chi_1^2 \geq t)$ for $t > 0$. This is related to the familiar result where the p -value for a two-sided alternative concerning a scalar parameter is equal to twice the p -value for the corresponding one-sided alternative.

2. Let $p = 2$. Then (see Example 3.3.4)

$$\begin{aligned} w_0(2, V) &= (1/2)\pi^{-1} \cos^{-1}(\rho_{12}), & w_1(2, V) &= (1/2), \text{ and} \\ w_2(2, V) &= (1/2) - (1/2)\pi^{-1} \cos^{-1}(\rho_{12}) \end{aligned} \quad (3.25)$$

where ρ_{12} is the correlation coefficient $v_{12}\{v_{11}v_{22}\}^{-1/2}$.

3. Let $p = 3$. Then, (for example, see Kudo, 1963, pp 414 - 415)

$$\begin{aligned} w_3(3, V) &= (4\pi)^{-1}(2\pi - \cos^{-1}\rho_{12} - \cos^{-1}\rho_{13} - \cos^{-1}\rho_{23}), \\ w_2(3, V) &= (4\pi)^{-1}(3\pi - \cos^{-1}\rho_{12,3} - \cos^{-1}\rho_{13,2} - \cos^{-1}\rho_{23,1}), \\ w_1(3, V) &= (1/2) - w_3(3, V), \\ w_0(3, V) &= (1/2) - w_2(3, V), \end{aligned} \quad (3.26)$$

where ρ_{ij} is the correlation coefficient $v_{ij}\{v_iv_{ji}\}^{-1/2}$, and ρ_{ijk} is the partial correlation coefficient $(\rho_{ij} - \rho_{ik}\rho_{jk})/\{(1 - \rho_{ik}^2)(1 - \rho_{jk}^2)\}^{1/2}$. ■

4. Let $V = I$. Then, by part 2 of Proposition 3.6.1,

$$w_i(p, I) = 2^{-p}p!/i!(p-i)!) \text{ for every } p \text{ and } i. \quad (3.27)$$

When $p = 4$, all the weights can be expressed as functions of $w_4(4, V)$ and V (for example, see Wolak (1987)). However, numerical integration or Monte Carlo simulation needs to be used to compute $w_4(4, V)$. Therefore, these formulas may not be that helpful, because it may be just as easy to compute all the weights or the tail probability of the $\bar{\chi}^2$ distribution by simulation.

The simulation approaches discussed in this section require a computer program for computing $\min\{(\mathbf{Z} - \boldsymbol{\theta})^T \mathbf{V}^{-1}(\mathbf{Z} - \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathcal{C}\}$. If \mathcal{C} involves only linear inequalities, we can use a standard quadratic program (see Section 2.1.2.2). If \mathcal{C} involves nonlinear inequalities in $\boldsymbol{\theta}$, then a general-purpose nonlinear optimization program may be used; such programs are available in IMSL, MATLAB, and NAG. If \mathcal{C} involves nonlinear inequality constraints and minimization of $(\mathbf{Z} - \boldsymbol{\theta})^T \mathbf{V}^{-1}(\mathbf{Z} - \boldsymbol{\theta})$ subject to the constraint $\boldsymbol{\theta} \in \mathcal{C}$ is difficult to be incorporated in a simulation program, then another approach worth considering is to compute bounds for the critical values and/or p -values, for example, by the following method. Let \mathcal{P}_1 and \mathcal{P}_2 be two polyhedrals, such that $\mathcal{P}_1 \subset \mathcal{C} \subset \mathcal{P}_2$. Since $\bar{\chi}^2(\mathbf{V}, \mathcal{C}) \leq \bar{\chi}^2(\mathbf{V}, \mathcal{P}_1) \leq \bar{\chi}^2(\mathbf{V}, \mathcal{P}_2) \leq \bar{\chi}^2(\mathbf{V}, \mathcal{P}_2)$ we have that

$$\text{pr}\{\bar{\chi}^2(\mathbf{V}, \mathcal{P}_1) \geq c\} \leq \text{pr}\{\bar{\chi}^2(\mathbf{V}, \mathcal{C}) \geq c\} \leq \text{pr}\{\bar{\chi}^2(\mathbf{V}, \mathcal{P}_2) \geq c\}.$$

This approach is useful for obtaining lower and upper bounds for the tail probability of a chi-bar-square distribution; this method was used in Silvapulle (1992b).

••• RESULTS ON CHI-BAR-SQUARE WEIGHTS

Often, the constraints defining \mathcal{C} are linear and independent. In this case, $w_i(p, V, C)$ can be expressed as $w_i(k, W, \mathbb{R}^{+k})$ for some W and k . Therefore, to simplify notation, we denote $w_i(p, V, C)$ by $w_i(p, V)$ when C is the positive orthant:

$$w_i(p, V) = w_i(p, V, \mathbb{R}^{+p}). \quad (3.28)$$

The main results concerning $w_i(p, V, C)$ are stated in this section; the proofs are not given here but may be found in Kudo (1963), Perlman (1969), Wolak (1987b), Shapiro (1985, 1988), and Silvapulle (1996a).

Proposition 3.6.1 Let C be a closed convex cone in \mathbb{R}^p and V be a $p \times p$ nonsingular covariance matrix. Then we have the following:

1. Let $C = \{x : a_i^T x \geq 0, i = 1, \dots, k\}$ for some a_1, \dots, a_k , $J = \{j : a_j^T \Pi_V(Z|C) = 0\}$, and $\phi = \{\theta : a_j^T \theta = 0, \forall j \in J\}$. Then $w_i(p, V, C) = \text{pr}(\text{The linear space } \phi \text{ is of dimension } i)$. (3.29)
- This can be restated as follows: Suppose that C is a polyhedral, $Z \sim N(\mathbf{0}, V)$ and let $F(Z, V, C)$ denote the face of C such that $\Pi_V(Z|C)$ lies in the relative interior of $F(Z, V, C)$; for a definition of a face of a cone and relative interior of a face see the Appendix to this chapter (page 123). Then $w_i(p, V, C) = \text{pr}\{\text{dimension of the linear space spanned by } F(Z, V, C) \text{ is } i\}$.
- Let $Z \sim N(\mathbf{0}, V)$ and C be the nonnegative orthant. Then, $w_i(p, V) = \text{pr}\{\Pi_V(Z|C) \text{ has exactly } i \text{ positive components}\}$.
- $\sum_0^p (-1)^i w_i(p, V, C) = 0$.
- $0 \leq w_i(p, V, C) \leq 0.5$.
- Let C^o denote the polar cone, $\{x \in \mathbb{R}^p : x^T V^{-1} y \leq 0 \text{ for every } y \in C\}$, of C with respect to the inner product $\langle x, y \rangle = x^T V^{-1} y$. Then $w_i(p, V, C^o) = w_{p-i}(p, V, C)$.
- Let $C = \{\theta \in \mathbb{R}^p : R\theta \geq 0\}$ where R is a $p \times p$ nonsingular matrix. Then $\tilde{\chi}^2(V, C) = \tilde{\chi}^2(RV R^T, \mathbb{R}^{+p})$ and $w_i(p, V, C) = w_i(p, RV R^T)$.
- $w_i(p, V) = w_{p-i}(p, V^{-1})$.
- Let $C = \{\theta \in \mathbb{R}^p : R\theta \geq 0\}$ where R is a $k \times p$ matrix of rank $k (\leq p)$. Then $w_{p-k+i}(p, V, C) = \begin{cases} w_i(k, RV R^T) & \text{for } i = 0, \dots, k \\ 0 & \text{otherwise} \end{cases}$
- Let $C = \{\theta \in \mathbb{R}^p : R_1 \theta \geq 0, R_2 \theta = 0\}$ where R_1 is $s \times p$, R_2 is $t \times p$, $s + t \leq p$, $[R_1^T, R_2^T]$ is of full rank, and $A = R_1 V R_1^T - (R_1 V R_2^T)(R_2 V R_2^T)^{-1}(R_2 V R_1^T)$. Then $w_{p-s-t+j}(p, V, C) = \begin{cases} w_j(s, A) & \text{for } j = 0, \dots, s \\ 0 & \text{otherwise.} \end{cases}$
- Let R be a $r \times p$ matrix of rank r , R_1 be a $q \times p$ submatrix of R , $\mathcal{M} = \{\beta : R\beta = 0\}$, $C = \{\beta : R_1 \beta \geq 0\}$, and \mathcal{M}^\perp be the orthogonal complement of \mathcal{M} with respect to $\langle x, y \rangle_V = x^T V^{-1} y$. Then $w_{r-q+j}(p, V, C \cap \mathcal{M}^\perp) = \begin{cases} w_j(q, R_1 V R_1^T) & \text{for } j = 0, \dots, q \\ 0 & \text{otherwise.} \end{cases}$

11. Let C denote the correlation matrix corresponding to V . Then $\tilde{\chi}^2(V, C) = \tilde{\chi}^2(C, C)$ and $w_i(p, V, C) = w_i(p, C, C)$ for every i .

Proof: For (1) see Shapiro (1985); (2) follows from (1) but it has also been known for quite some time (for example, see Perlman (1969)). See (5.3), (5.4), (5.5), (5.8) and (5.10) in Shapiro (1988) for (6), (7), (8), (9), and (10) respectively. ■

Kudo (1963) showed that $w_i(p, V) = \sum p\{(\Lambda_N)^{-1}\} p(\Lambda_{M:N})$ where $M \subset \{1, \dots, p\}$, the sum extends over all M with exactly i elements, $N = \{1, \dots, p\} \setminus M$, Λ_M is the covariance matrix of the vector $\{Z_i : i \in M\}$, $\Lambda_{M:N}$ is the same under $Z_j = 0$ for $j \notin M$, and $p(A)$ is $\text{pr}(X \geq 0)$ when $X \sim N(\mathbf{0}, A)$. Thus, $w_i(p, V)$ can be expressed in terms of orthant probabilities. Schervish (1983) provided a program for computing probabilities of rectangles corresponding to multivariate normal; this can also be used for computing the orthant probabilities. Bohrer and Chow (1978) provided a computer program based on this result for computing $w_i(p, V)$, for $p \leq 10$. Another FORTRAN program is given in Sun (1988a). Kudo (1963) also discussed its use for computations. Formulas for the chi-bar square weights when V is a diagonal matrix are given in Section 2.4 of Robertson et al. (1988) for some special cases.

The following useful corollary to part 5 of the foregoing proposition follows from Theorem 3.4.2; it says that when C is replaced by its polar cone the weights appear in the reverse order.

Corollary 3.6.2 Let C and V be as in Proposition 3.6.1. Then $\text{pr}\{\tilde{\chi}^2(V, C^o) \leq c\} = \sum_{i=0}^p w_{p-i}(p, V, C) \text{pr}(\chi_i^2 \leq c)$. ■

3.7 LRT FOR TYPE A PROBLEMS: V IS KNOWN

Section 3.3 provided detailed discussions on the distribution of LRT in two dimensions, and the discussions therein showed how chi-bar-square distributions arise. The main result in Section 3.4 showed that the LRT for testing $\theta = 0$ against $\theta \in C$ based on a single observation of X , where $X \sim N(\theta, V)$, is the $\tilde{\chi}^2$ statistic

$$X^T V^{-1} X - \min_{\theta \in C} (X - \theta)^T V^{-1} (X - \theta)$$

and that its null distribution is the chi-bar-square distribution

$$\sum_{i=0}^p w_i(p, V, C) \text{pr}(\chi_i^2 \leq c).$$

A similar result holds even if the null parameter space, $\{\mathbf{0}\}$, is replaced by a linear space; this section provides a discussion of such results.

Let \mathcal{M} be a linear space contained in C . Suppose that we are interested to test

$$H_0 : \theta \in \mathcal{M} \text{ against } H_1 : \theta \in C$$

based on a single observation of \mathbf{X} , where θ is a location parameter of \mathbf{X} , for example, it could be the mean of \mathbf{X} . A least squares statistic for testing H_0 against H_1 is

$$L_A = \min\{q(\mathbf{a}) : \mathbf{a} \in \mathcal{M}\} - \min\{q(\mathbf{a}) : \mathbf{a} \in \mathcal{C}\}, \quad (3.30)$$

where

$$q(\mathbf{a}) = (\mathbf{X} - \mathbf{a})^T \mathbf{V}^{-1} (\mathbf{X} - \mathbf{a})$$

and \mathbf{V} is a known positive definite matrix. If $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{V})$ then L_A is the LRT for testing H_0 against H_1 . A general result on the distribution of LRT for a Type A testing problem is given in the next result.

Theorem 3.7.1 Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{V})$ where \mathbf{V} is a positive definite matrix. Then, the LRT for the Type A testing problem, $H_0 : \boldsymbol{\theta} \in \mathcal{M}$ against $H_1 : \boldsymbol{\theta} \in \mathcal{C}$, is similar¹ and its null distribution is given by

$$\text{pr}(LRT \leq c \mid H_0) = \sum_{i=0}^p w_i(p, \mathbf{V}, \mathcal{C} \cap \mathcal{M}^\perp) \text{pr}(\chi_i^2 \leq c) \quad (3.31)$$

where $\mathcal{M}^\perp = \{\mathbf{x} : \mathbf{x}^T \mathbf{V}^{-1} \mathbf{y} = 0 \text{ for every } \mathbf{y} \in \mathcal{M}\}$ is the orthogonal complement of \mathcal{M} with respect to the inner product $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{V}} = \mathbf{x}^T \mathbf{V}^{-1} \mathbf{y}$.

Proof: The proof of this theorem is based on deeper results about projections onto to convex cones. These are discussed in the Appendix to this chapter. Here, we shall provide an intuitive explanation of the main ideas after a short proof.

It follows from parts (h) and (i) of Proposition 3.12.6 in the Appendix that

$$\begin{aligned} LRT &= \min\{\|\mathbf{X} - \boldsymbol{\theta}\|_{\mathbf{V}}^2 : \boldsymbol{\theta} \in \mathcal{M}\} - \min\{\|\mathbf{X} - \boldsymbol{\theta}\|_{\mathbf{V}}^2 : \boldsymbol{\theta} \in \mathcal{C}\} \\ &= \|\mathbf{X}\|_{\mathbf{V}}^2 - \min\{\|\mathbf{X} - \boldsymbol{\theta}\|_{\mathbf{V}}^2 : \boldsymbol{\theta} \in \mathcal{C} \cap \mathcal{M}^\perp\}. \end{aligned}$$

The last expression on the right-hand side is also the LRT for testing $\boldsymbol{\theta} = \mathbf{0}$ against $\boldsymbol{\theta} \in \mathcal{C} \cap \mathcal{M}^\perp$. Therefore, the proof follows from Theorem 3.4.2. ■

Figure 3.10 provides an intuitive picture of the ideas underlying the foregoing proof in three-dimensions. Let O be the origin and OR be the θ_1 axis; θ_2 and θ_3 axes are not shown. Let \mathcal{M} be the linear space spanned by the line OR ; thus, \mathcal{M} is simply the θ_1 axis. The rectangle $OPQR$ spans a plane; similarly, the rectangle $OSTR$ spans another plane. Let \mathcal{C} be the convex cone generated by these two planes as shown in the figure; thus, the cone \mathcal{C} has a wedge shape with its spine being OR . Clearly, \mathcal{M}^\perp is the plane perpendicular to the θ_1 axis and passes through O . Therefore, \mathcal{M}^\perp is spanned by the θ_2 axis and the θ_3 axis, and $\mathcal{C} \cap \mathcal{M}^\perp$ is the cone SOP . Let OA represent the vector \mathbf{X} . Let B be the projection of A onto C ; more precisely, let OB be the projection of OA onto C . Suppose that B is an interior point of $OPQR$. Then B is also the projection of A onto the plane spanned by $OPQR$. Therefore,

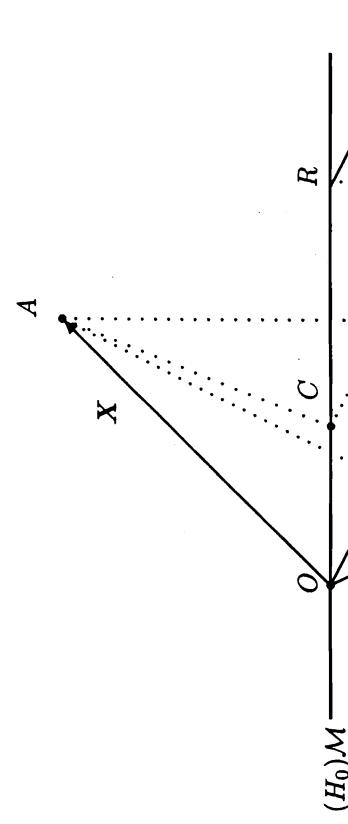


Fig. 3.10 Projections onto \mathcal{C} , \mathcal{M} , and $\mathcal{C} \cap \mathcal{M}^\perp$ where $\mathcal{M} \subset \mathcal{C}$.

AB is orthogonal to $OPQR$. Let D and C be the projections of A onto $\mathcal{C} \cap \mathcal{M}^\perp$ and \mathcal{M} respectively. Then D and C are also the projections of B onto $\mathcal{C} \cap \mathcal{M}^\perp$ and \mathcal{M} respectively. Now, we have the following :

$$AB \perp BD, AB \perp BC, AC \perp OC, AD \perp OD.$$

It follows that

$$\begin{aligned} \text{LRT of } \boldsymbol{\theta} \in \mathcal{M} \text{ vs } \boldsymbol{\theta} \in \mathcal{C} &= \min\{\|\mathbf{X} - \boldsymbol{\theta}\|_{\mathbf{V}}^2 : \boldsymbol{\theta} \in \mathcal{M}\} - \min\{\|\mathbf{X} - \boldsymbol{\theta}\|_{\mathbf{V}}^2 : \boldsymbol{\theta} \in \mathcal{C}\} \\ &= \|\mathbf{AC}\|^2 - \|\mathbf{AB}\|^2 = \|\mathbf{BC}\|^2 = \|\mathbf{DO}\|^2 = \|\mathbf{AO}\|^2 - \|\mathbf{AD}\|^2 \\ &= \|\mathbf{X}\|_{\mathbf{V}}^2 - \min\{\|\mathbf{X} - \boldsymbol{\theta}\|_{\mathbf{V}}^2 : \boldsymbol{\theta} \in \mathcal{C} \cap \mathcal{M}^\perp\}. \end{aligned}$$

This is also the LRT for testing $\boldsymbol{\theta} = \mathbf{0}$ against $\boldsymbol{\theta} \in \mathcal{C} \cap \mathcal{M}^\perp$. Since $\mathcal{C} \cap \mathcal{M}^\perp$ is a closed convex cone, its null distribution is a chi-bar-square and is given by (3.31).

In general, the weights of the chi-bar-square distribution in (3.31) depend on the parameter spaces \mathcal{M} and \mathcal{C} . There is no easy way to compute these weights for an arbitrary $\{\mathcal{C}, \mathcal{M}, \mathbf{V}\}$. The simulation procedure indicated in Section 3.5 is available as a general purpose procedure for computing $\text{pr}(LRT \leq c \mid H_0)$. Fortunately, often in practice, the constraints defining \mathcal{C} turn out to be linear and independent. In this case, we have the following simpler result, which follows from Proposition 3.6.1.

Corollary 3.7.2 Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{V})$ where \mathbf{V} is a positive definite matrix, R be a matrix of order $r \times p$, $\text{rank}(R) = r \leq p$, and let R_1 be a submatrix of R of order

¹Here, similar means that the null distribution of the test statistic is the same at every point in the null parameter space.

$q \times p$. Let the null and alternative hypothesis be $H_0 : R\theta = \mathbf{0}$ and $H_1 : R_1\theta \geq \mathbf{0}$ respectively. Then, the LRT is similar and its null distribution is given by

$$\text{pr}(LRT \leq c|H_0) = \sum_{i=0}^q w_i(q, R_1 V R_1^T) \text{pr}(\chi_{r-q+i}^2 \leq c)$$

Note that the chi-bar-square weights appearing in this corollary take the simpler form $w_i(q, R_1 V R_1^T)$, which by definition is equal to $w_i(q, R_1 V R_1^T, \mathbb{R}^{+q})$. Further, the number of terms in the foregoing $\tilde{\chi}^2$ distribution depends on the number of inequalities in H_1 only, not on the dimension p of θ . If the alternative hypothesis has only independent linear inequalities, and the number of them does not exceed 3, then we can use the explicit formulas for weights given in (3.24), (3.25), and (3.26). If the alternative hypothesis does not have any inequality constraints then $q = 0$ and hence we recover the classical result, $\text{pr}(LRT \leq c|H_0) = \text{pr}(\chi_r^2 \leq c)$.

If X is not normal but the distribution of $X - \theta$ does not depend on any unknown parameters, then L_A in (3.30) is still a suitable statistic for testing $H_0 : \theta \in \mathcal{M}$ against $H_1 : \theta \in C$; it is a generalized least squares based statistic. For this we have the following:

Proposition 3.7.3 Suppose that the distribution of $X - \theta$ does not depend on any unknown parameters. Then the null distribution of L_A in (3.30) does not depend on the particular value of θ in the null parameter space; ie. the test is similar.

Proof: It may be verified easily that $C + \alpha = C$ for any $\alpha \in \mathcal{M}$, where $C + \alpha$ is defined as $\{x + \alpha : x \in C\}$; and $\mathcal{M} + \alpha = \mathcal{M}$ for any $\alpha \in \mathcal{M}$. Let θ be an element of \mathcal{M} and let $Y = (X - \theta)$. Then $L_A = \min\{q^*(a) : a \in \mathcal{M}\} - \min\{q^*(a) : a \in C\}$, where $q^*(a) = (Y - a)^T V^{-1} (Y - a)$, and the distribution of Y does not depend on any unknown parameters. Therefore, if H_0 is true then the distribution of L_A is the same at every value in the null parameter space, \mathcal{M} . ■

In view of the foregoing Proposition, even if X is not normal, we can talk about the null distribution of L_A without specifying the null value of θ , and for the purposes of computing the p-value/critical value of L_A , we may assume that $\theta = \mathbf{0}$. With t denoting the observed value of L_A , we have that

$$\text{p-value} = \text{pr}(L_A \geq t | H_0) = \text{pr}(L_A \geq t | \theta = \mathbf{0}). \quad (3.32)$$

A consequence of this is that, even if X is not normal but the distribution of $X - \theta$ is known and does not depend on any unknown parameters, then exact finite sample p-value corresponding to L_A can be computed by simulation with pseudo-random observations generated at $\theta = \mathbf{0}$.

The foregoing results are useful in large samples even if the population distribution is not normal. Let us suppose that Z_1, \dots, Z_n are independently and identically distributed as Z , with $E(Z) = \theta$ and $\text{cov}(Z) = W$; it is assumed that the distribution of Z may not be normal, W is known, and that the distribution of $(Z - \theta)$ does not depend on any unknown parameters. Suppose that we are interested to test $H_0 : \theta \in \mathcal{M}$ against $H_1 : \theta \in C$. Then, with X denoting $n^{-1} \sum Z_i$, and assuming

that n is large, we have that X is approximately $N(\theta, V)$, where $V = n^{-1}W$. Therefore, L_A in (3.30) is still a suitable statistic for testing H_0 against H_1 , although it may not be the likelihood ratio statistic. Further, for large samples, its null distribution is approximately $\tilde{\chi}^2$. In particular, the results in Theorem 3.7.1 and Corollary 3.7.2 are applicable as large sample approximations. ■

Another large sample application of the foregoing may take the following form.

Let θ denote a parameter associated with a general statistical model, and let θ_0 denote its true value. For large samples, we typically have $\sqrt{n}(\hat{\theta} - \theta_0) \approx N(0, V)$. Assume that \hat{V} is a consistent estimator of V . Let

$$\hat{T} = \min_{\theta \in \mathcal{M}} n(\hat{\theta} - \theta)^T \hat{V}^{-1} (\hat{\theta} - \theta) - \min_{\theta \in C} n(\hat{\theta} - \theta)^T \hat{V}^{-1} (\hat{\theta} - \theta).$$

If the \hat{V} in the foregoing definition of \hat{T} is replaced by its probability limit, V , then the resulting change on \hat{T} is of order $o_p(1)$. Let $Z_n = \sqrt{n}(\hat{\theta} - \theta_0)$ and $Z \sim N(0, V)$. Suppose that $\theta_0 \in \mathcal{M}$ and $Z_n \xrightarrow{d} Z$. Then, we have that

$$\begin{aligned} \hat{T} &= \min_{\theta \in \mathcal{M}} (Z_n - \theta)^T V^{-1} (Z_n - \theta) - \min_{\theta \in C} (Z_n - \theta)^T V^{-1} (Z_n - \theta) + o_p(1) \\ &\xrightarrow{d} \min_{\theta \in \mathcal{M}} (Z - \theta)^T V^{-1} (Z - \theta) - \min_{\theta \in C} (Z - \theta)^T V^{-1} (Z - \theta). \end{aligned} \quad (3.33)$$

The distribution of this limit is given by (3.31); if the constraints in H_1 are linear and independent then the limiting distribution is given by Corollary 3.7.2. A difficulty in applying this general result will arise if the weights, $w_i(p, V, \mathcal{M}^\perp \cap C)$ in (3.31) depend on the unknown parameters defining V . Methods of dealing with such nuisance parameters will be discussed in more detail in the next chapter. For now, we note that if V depends on unknown nuisance parameters then a standard procedure in statistical inference is to define the p-value as the supremum of the tail probability over the nuisance parameter. ■

$$\text{p-value} = \sup_V \text{pr}(\hat{T} \geq t_{obs} | H_0) = \sup_V \text{pr}\{\tilde{\chi}^2(V, C \cap \mathcal{M}^\perp) \geq t_{obs}\}.$$

Example 3.7.1 Effect of in-breeding on growth (Kudo (1963))

In a large scale survey of child health carried out in Hiroshima and Nagasaki (see Neel and Schull (1954)), approximately 10000 children born of related as well as unrelated parents were studied. One facet of the study involved obtaining four anthropometric measurements as indices of growth and development. These were weight (y_1 gm), height (y_2 mm), head girth (y_3 mm), and chest girth (y_4 mm). The multivariate linear regression model, $E(Y) = \mu + t\alpha - c\beta$ was estimated where $Y = (Y_1, Y_2, Y_3, Y_4)^T$, t is the age in months at examination, c is the in-breeding

coefficient (see Neel and Schull (1954), p71) of the individual, and $\{\mu, \alpha, \beta\}$ are unknown parameters. This analysis does not allow for difference between the two cities, Hiroshima and Nagasaki; preliminary analysis failed to detect any significant difference between the two cities. For now, we assume that there is no difference between the two cities; the method of analysis given below would not be affected even if there is a difference. The null and alternative hypotheses of interest may be stated as

$$H_0 : \text{In-breeding has no effect on growth,}$$

$$\text{and } H_1 : \text{In-breeding has a negative effect on growth.}$$

respectively. The foregoing statement of the testing problem is broad, and it needs to be formulated in terms of some parameters so that it can be tested statistically. We formulate the problem as

$$H_0 : (\beta_1, \beta_2, \beta_3, \beta_4) = (0, 0, 0, 0) \text{ and } H_1 : (\beta_1, \beta_2, \beta_3, \beta_4) \geq (0, 0, 0, 0). \quad (3.34)$$

The null hypothesis will be rejected if there is evidence that at least one β_i is positive. The least squares estimate of β is $\hat{\beta} = (2.710, 0.775, 0.068, 0.329)^T$; the upper triangle of an estimate of the covariance matrix of $\hat{\beta}$ is $[1.199, 0.2164, 0.0839, 0.1841 | 0.0850, 0.0179, 0.0338 | 0.0232, 0.0137 | 0.0532]$, where “r” is used to separate different rows. The estimated covariance matrix is based on 10483 degrees of freedom. Now, since $\hat{\beta}$ is approximately $N(\beta, V)$, we wish to apply Corollary 3.7.2 with X replaced by $\hat{\beta}$. Note that we have only an estimate of the covariance matrix of $\hat{\beta}$. However, because the estimate is based on very large number of observations (more than 10000 degrees of freedom), it is reasonable to apply the results of the earlier sections with the estimated covariance matrix of $\hat{\beta}$ being treated as the exact covariance matrix of $\hat{\beta}$; essentially we are applying (3.33). Since each component of $\hat{\beta}$ is positive, the estimate of β based on the single observation of $\hat{\beta}$ subject to $\beta \geq 0$ is also $\hat{\beta}$. The sample value of the test statistic is

$$L_A = \|\hat{\beta}\|_V^2 - \min\{(\hat{\beta} - \beta)^T V^{-1}(\hat{\beta} - \beta) : \beta \geq 0\} = \|\hat{\beta}\|_V^2 = 9.3088.$$

Therefore, the p -value, $\text{pr}(L_A \geq 9.3088 | \beta = 0)$, is equal to

$$w_0(4, V) \text{pr}(\chi_0^2 \geq 9.3088) + \dots + w_4(4, V) \text{pr}(\chi_4^2 \geq 9.3088).$$

Now, we compute the weights $w_0(4, V), \dots, w_4(4, V)$ by simulation (see page 79). The computed values are $w_0 = 0.0062, w_1 = 0.0682, w_2 = 0.2700, w_3 = 0.4355$, and $w_4 = 0.2201$; a FORTRAN computer program took a few seconds to compute these weights by simulation using 50000 samples. Now,

$$\begin{aligned} p\text{-value} &= 0.0682 \times \text{pr}(\chi_1^2 \geq 9.3088) + 0.27 \times \text{pr}(\chi_2^2 \geq 9.3088) \\ &\quad + 0.4355 \times \text{pr}(\chi_3^2 \geq 9.3088) + 0.2201 \times \text{pr}(\chi_4^2 \geq 9.3088) = 0.026. \end{aligned}$$

Therefore, it appears that there is sufficient evidence to accept that in-breeding has a negative effect on growth. Let us note that the p -value could have been computed using Simulation 1 without computing the weights first.

Another approach to formulating and testing the foregoing hypothesis is to carry out the multiple hypotheses tests of $\beta_1 = 0$ vs $\beta_1 > 0$, $\beta_2 = 0$ vs $\beta_2 > 0$, $\beta_3 = 0$ vs $\beta_3 > 0$, and $\beta_4 = 0$ vs $\beta_4 > 0$. There are advantages and disadvantages of this approach. The main advantage is that the multiple tests attempt to identify which of β_1, \dots, β_4 are positive, and hence which aspects of growth are affected and which are not by in-breeding. However, a disadvantage of the multiple testing formulation is that they require an adjustment to the individual levels of significance, such as Bonferroni-type adjustment, which is known to result in low power for the combined tests. ■

Example 3.7.2 Meta analysis: Effect of sodium on blood pressure (Follmann (1996b)).

A large number of studies have been carried out to evaluate the effect of sodium reduction in reducing blood pressure. Let $X^T = (X_1, X_2) = (\text{systolic blood pressure, diastolic blood pressure})$. The extent of evidence in favor of the positive effect of sodium reduction in reducing the blood pressure vary considerably across various studies. What is desired is an overall probability statement that summarizes the strength of evidence in favor of the hypothesis that sodium reduction reduces blood pressure. To this end, Cutler et al. (1991) collected summary information available in a large number of published studies. It is worth pointing out that meta analysis requires considerable care in selecting the sample and ensuring that the observations are independent and that they are measuring the same parameter; we shall not consider such issues here. The relevant data are given in Table 1.13 (Cutler et al. (1991)).

Let us consider the data for the normotensive subjects. There are $n = 7$ observations. Let $\theta = (\theta_1, \theta_2)^T$ denote the mean reduction in blood pressure. Let us formulate the null and alternative hypotheses as $H_0 : \theta = 0$ and $H_1 : \theta \geq 0$, respectively. This is a Type A testing problem with two inequality constraints in H_1 . Because the sample size in each of the seven studies with normotensive subjects is reasonably large, let us assume that the estimated standard error of $\hat{\theta}_{ij}$ is its true standard deviation. From iith published study we obtain an estimate $\hat{\theta}_i$ and an estimated standard error for each component of $\hat{\theta}_i$. Because these standard errors are not the same across different studies, we apply weighted least squares to estimate θ ; each $\hat{\theta}_{ij}$ is weighted by the inverse of its variance. This leads to the following estimates:

$$\hat{\theta} = (0.17, 1.04)^T; \hat{V} = (0.1524, 0.0974 | 0.0974, 0.1510).$$

Now, the projection of $\hat{\theta}$, with respect to $\|\cdot\|_{\hat{V}}$, onto H_0 is 0 and that onto the positive orthant is $\hat{\theta}$ itself because both components of $\hat{\theta}$ are positive. Direct substitution provides $\hat{\theta}^T \hat{V}^{-1} \hat{\theta} = 9.97$. Let ρ denote the correlation coefficient corresponding to V ; thus it is the correlation coefficient between systolic and diastolic blood pressures. If ρ were known, a large sample p -value would be

$$[0.5\text{pr}\{\chi_1^2 \geq 9.97\} + (2\pi)^{-1} \cos^{-1}(\rho) \text{pr}\{\chi_2^2 \geq 9.97\}].$$

Since ρ is unknown,

$$p\text{-value} \simeq \sup_{\rho} [0.5\text{pr}\{\chi_1^2 \geq 9.97\} + (2\pi)^{-1} \cos^{-1}(\rho) \text{pr}\{\chi_2^2 \geq 9.97\}].$$

Let us make the reasonable assumption that the reductions in systolic and diastolic blood pressures are nonnegatively correlated; therefore, $\cos^{-1} \rho \leq \pi/2$. Hence

$$\sup_{\rho > 0} [(1/2) \text{pr}\{\chi_1^2 \geq 9.97\} + \{(2\pi)^{-1} \cos^{-1} \rho\} \text{pr}\{\chi_1^2 \geq 9.97\}] \leq 0.001.$$

Therefore, a large sample estimate of the *p*-value does not exceed 0.001; there is sufficient evidence to reject H_0 .

The results presented so far for testing $H_0 : \theta = 0$ against $H_1 : \theta \geq 0$ say that there is sufficient statistical evidence to reject $\theta = 0$ but this alone does not mean that there is sufficient statistical evidence to accept $\theta \geq 0$. Suppose that we are prepared to assume *a priori* that a reduction in sodium intake cannot increase mean systolic or diastolic blood pressure, which translates to $\theta \geq 0$. Now, since we have already noted that there is sufficient statistical evidence to reject $\theta = 0$ we can conclude that there is sufficient statistical evidence to accept $\theta \geq 0$.

Now, let us suppose that we do not wish to assume *a priori* that $\theta \geq 0$ but would like to establish that $\theta \geq 0$. This type of problems will be studied in a later chapter under the topic *sign testing*; see also *combination therapy*. It will be seen that we can formulate the inference problem as test of

$$H_0 : \theta_1 \leq 0 \text{ or } \theta_2 \leq 0 \quad \text{vs} \quad H_1 : \theta_1 > 0 \text{ and } \theta_2 > 0.$$

Once the problem is formulated in this form, the so-called *min test* can be applied. In this formulation if the *p*-value turns out to be small then we can conclude that there is sufficient statistical evidence not only to reject $H_0 : \theta_1 \leq 0$ or $\theta_2 \leq 0$ but also to accept $H_1 : \theta_1 > 0$ and $\theta_2 > 0$. ■

3.8 LRT FOR TYPE B PROBLEMS: V IS KNOWN

Let $X \sim N(\theta, V)$, where V is a given positive definite matrix. We consider the Type B problem of testing

$$H_1 : \theta \in C \text{ against } H_2 : \theta \notin C. \quad (3.35)$$

In contrast to Type A problems, a main feature of Type B problems is that the null hypothesis involves inequalities. In this case, some important issues arise. To illustrate some of these, we first consider a simple bivariate case.

Let $X = (X_1, X_2)^T$ be a bivariate random variable with density function $f(x - \theta)$ for some f and θ , and let the null and alternative hypotheses be

$$H_1 : (\theta_1, \theta_2) \geq (0, 0) \quad \text{and} \quad H_2 : (\theta_1, \theta_2) \not\geq (0, 0),$$

respectively. This is a special case of (3.35) with $C = \{\theta : \theta \geq 0\}$. Let

$$L_B = \min\{\|X - \theta\|^2 : \theta \geq 0\} - \min\{\|X - \theta\|^2 : \theta \in \mathbb{R}^2\}, \quad (3.36)$$

where $\|x\|^2 = x^T x$. Let $\tilde{\theta}$ be the projection of X onto the positive orthant. Then

$$L_B = \|X - \tilde{\theta}\|^2.$$

If $X \sim N(\theta, I)$ then $\tilde{\theta}$ is the *mle* of θ and L_B is the *LRT*; however, for now we shall not assume that X is $N(\theta, I)$. With ℓ_B denoting the sample value of L_B , one is tempted to define the *p*-value as $\text{pr}(L_B \geq \ell_B \mid H_0)$. A difficulty that arises with this is that this probability depends on the particular null value, θ , which may be anywhere in the null parameter space $\{\theta : \theta \geq 0\}$. Thus, $\text{pr}_\theta(L_B \geq \ell_B \mid \theta \in H_0)$ is not just a fixed number on the null parameter space, but a function of θ , and hence does not define a *p*-value. In this case, a reasonable approach to overcome the difficulty appears to be not to reject null hypothesis if there is at least one value in the null parameter space with which the data are consistent; or equivalently reject H_0 if the data are inconsistent with every value of θ in the null parameter space. In fact, the unusual procedure is to define the *p*-value as $\sup_{\theta \in H_0} \text{pr}_\theta(L_B \geq \ell_B)$. At first, it may seem to be a formidable task to evaluate this supremum. Fortunately, it is reached at $\theta = 0$, and therefore

$$\text{p-value} = \sup_{\theta \in H_0} \text{pr}_\theta(L_B \geq \ell_B) = \text{pr}_0(L_B \geq \ell_B). \quad (3.37)$$

The proof of this is easy, and is given below in Theorem 3.8.1.

It is implicit in the statement (3.37) that the strength of evidence against H_0 and in favor of H_1 based on $L_B = \ell_B$ depends on the assumed true value of θ in the null parameter space and that it is *least* when $\theta = 0$. For this reason, $\theta = 0$ is called the *least favorable null value* of L_B (also, the *least favorable null configuration of L_B*), and the distribution of L_B at this value is called the *least favorable null distribution of L_B* . The foregoing result holds in higher dimensions as well for Type B testing problems even if the distribution of X is not normal; this is stated below.

Theorem 3.8.1 Assume that the distribution of $X - \theta$ does not depend on any unknown parameters. Let the null and alternative hypotheses be $H_1 : \theta \in C$ and $H_2 : \theta \notin C$, respectively, where C is a closed convex cone. Let V be a given positive definite matrix, and let $L_B = \min\{\|X - \alpha\|^2 V^{-1} (X - \alpha) : \alpha \in C\}$. Then, $\text{pr}_\theta(L_B \geq c \mid \theta \in C) \leq \text{pr}_0(L_B \geq c)$, and therefore the least favorable null value of L_B is 0 and

$$\text{p-value} = \sup_{\theta \in C} \text{pr}_\theta(L_B \geq \ell_B) = \text{pr}_0(L_B \geq \ell_B) \quad (3.38)$$

where ℓ_B is the sample value of L_B . If C contains a linear space M then every point in M is also a least favorable null value for L_B .

Proof : The proof follows a set containment argument. The following proof is illustrated in Fig. 3.6 for the special case, $p = 2$, $V = I$ and $C = \mathbb{R}^{+2}$. Suppose that $\theta \in C$ and let $Z = X - \theta$. Then, for $c > 0$,

$$\begin{aligned} \text{pr}_\theta(L_B \geq c) &= \text{pr}_\theta[\|Z - C\|_V^2 \geq c] \\ &= \text{pr}[\|(Z + \theta) - C\|_V^2 \geq c] \\ &\leq \text{pr}[\|Z - C\|_V^2 \geq c]; \end{aligned}$$

The last step follows because $C \subset C - \theta$ and hence the distance from Z to C is at least as large as that to $C - \theta$. Therefore, we have that $\text{pr}_\theta(L_B \geq c \mid \theta \in C) \leq \text{pr}_0(L_B \geq c)$. If $\theta \in M$ then $C = C - \theta$ and hence the last part of the theorem follows. ■

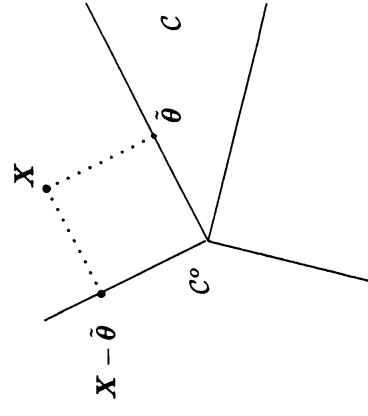


Fig. 3.13 Likelihood ratio test of $\theta \in C$ against $\theta \notin C$.

that the probability of X falling in the critical region and $\|Z\| \leq d$ decreases as θ moves along any straight line from the origin into the nonnegative orthant. This is true for any $d > 0$ and therefore the probability of X falling in the critical region is a maximum when $\theta = 0$, and hence the least favorable null value of L_B is 0.

Least favorable null distribution of LRT when $p = 2$

Now, let us derive the least favorable null distribution of the LRT in two dimensions. Some important results do not become clear when C is the nonnegative orthant, \mathbb{R}^+ . Therefore, we consider an arbitrary C . Let $\mathbf{X} \sim N(\boldsymbol{\theta}, I)$, $H_1 : \boldsymbol{\theta} \in C$ and $H_2 : \boldsymbol{\theta} \notin C$; this set-up is the same as that in Example 3.3.3 in Section 3.3 and the details therein would be helpful here as well. For testing H_1 vs H_2 , we have that

$$LRT = \min\{\|\mathbf{X} - \boldsymbol{\theta}\|^2 : \boldsymbol{\theta} \in C\} - \min\{\|\mathbf{X} - \tilde{\boldsymbol{\theta}}\|^2 : \boldsymbol{\theta} \in \mathbb{R}^2\} = \|\mathbf{X} - \tilde{\boldsymbol{\theta}}\|^2,$$

where $\tilde{\boldsymbol{\theta}}$ is the *MLE* of $\boldsymbol{\theta}$ subject to $\boldsymbol{\theta} \in C$. Note that $(\mathbf{X} - \tilde{\boldsymbol{\theta}})$ is also the point in C° that is closest to \mathbf{X} . Fig. 3.13 shows this for the present two dimensional case; this also holds in \mathbb{R}^p (see Proposition 3.12.4). Since C° is a closed convex cone in \mathbb{R}^2 , $\|\mathbf{X} - \tilde{\boldsymbol{\theta}}\|^2$ is the LRT for testing

$$H_0^* : \boldsymbol{\theta} = 0 \text{ against } H_1^* : \boldsymbol{\theta} \in C^\circ.$$

We just noted that $\|\mathbf{X} - \tilde{\boldsymbol{\theta}}\|^2$ is also the LRT for testing $H_1 : \boldsymbol{\theta} \in C$ against $H_2 : \boldsymbol{\theta} \notin C$ based on a single observation of \mathbf{X} . Therefore, the least favorable null distribution of the LRT for testing $H_1 : \boldsymbol{\theta} \in C$ against $H_2 : \boldsymbol{\theta} \notin C$, and the null distribution of the LRT for testing $H_0^* : \boldsymbol{\theta} = 0$ against $H_1^* : \boldsymbol{\theta} \in C^\circ$ are the same.

Since the angles (in radians) of C and C° at their vertices sum to π , we conclude from Example 3.3.3 in Section 3.3 that

$$\begin{aligned} \text{pr}(LRT \geq c | \boldsymbol{\theta} \in C) &\leq \text{pr}(LRT \geq c | \boldsymbol{\theta} = 0) \\ &= (0.5 - q)\text{pr}(\chi_0^2 \geq c) + 0.5\text{pr}(\chi_1^2 \geq c) + \text{qpr}(\chi_2^2 \geq c) \end{aligned} \quad (3.39)$$

where $q = \gamma/(2\pi)$, and γ is the angle (in radians) of the cone C° at the vertex.

Note that for this Type B testing problem, the weights in the relevant chi-bar-square

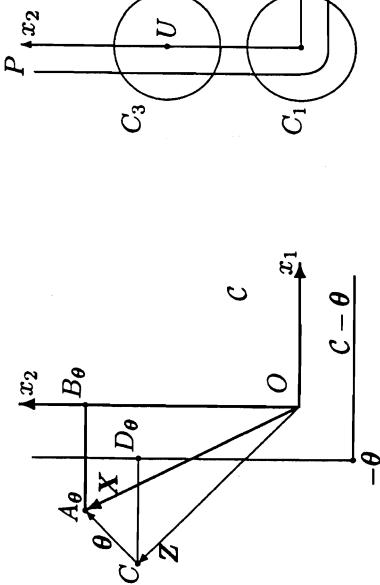


Fig. 3.11 Least favorable null value for $H_1 : \boldsymbol{\theta} \geq 0$ vs $H_2 : \boldsymbol{\theta} > 0$ is the origin

Least favorable null value when $p = 2$

The foregoing proof is illustrated in Fig. 3.11 for the special case $p = 2$, $V = I$ and $C = \mathbb{R}^{+2}$. Let Z be a random variable and let it be represented by OC . Let us think of \mathbf{X} as the sum of the deterministic part $\boldsymbol{\theta}$ in \mathbb{R}^{+2} and the random part Z that is unrelated to $\boldsymbol{\theta}$. Thus, the position of the point C is not related to $\boldsymbol{\theta}$. The vector X is represented by OA_θ where the position of A_θ depends on $\boldsymbol{\theta}$. Let B_θ and D_θ be the points in \mathbb{R}^{+2} and $\mathbb{R}^{+2} - \boldsymbol{\theta}$ that are closest to A_θ and C , respectively. Then $L_B = \|A_\theta B_\theta\|^2 = \|CD_\theta\|^2$, because the position of C is unrelated $\boldsymbol{\theta}$, it follows that $\|CD_\theta\|^2$ is a maximum when $\boldsymbol{\theta} = 0$. Therefore, $\boldsymbol{\theta} = 0$ is the least favorable null value.

It would be instructive to provide a different geometric interpretation of the foregoing proof. To this end, let us rewrite the proof of the previous theorem slightly differently. Let c and d be positive numbers. Suppose that $\boldsymbol{\theta} \in C$ and let $Z = \mathbf{X} - \boldsymbol{\theta}$. Now, the critical region is of the form $\mathcal{A} = \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x} - C\| \geq c\}$. Again, $C \subset C - \boldsymbol{\theta}$ for any $\boldsymbol{\theta} \in C$, and we have that

$$\begin{aligned} \text{pr}_{\boldsymbol{\theta}}(\|\mathbf{Z} - C\|_V \geq c \text{ and } \|\mathbf{Z}\|_V \leq d) &= \text{pr}(\|\mathbf{Z} - (C - \boldsymbol{\theta})\|_V \geq c \text{ and } \|\mathbf{Z}\|_V \leq d) \\ &\leq \text{pr}(\|\mathbf{Z} - C\|_V \geq c \text{ and } \|\mathbf{Z}\|_V \leq d) = \text{pr}_0(\|\mathbf{X} - C\|_V \geq c \text{ and } \|\mathbf{Z}\|_V \leq d). \end{aligned}$$

Now, the proof follows by taking the limit $d \rightarrow \infty$.

This provides a different geometric interpretation of the main idea that underlies Theorem 3.8.1. Let us illustrate this using the same simple bivariate example. Since the test statistic is $L_B = \|\mathbf{X} - \mathbb{R}^{+2}\|^2$, the critical region is of the form $\mathcal{A} = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x} - \mathbb{R}^{+2}\| \geq c\}$ for some $c > 0$; this is the region to the left/bottom of the curve PQ in Fig. 3.12. Let C_1 be a circle of radius d with center at the origin. Let $\boldsymbol{\theta}$ be the point in C as shown in Figure 3.12. Let the circle C_2 be obtained by shifting the center of C_1 to $\boldsymbol{\theta}$. Note that the part of the interior of C_1 in the critical region is larger than that corresponding to C_2 , for any $\boldsymbol{\theta}$ in C . Now, it is easily seen

distribution appear in the reverse order to that of the corresponding Type A problem of testing $\theta = 0$ against $\theta \in C$ (see Example 3.3.3).

In the foregoing example, if we were to assume that the covariance matrix of X is V rather than I , then the details of Example 3.3.4 are relevant and (3.39) would still hold with q as in (3.15).

This example in two-dimensions illustrates some important results that hold in higher dimensions as well; this is stated in the next theorem.

Theorem 3.8.2 Let $X \sim N(\theta, V)$ where V is a $p \times p$ positive definite matrix, and let the null and alternative hypotheses be

$$H_1 : \theta \in C \quad \text{and} \quad H_2 : \theta \notin C,$$

respectively. Then the least favorable null value of LRT is $\theta = 0$. The least favorable null distribution of LRT is $\bar{\chi}^2(V, C^\circ)$ and

$$\Pr(LRT \leq c \mid \theta = 0) = \sum_{i=0}^p w_{p-i}(p, V, C) \Pr(\chi_i^2 \leq c).$$

Proof: It follows from Propositions 3.4.1 (page 77) and 3.12.4 (page 116) that

$$\begin{aligned} LRT &= \min\{(\mathbf{X} - \boldsymbol{\alpha})^T V^{-1}(\mathbf{X} - \boldsymbol{\alpha}) : \boldsymbol{\alpha} \in C\} = \|\Pi_V(\mathbf{X} \mid C^\circ)\|_V^2 \\ &= \mathbf{X}^T V^{-1} \mathbf{X} - \min\{(\mathbf{X} - \boldsymbol{\theta})^T V^{-1}(\mathbf{X} - \boldsymbol{\theta}) : \boldsymbol{\theta} \in C^\circ\} \\ &= \bar{\chi}^2(V, C^\circ). \end{aligned}$$

Now, the proof follows from Corollary 3.6.2 on page 83. ■

As in Type A testing problems, if the inequalities defining C are all linear and independent, then the weights of the $\bar{\chi}^2$ distribution in the previous theorem take the simpler forms corresponding to the nonnegative orthant. The following corollary can be deduced from Proposition 3.6.1 (page 82).

Corollary 3.8.3 Let $X \sim N(\theta, V)$. Let the null and alternative hypotheses be

$$H_1 : R_1 \boldsymbol{\theta} \geq 0, R_2 \boldsymbol{\theta} = 0 \quad \text{and} \quad H_2 : \boldsymbol{\theta} \text{ is not restricted},$$

where R_1 is $s \times m$, R_2 is $t \times m$, and the rank of $[R_1^T, R_2^T]$ is $(s+t)$. Then

$$LRT = \min\{(\mathbf{X} - \boldsymbol{\alpha})^T V^{-1}(\mathbf{X} - \boldsymbol{\alpha}) : R_1 \boldsymbol{\alpha} \geq 0, R_2 \boldsymbol{\alpha} = 0\},$$

and the least favorable null distribution of LRT is

$$\Pr(LRT \leq c \mid \boldsymbol{\theta} = 0) = \sum_0^s w_{s-i}(s, A) \Pr(\chi_{t+i}^2 \leq c),$$

where $A = R_1 V R_1^T - (R_1 V R_2^T)(R_2 V R_2^T)^{-1}(R_2 V R_1^T)$.

In some situations, it is possible that linear equality constraints may be present in both the null and the alternative hypotheses as well. For example, suppose that the testing problem is

$$H_1 : \boldsymbol{\theta} \in C \quad \text{vs} \quad H_2 : \boldsymbol{\theta} \in \mathcal{L}$$

where \mathcal{L} is a linear subspace of \mathbb{R}^p and $C \subset \mathcal{L}$. In this case, the testing problem can be reformulated in the form discussed in this section. To this end, let $k = \dim(\mathcal{L})$ and A be a matrix such that $\text{rank}(A) = k (\leq p)$ and $\mathcal{L} = \{A\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^k\}$. Then the hypotheses take the form $H_1 : \boldsymbol{\beta} \in \mathcal{P}$ and $H_2 : \boldsymbol{\beta} \notin \mathcal{P}$, where \mathcal{P} is the closed convex cone, $\{\boldsymbol{\beta} : A\boldsymbol{\beta} \in C\}$. Now, with L denoting $-2\log(\text{likelihood})$, we have the following for $\boldsymbol{\theta} \in \mathcal{L}$:

$$\begin{aligned} L &= (\mathbf{X} - \boldsymbol{\theta})^T V^{-1}(\mathbf{X} - \boldsymbol{\theta}) = (\mathbf{X} - A\boldsymbol{\beta})^T V^{-1}(\mathbf{X} - A\boldsymbol{\beta}) \\ &= (\mathbf{X} - A\hat{\boldsymbol{\beta}})^T V^{-1}(\mathbf{X} - A\hat{\boldsymbol{\beta}}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T V^{-1}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \end{aligned}$$

where $\hat{\boldsymbol{\beta}} = (A^T V^{-1} A)^{-1} A^T V^{-1} \mathbf{A}$. Therefore, the problem is equivalent to testing $H_1 : \boldsymbol{\beta} \in \mathcal{P}$ against $H_2 : \boldsymbol{\beta} \notin \mathcal{P}$ based on a single observation of $\hat{\boldsymbol{\beta}}$ where $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, W)$. This is of the same form as that considered in the foregoing theorem and the corollary.

Finally, let L_A denote the LRT for testing $\boldsymbol{\theta} = 0$ vs $\boldsymbol{\theta} \in C$ and L_B denote the LRT for testing $\boldsymbol{\theta} \in C$ vs $\boldsymbol{\theta} \notin C$ based on a single observation of \mathbf{X} from $N(\boldsymbol{\theta}, V)$. Then, by using arguments similar to the proof of Lemma 3.13.6 (page 129), it can be deduced that

$$\Pr(L_A \in A \text{ and } L_B \in B \mid \boldsymbol{\theta} = 0) = \sum_{i=0}^p w_i(p, V, C) \Pr(\chi_{p-i}^2 \in B).$$

This was obtained by Raubertas et al. (1986). There may be some scenarios in which this may be useful. For example, if one wishes to reject $\boldsymbol{\theta} = 0$ and accept $\boldsymbol{\theta} \in C \setminus \{0\}$ if L_A is large and L_B is small, this result may have some use. However, a sketch of the critical region would show that it has an unappealing shape when $V \neq I$.

3.9 TESTS ON THE LINEAR REGRESSION PARAMETER

In this section it will be shown that most of the normal theory results in the previous sections extend to the normal theory linear model as well. It will be seen that most of the results concerning the LRT on the comparison of normal means studied in the previous chapter are special cases of the results presented in this section. Let us consider the linear model,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{E} \quad (3.40)$$

where \mathbf{Y} is an $N \times 1$ vector of observations, \mathbf{X} is a known $N \times p$ matrix of rank p with $p < N$, $\boldsymbol{\theta}$ is a $p \times 1$ vector of unknown parameters and \mathbf{E} has mean 0 and

covariance matrix $\sigma^2 U$, U is known and σ is unknown. Assume that the distribution of E does not depend on any unknown parameters. Let

$$\hat{\theta} = (\mathbf{X}^T \mathbf{U}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U}^{-1} \mathbf{Y}, \quad W = (\mathbf{X}^T \mathbf{U}^{-1} \mathbf{X})^{-1} \\ Q(\theta) = (\mathbf{Y} - \mathbf{X}\theta)^T \mathbf{U}^{-1} (\mathbf{Y} - \mathbf{X}\theta), \quad q(\theta) = (\hat{\theta} - \theta)^T W^{-1} (\hat{\theta} - \theta). \quad (3.41)$$

Using the normal equation, it is easily seen that

$$Q(\theta) = (\mathbf{Y} - \mathbf{X}\hat{\theta})^T \mathbf{U}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\theta}) + q(\theta). \quad (3.42)$$

Let the null and alternative hypotheses be

$$H_a : \theta \in C_a \text{ and } H_b : \theta \in C_b,$$

respectively, where $C_a \subset C_b$. Let θ^a and θ^b be the points at which $Q(\theta)$ is minimized over C_a and C_b , respectively. To define a test statistic based on generalized least squares without assuming normality, first note that $Q(\theta)$ is a generalized sum of squares. Therefore, $\{Q(\theta^a) - Q(\theta^b)\}$ is the reduction in this generalized sum of squares; it is a measure of the discrepancy between H_a and H_b . However, as it stands, this reduction cannot be used as a statistic for testing H_a against H_b because its null distribution depends on the unknown scale parameter σ^2 ; but, the null distribution of $\{Q(\theta^a) - Q(\theta^b)\}/\sigma^2$ does not depend on σ^2 . Therefore, it provides a good starting point for constructing a test statistic. As in the classical least squares theory, the idea is to replace the σ^2 in the denominator of $\{Q(\theta^a) - Q(\theta^b)\}/\sigma^2$ by a suitable estimate of it so that the distribution of the resulting statistic is independent of σ^2 . One such estimator turns out to be $Q(\theta^a)$, and we define the following statistic:

$$\bar{E}^2 = \{Q(\theta^a) - Q(\theta^b)\}/Q(\theta^a). \quad (3.43)$$

In the analysis of the classical normal theory linear models with no inequality constraints, tests on regression parameters are usually carried out using F -tests. A similar test statistic can also be developed when the hypotheses involve inequality constraints. Let

$$S^2 = \nu^{-1} (\mathbf{Y} - \mathbf{X}\hat{\theta})^{-1} \mathbf{U}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\theta})$$

with $\nu = (N - p)$, the usual error mean square. Motivated by $\{Q(\theta^a) - Q(\theta^b)\}/\sigma^2$ as a measure of discrepancy between H_a and H_b , we define

$$\bar{F} = \{Q(\theta^a) - Q(\theta^b)\}/S^2. \quad (3.44)$$

In the definition of the traditional unrestricted/two-sided F -statistic for testing equality constraints, the numerator of the F -statistic is $\{Q(\theta^a) - Q(\theta^b)\}/k$ where k is the number of equality constraints imposed by the null hypothesis. However, we did not introduce the divisor k in the definition of \bar{F} because it does not really simplify anything.

The difference between the definitions of \bar{E}^2 and \bar{F} is that \bar{E}^2 uses the dispersion of the restricted residuals, $\mathbf{Y} - \mathbf{X}\theta^a$, to estimate σ^2 while \bar{F} uses the dispersion of

the unrestricted residuals, $\mathbf{Y} - \mathbf{X}\hat{\theta}$. Because of the close similarity between these statistics, we would not expect the differences between them to be substantial, at least in large samples. Wright (1988) provided some calculations to compare the powers of \bar{E}^2 with \bar{F} in one-way lay-out models.

The basic idea of the \bar{F} statistic appears to be due to Kudo (1963) who suggested it when the observations are *iid* from $N(\theta, \sigma^2 U)$ and stated its null distribution.

Wolak (1987a) applied it for inference with linearly independent constraints in the linear model. Wright (1988) considered it for ANOVA models, and Silvapulle (1996a) obtained results for the case of testing for or against inequality constraints in a more general context. Dufour (1989) obtained bounds when the parameter spaces may not be cones. Wolak (1987a) also provided a detailed discussion of \bar{F} and how it relates to Kuhn-Tucker multipliers. Large sample results relating to \bar{F} will be obtained in the next chapter.

So far, the discussions have been motivated by generalized least squares, without assuming that E has any known distributional form. Now, suppose that $E \sim N(0, \sigma^2 U)$. Then the loglikelihood $L(\theta, \sigma^2)$ takes the form

$$L(\theta, \sigma^2) = -(1/2)\sigma^{-2}Q(\theta) - (1/2)np \log(\sigma^2) + Const. \quad (3.45)$$

For any given θ , $L(\theta, \sigma^2)$ is maximized at $\sigma^2 = (np)^{-1}Q(\theta)$. Therefore, the profile (or concentrated) loglikelihood $\bar{L}(\theta)$ is obtained by substituting $\sigma^2 = (np)^{-1}Q(\theta)$ in $L(\theta, \sigma^2)$. This leads to

$$\bar{L}(\theta) = \max L(\theta, \sigma^2) = -(1/2)np \log Q(\theta) + Const.$$

Therefore, θ^a and θ^b , which were introduced as generalized least squares estimates, are also the *MLE*'s of θ under H_a and H_b , respectively. Now, the likelihood ratio statistic for testing H_a against H_b is given by

$$LRT = np \log \{Q(\theta^a)/Q(\theta^b)\}. \quad (3.46)$$

Since

$$\bar{E}^2 = [1 - \exp\{-LRT/(np)\}],$$

LRT is an increasing function of \bar{E}^2 . Therefore, the \bar{E}^2 -test, which rejects H_a for large values of \bar{E}^2 , is equivalent to the likelihood ratio test when $E \sim N(0, \sigma^2 U)$.

3.9.1 Null Distributions

Many of the results that we established for \bar{E}^2 statistics also have corresponding results for \bar{E}^2 and \bar{F} . In what follows, the suffices A and B indicate that the test statistics correspond to Type A and Type B problems, respectively. The proof of the next result is given in the Appendix (page 130).

Theorem 3.9.1 *Let the linear model be as in (3.40), and assume that the distribution of E/σ does not depend on any unknown parameters. Then we have the following:*

- (a) *For the Type A testing problem,*

$$H_0 : \theta \in M \quad \text{against} \quad H_1 : \theta \in C,$$

the null distributions of \bar{E}^2 and \bar{F} do not depend on σ or on the value of θ in the null parameter space; in other words, the tests are similar.
 (b) For the Type B testing problem,

$$H_0: \theta \in C \quad \text{against} \quad H_1: \theta \notin C,$$

a least favorable null value of \bar{E}^2 and of \bar{F} is $\theta = 0$; if the null hypothesis also includes the linear equality constraint, $R\theta = 0$, then any θ satisfying $R\theta = 0$ is also a least favorable null value. ■

Suppose that the testing problem is as in part (a) or part (b) of the foregoing Theorem. An important consequence of this result is that the critical values and p -values can be computed by simulation for a large class of error distributions that includes the normal distribution; the simulations can be carried out at any fixed value of (θ, σ) . The simulation steps for computing the p -values of \bar{F} and \bar{E}^2 corresponding to any error distribution are given below.

Computation of the p -values for \bar{F} and \bar{E}^2 by simulation

Suppose that the conditions of Theorem 3.9.1 are satisfied. Let $F(\epsilon/\sigma)$ denote the cdf of the error E where σ may be unknown, but F is assumed known; for example, the distribution of E may be $N(0, \sigma^2 U)$. Now, the following steps would compute the p -value for \bar{F} (respectively, for \bar{E}^2).

1. Generate one observation of \mathbf{Y} from $F(\epsilon)$; this is same as generating one observation \mathbf{E} from $F(\epsilon)$ and then computing \mathbf{Y} as in (3.40) with $\theta = 0$.
2. Compute the \bar{F} (respectively, \bar{E}^2) statistic.

3. Repeat the previous two steps N times (say $N = 10000$), and estimate the p -value by M/N where M is the number of times the \bar{F} (respectively, \bar{E}^2) statistic in the second step exceeded its sample value. ■

Note that in the first step of the simulation, the observations may be generated from a distribution with any value for the common scale parameter because, in view of Theorem 3.9.1, the null distributions of \bar{F} and \bar{E}^2 do not depend on σ . Thus, if $E \sim N(0, \sigma^2 U)$ where σ is unknown, then it suffices to generate \mathbf{Y} from $N(0, U)$ in the first step of the simulation for computing the p -value. The level- α critical value can also be computed by ordering the N pseudo-random values of the statistic computed in step 2, and then finding upper level- α quantile of the empirical distribution of the N values.

The foregoing results would be particularly useful in model selection. For example, suppose that we wish to test H_0 vs H_1 where

$$\begin{aligned} H_0: y &= \theta_0 + \theta_1 x_1 + \dots + \theta_q x_q + e, \\ \text{and } H_1: y &= \theta_0 + \theta_1 x_1 + \dots + \theta_q x_q + \theta_{q+1} x_{q+1} + \dots + \theta_p x_p + e. \end{aligned}$$

Suppose also that some of $\{\theta_{q+1}, \dots, \theta_p\}$ are known to be nonnegative. Then the task of comparing the smaller model (H_0) to the larger one (H_1) reduces to a testing problem of the type just discussed.

If we assume that the error distribution is normal, then we can derive the null distributions of \bar{E}_A^2 and \bar{F}_A . The main result is given in the next theorem, and the proof is given in the Appendix (see page 131). The derivations use the same basic approach and constructions as for deriving the χ^2 distribution.

Theorem 3.9.2 Let the linear model be as in (3.40), and assume that $E \sim N(\mathbf{0}, \sigma^2 U)$.

(a) Let the null and alternative hypotheses be

$$H_0: \theta \in \mathcal{M} \quad \text{and} \quad H_1: \theta \in \mathcal{C}$$

respectively, where $\dim(\mathcal{M}) = q$. Then the null distributions of \bar{E}^2 and \bar{F} are given by

$$pr(\bar{E}_A^2 \leq c \mid H_0) = \sum_{i=0}^p w_i(p, W, \mathcal{C} \cap \mathcal{M}^\perp) pr[\beta\{i/2, (N-q-i)/2\} \leq c], \quad (3.47)$$

$$pr(\bar{F}_A \leq c \mid H_0) = \sum_{i=0}^p w_i(p, W, \mathcal{C} \cap \mathcal{M}^\perp) pr(iF_{i,\nu} \leq c)$$

where $\beta(a, b)$ is the beta distribution with parameters a and b .

(b) Let the null and alternative hypotheses be

$$H_0: \theta \in \mathcal{C} \quad \text{and} \quad H_1: \theta \notin \mathcal{C}$$

respectively. Then $\theta = 0$ is a least favorable null value for \bar{E}_B^2 and \bar{F}_B . Further, the least favorable null distributions of \bar{E}_B^2 and \bar{F}_B are:

$$pr(\bar{E}_B^2 \leq c \mid \theta = 0) = \sum_{i=0}^p w_{p-i}(p, W, \mathcal{C}) pr[\beta\{i/2, (N-p)/2\} \leq c],$$

$$pr(\bar{F}_B \leq c \mid \theta = 0) = \sum_{i=0}^p w_{p-i}(p, W, \mathcal{C}) pr(iF_{i,\nu} \leq c). \blacksquare$$

As in Sections 3.7 and 3.8 on results concerning $\bar{\chi}^2$ statistics, if the constraints are linear and independent then the foregoing results take simpler forms because the weights can be expressed in terms of those corresponding to the nonnegative orthant using the results in Proposition 3.6.1; this is stated in the next corollary.

Corollary 3.9.3 Let the linear model be as in (3.40), and assume that $E \sim N(\mathbf{0}, \sigma^2 U)$.

(a) Let R be a matrix of order $r \times p$, $r \leq p$, $\text{rank}(R) = r$, R_1 be a submatrix of R of order $q \times p$. Let the null and alternative hypotheses be

$$H_0: R\theta = \mathbf{0} \quad \text{and} \quad H_1: R_1\theta \geq \mathbf{0}.$$

Then we have the following for the null distributions of \bar{E}_A^2 and \bar{F}_A :

$$pr(\bar{E}_A^2 \leq c \mid R\theta = \mathbf{0}) = \sum_{i=0}^q w_i(q, R_1 V R_1^T) pr[\beta\{(r-q+i)/2, (N-p+q-i)/2\} \leq c]$$

$$pr(\bar{F}_A \leq c \mid R\theta = \mathbf{0}) = \sum_{i=0}^q w_i(q, R_1 V R_1^T) pr[(r-q+i)F_{r-q+i,\nu} \leq c].$$

(b) Now, let the null and alternative hypotheses be

$$H_1: R_1\theta \geq \mathbf{0}, \quad R_2\theta = \mathbf{0} \quad \text{and} \quad H_2: \theta \text{ is not restricted},$$

where R_1 is $s \times p$, R_2 is $t \times p$, and $[R_1^T, R_2^T]$ has full rank. Then, any θ satisfying $R_1\theta = 0$ and $R_2\theta = 0$ is a least favorable null value for \bar{E}_B^2 and \bar{F}_B . Further,

$$\Pr(\bar{F}_B \leq c|\theta = 0) = \sum_{i=0}^s w_{s-i}(s, A) \Pr[(t+i)/2, (N-k)/2] \leq c, \quad \blacksquare.$$

The foregoing results for Type B testing problems are useful even if some equality constraints on the regression parameters are present under both the null and alternative hypotheses. For example, consider the Type B testing problem and suppose that some linear equality constraints are present in the null and alternative hypotheses. Then the hypotheses take the form $H_1 : \theta \in C$ and $H_2 : \theta \in \mathcal{L}$, where $\mathcal{L} = \{A\beta : \beta \in \mathbb{R}^k\}$ for some matrix A of order $p \times k$ and $\text{rank}(A) = k$. Let $\mathcal{P} = \{\beta : A\beta \in C\}$. Then the linear model takes the form $Y = B\beta + E$, where $B = XA$, and the null and alternative hypotheses take the standard form, $H_0 : \beta \in \mathcal{P}$ and $H_1 : \beta \notin \mathcal{P}$. Now we can apply the foregoing results.

3.10 TESTS WHEN V IS UNKNOWN (PERLMAN'S TEST AND ALTERNATIVES)

In the foregoing sections we considered the cases when the covariance matrix V is either known or of the form $\sigma^2 U$ where U is known and σ is unknown. Now we consider the case when V is completely unknown.

Let X_1, \dots, X_n be n independently and identically distributed observations from $N(\theta, V)$ where V is unknown. Exact results for the likelihood ratio test were obtained by Perlman in his seminal paper, Perlman (1969). Corresponding results for a closely related test were obtained by Silvapulle (1995); conditional version of this test was introduced by Perlman and Wu (2002a). Other procedures include Tang (1994), Wang and McDermott (1993a, 1998b), Tang et al. (1989a), O'Brien (1984), and Follmann (1996b).

Let us consider the general problem of testing

$$H_a : \theta \in C_a \text{ against } H_b : \theta \in C_b \quad (3.49)$$

where C_a and C_b are *one-sided* closed cones; a set A is said to be *one-sided* if there exists a such that $A \subset \{x : a^T x \geq 0\}$. In this section, the cones considered are all *one-sided*. Since V is unknown, the loglikelihood up to a constant is

$$\ell(\theta, V) = \sum_{i=1}^n (-1/2)\{\log |V| + (X_i - \theta)^T V^{-1} (X_i - \theta)\}. \quad (3.50)$$

Further, we have that

$LRT = 2[\max\{\ell(\theta, V) : \theta \in C_b, V > 0\} - \max\{\ell(\theta, V) : \theta \in C_a, V > 0\}]$, where $V > 0$ means that V is positive definite. Note that the maximization needs to be carried out over θ and over all the positive definite matrices, $V > 0$, because now

V is an unknown parameter; therefore, this LRT is not the same as that for the case when V is known. Let

$$S = n^{-1} \sum (X_i - \bar{X})(X_i - \bar{X})^T$$

denote the sample covariance matrix, and let

$$\mathcal{U}(C_a, C_b) = \{\|\Pi_S(\bar{X}|C_b)\|_S^2 - \|\Pi_S(\bar{X}|C_a)\|_S^2\} \{1 + \|\bar{X} - \Pi_S(\bar{X}|C_b)\|_S^2\}^{-1}. \quad (3.51)$$

Now, we have the following that simplifies the LRT.

Proposition 3.10.1 (Perlman (1969), Theorem 5.1). $\mathcal{U}(C_a, C_b)$ is a strictly increasing function of LRT for testing $H_a : \theta \in C_a$ against $H_b : \theta \in C_b$.

It follows that the likelihood ratio test is equivalent to rejecting H_a for large values of $\mathcal{U}(C_a, C_b)$. Therefore, to apply the LRT, we do not have to compute the maximum of the function, $\ell(\theta, V)$, or any other function of V over the complicated parameter space, $\{V > 0\}$. If the constraints in C_a and C_b are linear, then $\mathcal{U}(C_a, C_b)$ can be computed using only a quadratic program. The foregoing result does not require observations to be *iid* but it does require an observation of Y from $N(\sqrt{n}\theta, V)$, and an observation S from the Wishart distribution $W(n-1, V)$ where Y and S are independent.

Now, let us introduce Perlman's test as follows:

Reject $H_0 : \theta \in C_a$ in favor of $H_1 : \theta \in C_b$ for large values of $\mathcal{U}(C_a, C_b)$.

Before we study properties of this test, let us introduce another class of tests. In classical multivariate analysis, tests of $\theta = 0$ against $\theta \neq 0$ are carried out using the well-known Hotelling's T^2 statistic, $\bar{X}^T S^{-1} \bar{X}$. For this unrestricted testing problem, the LRT reduces to $\bar{X}^T S^{-1} \bar{X}$, the Hotelling's T^2 . To motivate another test, let us note that, since S is positive definite, $\|\cdot\|_S^2 = \bar{X}^T S^{-1} \bar{X}$ is the squared distance between \bar{X} and the null parameter space. A similar statistic is available for the two-sample problem. Based on an idea of Shorack (1967), generalizations of the statistic $\bar{X}^T S^{-1} \bar{X}$ have been proposed when there are inequalities in θ (see Silvapulle (1995)). The basic form of the statistic is quite simple:

$$T^{*2} = \|\bar{X} - C_a\|_S^2 - \|\bar{X} - C_b\|_S^2. \quad (3.52)$$

In other words, T^{*2} is the difference between “the squared distance from \bar{X} to the null parameter space” and “the squared distance from \bar{X} to the alternative parameter space.” Let us remark that the foregoing T^{*2} and the Hotelling's T^2 have similarities with respect to their algebraic forms; but, Hotelling's T^2 is based on likelihood ratio and the T^{*2} in (3.52) is not.

For the special case of testing $H_0 : \theta = 0$ against $H_1 : \theta \geq 0$, Sen and Tsai (1999), obtained a union-intersection test. Their test statistic, which they denote by T_n^* , is the same as T^{*2} ; hence the latter can also be motivated from the union-intersection principle.

Perlman (1969) obtained the null distribution of $\mathcal{U}(C_a, C_b)$ for some important testing problems. Several corresponding results for the T^{*2} were obtained by Silvapulle (1995). The important ones are discussed below. Let U denote either the $\mathcal{U}(C_a, C_b)$ in (3.51) or the T^{*2} in (3.52). Let u denote the sample value of U . In general, the null distribution of U depends on (θ, V) . Therefore, $\text{pr}\{U \geq u \mid H_0, V\}$ is not an operational p -value. A suitable test procedure is,

reject H_0 if $\sup\{\text{pr}\{U \geq u \mid \theta, V\} : \theta \in H_0, V > 0\}$ is small, (3.53)

where the supremum is taken over θ in H_0 and over all positive definite matrices V ; for example, see Bickel and Doksum (1977) page 170. Similarly, if $\inf\{\text{pr}\{U \geq u \mid H_0, V\} : \theta \in H_0, V > 0\}$ is large then do not reject H_0 . Otherwise, the test is not conclusive. In some cases, for example, in Type A testing problems, $\text{pr}\{U \geq u \mid H_0, V\}$ does not depend on the value of θ in the null parameter space but it depends on V . In such cases, we shall also consider the possibility of estimating $\text{pr}\{U \geq u \mid H_0, V\}$ by substituting \hat{V} for V , but it does require caution (see Section 4.3.2). The null distributions of LRT and T^{*2} for inequality constrained problems, involve the following two random variables:

$$G(i, j) = \chi_i^2 / \chi_j^2 \quad \text{and} \quad H(k, r, n) = (\chi_k^2 / \chi_{n-r}^2)(1 + \chi_{r-k}^2 / \chi_{n-r+k}^2),$$

where the different χ^2 variates are independent. For any given c , the tail probabilities $\text{pr}\{G(i, j) \geq c\}$ and $\text{pr}\{H(k, r, n) \geq c\}$ can be computed easily by simulation because generation of random numbers from independent χ^2 variates is straight forward. Consequently, it will be seen that, in terms of computational demands, the difference between T^{*2} and LRT is small; the choice between the two would need to be based on their statistical properties.

It is also worth noting that $(j/i)G(i, j)$ is the usual F distribution with (i, j) degrees of freedom which we denote by $F_{i,j}$. Therefore,

$$\text{pr}\{G(i, j) \geq c\} = \text{pr}\{(i/j)F_{i,j} \geq c\},$$

and hence the F -distribution can be used to compute the tail probabilities of $G(i, j)$.

3.10.1 Type A Testing Problem

To consider a special Type A problem, let us partition θ as $(\theta_1^T, \theta_2^T)^T$ where θ_1 is $q \times 1$ and θ_2 is $r \times 1$; let V, S and X be partitioned conformably. Now, consider the test of

$$H_0 : \theta_2 = 0 \text{ against } H_1 : \theta_2 \in \mathcal{P} \quad (3.54)$$

where \mathcal{P} is a one-sided closed convex cone. Then, it is easily shown that the null distributions of LRT and T^{*2} do not depend on θ in the null parameter space; the argument is virtually the same as for Proposition 3.7.3. A bounds test can be carried out using the bounds given in the next result.

Proposition 3.10.2 (Perlman (1969) Theorem 6.3, and Silvapulle (1995) Theorem 3.) Let the testing problem be as in (3.54), and let \mathcal{U} denote $\mathcal{U}(C_a, C_b)$ for this testing

problem. Then under H_0 ,

$$\begin{aligned} \inf_{\theta, V} \text{pr}\{\mathcal{U} \geq u \mid H_0, V\} &= (1/2)\text{pr}\{G(1, n-r) \geq u\}, \\ \inf_{\theta, V} \text{pr}\{T^{*2} \geq u \mid H_0, V\} &= (1/2)\text{pr}\{H(1, r, n) \geq u\}. \end{aligned} \quad (3.55)$$

If \mathcal{P} contains an open set of dimension r then

$$\begin{aligned} \sup_{\theta, V} \text{pr}\{\mathcal{U} \geq u \mid H_0, V\} &= \\ (1/2)[\text{pr}\{G(r-1, n-r) \geq u\} + \text{pr}\{G(r, n-r) \geq u\}], \\ \sup_{\theta, V} \text{pr}\{T^{*2} \geq u \mid H_0, V\} &= \\ (1/2)[\text{pr}\{H(r-1, r, n) \geq u\} + \text{pr}\{G(r, n-r) \geq u\}]. \end{aligned} \quad (3.56)$$

The results in (3.55) and (3.56) are adequate to carry out the test as in (3.53). These two results are still quite general because the parameter space \mathcal{P} in (3.54) may include nonlinear inequalities in θ . In most practical situations, the inequalities in θ are linear and independent. In this case we can obtain closed forms for the null distributions of \mathcal{U} and T^{*2} , these are given in the next result.

Proposition 3.10.3 (Perlman (1969) Corollary 7.6, and Silvapulle (1995)). Let the null and alternative hypotheses be

$$H_0 : \theta_2 = 0 \text{ and } H_1 : A\theta_2 \geq 0 \quad (3.57)$$

respectively, where A is a given nonsingular matrix. Let \mathcal{U} denote the $\mathcal{U}(C_a, C_b)$ for this testing problem. Then

$$\begin{aligned} \text{pr}\{\mathcal{U} \geq u \mid H_0, V\} &= \sum_{i=1}^r \text{pr}\{G(i, n-r) \geq u\}w_i(r, W), \\ \text{pr}\{T^{*2} \geq u \mid H_0, V\} &= \sum_{i=1}^r \text{pr}\{H(k, r, n) \geq u\}w_i(r, W). \end{aligned} \quad (3.58)$$

where $W = A^{-1}V_{22}(A^{-1})^T$, and V_{22} is the bottom-right submatrix of V obtained by partitioning V to conform with the partitioning $\theta = (\theta_1^T, \theta_2^T)$. ■

Note that the expressions on the right-hand side of (3.58) depend on W but not on θ .

In the null parameter space. Therefore, we have to use its infimum and supremum over $V > 0$, given in (3.55) and (3.56), to obtain bounds on the p -value and perform the test as in (3.53). These bounds can be quite far from those corresponding to the true W if the sample size is small. This is the price that we pay for not knowing anything about the covariance matrix, V . One approach to improve the situation is to narrow down the class of matrices allowed for V based on nonsample information, and then take the supremum of (3.58) over that smaller class. Finding and evaluating the exact supremum over such a restricted set of matrices would typically be difficult, but it should be possible to generate a well-spread finite number of matrices within that smaller class, evaluate (3.58) for each V and take the maximum as an estimate of the p -value. If the sample size is large, an estimate of (3.58) is obtained by substituting an estimate \hat{W} of W for W in $w_i(r, W)$. One needs to be cautious in using this as a p -value because if \hat{W} is far from W , the resulting estimate of the p -value could be a

poor estimate. A more desirable approach would be to compute the expressions on the right-hand side of (3.58) for a range of possible values of W and take the maximum of these values. Even this has a certain degree of subjective element, but it is still better than simply substituting an estimate for W . A formal way of doing this is to take the supremum over a confidence region for W and make some adjustments for using only a confidence region rather than all possible W . A rigorous development of this procedure is provided in the next chapter (see Silvapulle (1996b) and Berger and Boos (1994)).

The LRT is biased when \mathcal{P} contains a non-empty open set; see Perlman (1969), p 558. There are points in the alternative parameter space that are close to the null space at which the restricted test, (3.53), has less power than the unrestricted one. However, for the restricted testing problem, we would still prefer the restricted test to the unrestricted one because the former would have better overall performance (see Silvapulle (1995) and Perlman and Wu (2002a)). The results of a simulation study shows that LRT is neither dominated by nor dominates T^{*2} . For Type A problems, some calculations show that the T^{*2} is likely to be more powerful than the LRT for values for θ near the boundary; for values of θ in the central direction, LRT is likely to be more powerful than the T^{*2} (see Silvapulle (1995)). At this stage we do not have sufficient results to recommend one over the other.

3.10.2 Type B Testing Problem

Now, consider the Type B problem of testing

$$H_1: \theta_2 \in \mathcal{P} \text{ against } H_2: \theta_2 \notin \mathcal{P}, \quad (3.60)$$

where as in the previous subsection we assume that θ is partitioned as $(\theta_1^T, \theta_2^T)^T$ and \mathcal{P} is assumed to be a one-sided closed cone. In the notation introduced at the beginning of this section, $C_a = \mathcal{P} \subset C_b = \mathbb{R}^r$. Let \mathcal{U} denote $\mathcal{U}(C_a, C_b)$ for this testing problem and let u denote the sample value of \mathcal{U} . It follows from (3.51) and (3.52) that T^{*2} and \mathcal{U} are equal. The null distribution of \mathcal{U} depends on (θ_2, V) and therefore, the test is

$$\text{reject } H_0 \text{ if } \sup_{\theta_2 \in \mathcal{P}, V > 0} \text{pr}(\mathcal{U} \geq u \mid \theta, V) \text{ is small.} \quad (3.61)$$

The main result on the null distribution of \mathcal{U} is the following:

Proposition 3.10.4 (Perlman (1969), section 8). *Let \mathcal{U} denote $\mathcal{U}(C_a, C_b)$ for the testing problem (3.60). Then we have that,*

1. $\mathcal{U} = \min\{(\bar{X}_2 - \mathbf{a})^T \mathbf{S}_{22}^{-1} (\bar{X}_2 - \mathbf{a}) : \mathbf{a} \in \mathcal{P}\}$.

2. $\text{pr}(U \geq u \mid H_0, V) \leq \text{pr}(U \geq u \mid \theta = 0, V)$, for any V .

3. $\sup_{\theta_2 \in C, V > 0} \text{pr}(U \geq u \mid \theta, V) \leq$

$$(1/2)[\text{pr}\{G(r-1, n-r-1) \geq u\} + \text{pr}\{G(r, n-r) \geq u\}]. \quad (3.62)$$

4. If \mathcal{P} contains an open set of dimension r , then

$$\inf_{\theta_2 \in \mathcal{P}, V > 0} \text{pr}(LRT \geq t \mid H_0, V) = 0.$$

There are other procedures for testing when V is unknown. None of the tests is uniformly best. A conditional version of the T^{*2} was suggested by Perlman and Wu (2002a); they also compared the performance of several statistics for this testing problem. Likelihood based conditional tests have been proposed by Wang and McDermott (1998a,b); the implementation of this procedure is difficult and there appears to be some issues that need to be resolved (see Perlman and Wu (1999)). Other procedures have also been suggested, such as the simple test in Tang (1994). Although the test in Tang (1994) is uniformly more powerful than the likelihood ratio test, it does not follow that this test is better than the likelihood ratio test or the T^{*2} test. These will be discussed later.

3.10.3 Conditional Tests of $H_0: \theta = 0$ vs $H_1: \theta \geq 0$

An approach to dealing with the unknown nuisance parameter V is to eliminate it from the null distribution by conditioning on a sufficient statistic for V or modify the test so that the probability of Type I error, as a function of V , is as flat as possible. Tests of the former type were developed by Wang and McDermott (1998a,b), and those of the latter type were developed by Perlman and Wu (2002a). These two types of tests are discussed below in turn.

Wang-McDermott test:

Let X_1, \dots, X_n be iid as $X \sim N(\theta, V)$. Let the null and alternative hypotheses be

$$H_0: \theta = 0 \quad \text{and} \quad H_1: \theta \geq 0$$

respectively. Let $W = \sum X_i X_i^T, S = W - n \bar{X} \bar{X}^T, \hat{\theta}$ be the mle of θ subject to $\theta \geq 0$,

$$U = n \hat{\theta}^T S^{-1} \hat{\theta} \{1 + (\bar{X} - \hat{\theta})^T S^{-1} (\bar{X} - \hat{\theta})\}^{-1},$$

and u_{obs} be the observed value of U . It was shown earlier that the unconditional LRT of H_0 against H_1 is, reject H_0 for large values of U . For this test, the p -value is

$$\sup_{V > 0} \text{pr}_V(U \geq u_{obs}).$$

Under H_0 , W is a complete sufficient statistic for V , and hence the distribution of U conditional on $W = w$ does not depend on V or any other unknown parameters. Now, the level- α conditional LRT of Wang and McDermott is the following:

1. $\text{reject } H_0 \text{ if } U \geq c_\alpha(w) \text{ where } \text{pr}_0\{\mathcal{U} \geq c_\alpha(w) \mid W = w\} = \alpha.$
2. Note that the probability on the LHS does not depend on V . Therefore, the p -value for this conditional test is

$$\text{pr}_0(\mathcal{U} \geq u_{obs} \mid W = w).$$

It turns out that the density of \bar{X} under H_0 , conditional on $\mathbf{W} = \mathbf{w}$, has a simple closed form. Further, \mathcal{U} has a simple closed-form for the purposes of computing it. Therefore, the p -value in (3.62) can be computed by using a simple brute force numerical integration or Monte Carlo; an easy-to-implement algorithm for computing the foregoing p -value is given in Wang and McDermott (1998a).

The power function of this conditional test approaches 1 uniformly in (θ, V) as $\theta^T V^{-1} \theta \rightarrow \infty$; further, the test is consistent. However, it is biased; see Sen and Tsai (1999).

It also turns out that $c_\alpha^* \geq c_\alpha(\mathbf{w})$ where $c_\alpha^*(\mathbf{w})$ and $c_\alpha(\mathbf{w})$ are the level- α critical values of the unconditional and conditional LRTs, respectively. Since the test statistics for the two tests are the same, it follows that the conditional LRT is uniformly more powerful than the unconditional test. It does not automatically follow that the unconditional test is better because the probability of Type I error is higher for the unconditional test although both tests have the same size.

Conditional versions of the T^{*2} discussed in the previous subsections were studied by Wang and McDermott (1998b). Essentially the same ideas are applicable. Briefly the test is reject H_0 if $T^{*2} > c'_\alpha$ where c'_α satisfies $\text{pr}_{0,V}\{T^{*2} > c'_\alpha \mid \mathbf{W} = \mathbf{w}\} = \alpha$.

Perlman-Wu test:

Now, let us consider the class of conditional tests proposed by Perlman and Wu (2002a). These authors proposed a class of tests in which they condition on K , where K is the number of positive components of $\hat{\theta}$. In contrast to the statistic \mathbf{W} in the foregoing discussions, the statistic K is not sufficient for \mathbf{V} , and the resulting conditional test turns out to be not similar. However, Perlman and Wu (2002a) identify a particular member of this class of conditional tests that is nearly similar, in a sense to be made precise in the theorems stated below.

Let C_α denote the class of tests of H_0 against H_1 for which the test statistic is \mathcal{U} and they are conditional on K . Let $\mathbf{c} = (c_1, \dots, c_k)$ where $c_i > 0$ for $i = 1, \dots, p$. Now, define the conditional test of Perlman and Wu as follows:

Do not reject H_0 if $K = 0$, and reject H_0 if $K = k$ and $\mathcal{U} > c_k$ for $k = 1, \dots, p$.

Let $T(\mathbf{c})$ denote this test. The unconditional LRT is a special case of this corresponding to $c_1 = \dots = c_k$. Each choice of \mathbf{c} defines a conditional test and it is not similar. One way of choosing an optimal \mathbf{c} is to find that value for which the minimum probability of Type I error over the null parameter space is a maximum over all possible values of $\{\mathbf{c}\}$. Perlman and Wu (2002a) obtained results in this direction.

Theorem 3.10.5 (Perlman and Wu (2002a)). *For any α , $0 < \alpha < 1$, we have that*

$$\max_{T(\mathbf{c}) \in C_\alpha} \inf_{V > 0} \text{pr}_{0,V}\{T(\mathbf{c}) \text{ rejects } H_0\} = \alpha/2.$$

Further, this maximum is attained only by the test $T(\mathbf{c}_\alpha)$ where $\mathbf{c}_\alpha = (c_{1,\alpha}, \dots, c_{p,\alpha})^T$ and $\text{pr}(\chi_k^2 / \chi_{n-p}^2 \geq c_{k,\alpha}) = \alpha$ where χ_k^2 and χ_{n-p}^2 are independent ($k = 1, \dots, p$). ■

The same conditioning approach is applicable to the statistic T^{*2} . The Perlman and Wu (2002a) type conditional version of T^{*2} is the following:

Do not reject H_0 if $K = 0$, and reject H_0 if $K = k$ and $T^{*2}(\mathbf{c}) > c_k$

for $k = 1, \dots, p$. The optimal choice is given in the next theorem.

Theorem 3.10.6 (Perlman and Wu (2002a)). *For any α , $0 < \alpha < 1$, we have that*

$$\max_{T^{*2}(\mathbf{c}) \in C_\alpha} \inf_{V > 0} \text{pr}_{0,V}\{T^{*2}(\mathbf{c}) \text{ rejects } H_0\} = \alpha/2.$$

Further, this maximum is attained only by the test $T^{*2}(\mathbf{c}_\alpha^*)$ where $\mathbf{c}_\alpha^* = (c_{1,\alpha}^*, \dots, c_{p,\alpha}^*)^T$ and $\text{pr}[(\chi_k^2 / \chi_{n-p}^2) \{1 + (\chi_{p-k}^2 / \chi_{n-p-k}^2)\} \geq c_{k,\alpha}^*] = \alpha$ where the various χ^2 variates are independent ($k = 1, \dots, p$). ■

The foregoing conditional tests are unlikely to be admissible. The acceptance regions of $T'(\mathbf{c})$ and $T^{*2}(\mathbf{c})$ are neither monotonic nor convex. In fact, Perlman and Wu (2002a) conjecture that the conditional tests of Wang and McDermott (1998a,b) and the unconditional tests, T^{*2} and LRT, are also inadmissible. A simulation study in Perlman and Wu (2002a) compared the performance of the foregoing conditional and unconditional tests; they observed that overall, $T(\mathbf{c})$ and $T^{*2}(\mathbf{c})$ performed better.

3.11 OPTIMALITY PROPERTIES

In the classical unrestricted setting, the parameter spaces tend to be open and several tests tend to have various optimality properties such as UMP, UMPI, unbiasedness, consistency, monotonic power function, etc. However, these optimality properties do not automatically carry over when the parameter space is not open. Further, often different optimality/desirable properties cannot be achieved without sacrificing some other desirable properties. In this section we shall consider consistency and monotonicity of power function. A thorough investigation of all these relevant issues is outside the scope of this book.

3.11.1 Consistency of Tests

Let n denote the sample size and let a level- α test of $H_0 : \theta \in \Theta_0$ against $H_1 : \theta \notin \Theta_0$ be "reject H_0 if $T_n \geq c_n$ " where T_n is a test statistic and c_n is the critical value. As the sample size increases, one would expect that the information in the sample about θ would also increase, and, hence, if T_n is a good test statistic then the probability of rejecting H_0 would tend 1 as $n \rightarrow \infty$ for $\theta \notin \Theta_0$. This leads to the following definition. A test of $H_0 : \theta \in \Theta_0$ is said to be *consistent* at $\theta \notin \Theta_0$ if the probability of rejecting the null hypothesis tends to 1 as $n \rightarrow \infty$ when the true value of the parameter is θ . Tests under inequality constraints are usually consistent at points in the alternative parameter space, but not necessarily at every point in the parameter space. Let us consider a simple example in two-dimensions to illustrate the main idea.

Example 3.11.1 Consistency of the LRT in 2 dimensions.

Let X_1, \dots, X_n be iid as $N(\boldsymbol{\theta}, I)$, and let the null and alternative hypotheses be defined as

$$H_0 : \boldsymbol{\theta} = \mathbf{0} \text{ and } H_1 : \boldsymbol{\theta} \geq \mathbf{0},$$

respectively, where $\boldsymbol{\theta} \in \mathbb{R}^2$. Then, $LRT = n\{\|\bar{X}\|^2 - \|\bar{X} - \mathbb{R}^{+2}\|^2\}$. First, let us consider the case when the true value of $\boldsymbol{\theta}$ satisfies $\theta_1 > 0$ and $\theta_2 > 0$. Since $\bar{X} \xrightarrow{P} \boldsymbol{\theta}$, it follows that, with probability tending to 1, \bar{X} lies in an arbitrarily small neighborhood of $\boldsymbol{\theta}$ contained in \mathbb{R}^{+2} and $LRT = n\|\bar{X}\|^2$; therefore, $LRT \xrightarrow{P} \infty$ as $n \rightarrow \infty$. Hence, LRT is consistent at $\boldsymbol{\theta}$ when $\theta_1 > 0$ and $\theta_2 > 0$. Similarly, by considering the interior and the boundary of each of the other three quadrants separately, it is easily seen that

LRT is consistent at $\boldsymbol{\theta}$ if and only if $\boldsymbol{\theta} \notin \{\boldsymbol{\theta} \in \mathbb{R}^2 : \theta_1 \leq 0, \theta_2 \leq 0\}$.

In fact, it may be verified that $\lim_{n \rightarrow \infty} \Pr\{LRT \text{ rejects } H_0|\boldsymbol{\theta}\}$ is 1, less than 1, or zero according as $\boldsymbol{\theta}$ is not in Q_3 , lies on the boundary of Q_3 , or lies in the interior of Q_3 , where $Q_3 = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \theta_1 \leq 0, \theta_2 \leq 0\}$. ■

Now consider the LRT of $H_1 : \boldsymbol{\theta} \geq \mathbf{0}$ against $H_2 : \boldsymbol{\theta} \not\geq \mathbf{0}$. Again, by considering the interior and boundary of each quadrant separately, we have that

LRT is consistent at any $\boldsymbol{\theta} \notin \{\boldsymbol{\theta} \in \mathbb{R}^2 : \theta_1 \geq 0, \theta_2 \geq 0\}$;

further, $\lim_{n \rightarrow \infty} \Pr\{LRT \text{ rejects } H_1|\boldsymbol{\theta}\}$ is 1, less than 1, or zero according as $\boldsymbol{\theta}$ is not in \mathbb{R}^{+2} , lies on the boundary of \mathbb{R}^{+2} , or lies in the interior of \mathbb{R}^{+2} . ■

Essentially the same arguments can be used to study the consistency properties of more general tests of inequality constrained hypotheses (see Problem 3.29). For the rest of this section all inner products will be assumed to be with respect to a given positive definite matrix, V . Let X_1, \dots, X_n be iid as $\mathbf{X} \sim N(\boldsymbol{\theta}, V)$ and V is known. Then we have the following:

- (1) Let the null and alternative hypotheses be $H_0 : \boldsymbol{\theta} \in \mathcal{M}$ and $H_1 : \boldsymbol{\theta} \in C$, respectively. Then $LRT = n\|\bar{X} - \mathcal{P}\|^2$ where $\mathcal{P} = (C^\circ \oplus \mathcal{M}) = (C \cap \mathcal{M}^\perp)^\circ$; see Theorem 3.7.1 on page 84 and Proposition 3.12.6 on page 118. Since $(\bar{X} - \boldsymbol{\theta}) = o_p(1)$ and \mathcal{P} is a closed set, it follows that LRT is consistent if and only if $\boldsymbol{\theta} \notin \mathcal{P}$. As a special case, let $\mathcal{M} = \{0\}$. Then, since $LRT = n\|\bar{X} - C^\circ\|^2$, $(\bar{X} - \boldsymbol{\theta}) = o_p(1)$ and C° is a closed set, it follows that LRT is consistent if and only if $\boldsymbol{\theta} \notin C^\circ$.

- (2) Let the null and alternative hypotheses be $H_1 : \boldsymbol{\theta} \in C$ and $H_2 : \boldsymbol{\theta} \notin C$, respectively. Then $LRT = n\|\bar{X} - C\|^2$, $(\bar{X} - \boldsymbol{\theta}) = o_p(1)$ and C is a closed set. Therefore, LRT is consistent at any $\boldsymbol{\theta} \notin C$.

Now, let us consider the case when V is unknown, and the null and alternative hypotheses are $H_0 : \boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \boldsymbol{\theta} \in C$, respectively. In this case, the LRT statistic is consistent (see Wang and McDermott (1998a) p 383). Sen and Tsai (1999) note that essentially the same arguments are applicable to T^{*2} because $T^{*2} \geq LRT$; hence T^{*2} is also consistent. Wang and McDermott (1998) also showed that their conditional LRT is consistent.

It may also be of interest to study the behavior of the power when the effect of the sample size is absorbed into the value of $\boldsymbol{\theta}$. Problem 3.30 is a result along this line.

3.11.2 Monotonicity of the Power Function

For the simple univariate case of testing $H_0 : \mu = 0$ against $H_1 : \mu > 0$ based on a single observation of X where $X \sim N(\mu, 1)$, the power of a test with critical region $\{x \in \mathbb{R} : x \geq c\}$, for any fixed $c > 0$, increases to 1 as μ increases to ∞ . Similarly, it is also known that for testing $\boldsymbol{\theta} = \mathbf{0}$ against $\boldsymbol{\theta} \neq \mathbf{0}$ based on a single observation of \mathbf{X} from $N(\boldsymbol{\theta}, V)$, where V is a given positive definite matrix, the power of the LRT ($= \|\mathbf{X}\|_V^2$) at $k\Delta$ increases to 1 as k increases from zero to ∞ , for any fixed $\Delta \neq \mathbf{0}$ (see Anderson (1955)); this result holds even when the distribution of \mathbf{X} is not normal but symmetric and unimodal (see Theorem 3.11.4).

The result of Anderson (1955) (stated as Theorem 3.11.4 on page 111) is applicable when the acceptance region is convex and symmetric. Therefore, it is usually not applicable when there are order restrictions on the parameters because the acceptance region is usually not symmetric and further it may not be even convex. However, it turns out that the power function of several tests when there are order restrictions on $\boldsymbol{\theta}$ are monotonic along certain straight lines. The main results concerning the monotonicity of power are stated below, and the proofs are given for some of them; for the others, references are given where detailed proofs may be found.

There are essentially two approaches to study the monotonicity of the power function. The simplest one is called the *set containment argument*, similar to that in Lemma 8.2 of Perlman (1969). This uses the fact $C + \boldsymbol{\theta} \subset C$ when $\boldsymbol{\theta} \in C$ and hence the distance from any given point to $C + \boldsymbol{\theta}$ is not more than the distance from the same point to C . The second method is *geometric*. This approach starts with the fact that the power of a test is equal to the volume of a region under the probability density function and above the critical region and uses geometric properties of this volume. So far, general results using such a geometric approach have assumed that the density function is elliptically symmetric and unimodal; the set containment argument does not require the density function to be elliptically symmetric or unimodal. Therefore, the results based on set containment arguments are applicable to a broader class of distributions. However, there are cases in which the geometric approach works but not the set containment argument; in some cases the converse is true. Therefore, both approaches are useful and essential.

For the rest of this section, the inner product (and hence $\|\cdot\|_V$, $\Pi()$, etc) will be assumed to be with respect to a given positive definite matrix V , for tests based on a single observation.

Proposition 3.11.1 *Let $\mathbf{X} = \boldsymbol{\theta} + E$ where the distribution of E does not depend on any unknown parameters. Let V be a given positive definite matrix, $\pi_{01}(\boldsymbol{\theta})$ denote the power function of $\|\mathbf{X} - (C \cap \mathcal{M}^\perp)\|_V^2$ for testing $H_0 : \boldsymbol{\theta} \in \mathcal{M}$ against $H_1 : \boldsymbol{\theta} \in C$; similarly, let $\pi_{12}(\boldsymbol{\theta})$ denote the power function of $\|\mathbf{X} - C\|_V^2$ for testing $H_1 : \boldsymbol{\theta} \in C$ against $H_2 : \boldsymbol{\theta} \notin C$. Then, $\pi_{01}(\boldsymbol{\theta}) \leq \pi_{01}(\boldsymbol{\theta}')$ for $\boldsymbol{\theta} - \boldsymbol{\theta}' \in C^\circ$, and $\pi_{12}(\boldsymbol{\theta}) \geq \pi_{12}(\boldsymbol{\theta}')$ for $\boldsymbol{\theta}' - \boldsymbol{\theta} \in C$.*

Proof: The proof follows a set containment argument. To prove the first, let $\boldsymbol{\theta} - \boldsymbol{\theta}' \in C^\circ$. Then $C^\circ + (\boldsymbol{\theta} - \boldsymbol{\theta}') \subset C^\circ$, and hence $C^\circ - \boldsymbol{\theta}' \subset C^\circ - \boldsymbol{\theta}$. Now, $\pi_{01}(\boldsymbol{\theta}) =$

$$\begin{aligned} \text{pr}\{\|\mathbf{X} - \mathcal{C}^o\|_V^2 \geq c|\theta\} &= \text{pr}\{\|\mathbf{E} + \theta - \mathcal{C}^o\|_V^2 \geq c\} = \text{pr}\{\|\mathbf{E} - (\mathcal{C}^o - \theta)\|_V^2 \geq c\} \leq \\ \text{pr}\{\|\mathbf{E} - (\mathcal{C}^o - \theta')\|_V^2 \geq c\} &= \pi_{01}(\theta'). \end{aligned}$$

To prove the second, let $\theta' - \theta \in C$. Then $C + (\theta' - \theta) \subset C$, and hence $C - \theta \subset C - \theta'$. Now, $\pi_{12}(\theta) = \text{pr}\{\|\mathbf{X} - \mathcal{C}\|_V^2 \geq c|\theta\} = \text{pr}\{\|\mathbf{E} - (\mathcal{C} - \theta)\|_V^2 \geq c\} \geq \text{pr}\{\|\mathbf{E} - (\mathcal{C} - \theta')\|_V^2 \geq c\} = \pi_{12}(\theta')$. ■

A simple diagram can be used to illustrate this result. The first part of the result says that for testing $H_0 : \theta = \mathbf{0}$ against $H_1 : \theta \in C$, the power of the test at any point, say θ_0 , increases in any direction of $-\mathcal{C}^o$ and decreases in any direction of \mathcal{C}^o . Similarly, the second part of the result says that for testing $H_1 : \theta \in C$ against $H_2 : \theta \notin C$, the power of the test at any point, say θ_0 , increases in any direction of $-C$ and decreases in any direction of C .

Now, we state the main results that have been obtained using the geometric approach. Let h be a nonincreasing function on $[0, \infty)$, V be a positive definite matrix of order $p \times p$, and \mathbf{X} be a random variable with a density function $|V|^{-1/2}h(\|\mathbf{x} - \theta\|^2)$; a density function of this form is said to be *elliptically symmetric* and *unimodal* about θ . Let

$$p(\theta) = 1 - \int_A |V|^{-1/2}h(\|\mathbf{x} - \theta\|_V^2)dx, \quad A \subset \mathbb{R}^p, \theta \in \mathbb{R}^p.$$

For several tests, including the LRT, of $H_0 : \theta \in \mathcal{M}$ against $H_1 : \theta \in C$ and of $H_1 : \theta \in C$ against $H_2 : \theta \notin C$, the power of the test at θ takes the form $p(\theta)$, where A is the acceptance region. Therefore, monotonicity properties of functions of the form $p(\theta)$ are of interest. Anderson's theorem is applicable to establish the monotonicity of $p(\theta)$ when A is symmetric and convex. When there are inequality constraints, the acceptance region A is usually not symmetric and it may not be even convex. Mukerjee et al. (1986) studied the monotonicity properties of $p(\theta)$ when A satisfies certain conditions that include that A be convex but not necessarily symmetric. They showed that $p(k\theta_0)$ is monotonic in $|k|$ for $\theta_0 \in \mathbb{R}^p$, their main result (stated as Theorem 3.11.5 on page 111) leads to the following useful corollary.

Corollary 3.11.2 (Mukerjee et al. (1986)). *Let the distribution of \mathbf{X} be elliptically symmetric and unimodal about θ . Let $\theta_1 \in \mathbb{R}^p$ and $\theta_2 \in C$. Then as a function of k , $\pi_{12}(\theta_1 + k\theta_2)$ is nondecreasing for $k \in (-\infty, \infty)$ and $\pi_{12}(\theta_1 + k\bar{\theta}_1)$ is nondecreasing for $k \geq -1$ where $\bar{\theta}_1 = \Pi(\theta_1|\mathcal{C}^o)$.* ■

It is also shown in Robertson et al. (1988), page 106, that (a) $\pi_{01}(k\Delta)$ is nondecreasing in k , where $(k > 0)$ and $\Delta \in C$, and (b) $\pi_{12}(k\Delta)$ and $\pi_{12}^*(k\Delta)$ are nondecreasing in k , where $k > 0$, $\Delta \in C^o \oplus \mathcal{M}$ and π_{12}^* is the power function of the \bar{E}^2 -test.

For some \bar{E}^2 -tests, the acceptance region is not convex and hence the results of Mukerjee et al. (1986) are not applicable. Hu and Wright (1994b) studied the monotonicity of $p(\theta_0 + t\theta_1)$ in t for some given θ_0 and θ_1 and when A takes a special form; their regularity conditions do not require A to be convex or symmetric. Their result extends some of the results in Mukerjee et al. (1986), and is applicable for the \bar{E}^2 -test of H_0 against H_1 . A useful corollary of their main result (see Theorem 3.11.6) is the following:

Corollary 3.11.3 (Hu and Wright (1994b)). *Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_n$ are iid as $N(\theta, \sigma^2 U)$ where U is known and σ is unknown. Then, the power function of the LRT (namely, \bar{E}^2 -test) of $H_0 : \theta \in \mathcal{M}$ against $H_1 : \theta \in C$ is nondecreasing on each line segment starting at a point in \mathcal{M} and continuing in the direction of a vector in C . Consequently, this test is unbiased as well.* ■

Since $\bar{F} = \bar{\chi}^2/S^2$ and S^2 is independent of θ , monotonicity of $\bar{\chi}^2$ carry over to \bar{F} as well; hence, monotonicity of \bar{F} is not studied separately.

The main results concerning the monotonicity of $p(\theta)$ are stated below for completeness.

Theorem 3.11.4 (Anderson (1955)). *Let A be a convex and symmetric set in \mathbb{R}^p . Let f be a function such that (i) $f(x) = f(-x)$, (ii) $\{x : f(x) \geq u\}$ is convex for every $u > 0$, and (iii) $\int_A f(x)dx < \infty$ (in the Lebesgue sense). Then $\int_A f(x + ky) \geq \int_A f(x + y)dy$ for $0 \leq k \leq 1$.* ■

Theorem 3.11.5 (Mukerjee et al. (1986)). *Let $\theta \in \mathbb{R}^p$ and $S_\theta = \{b\theta : -\infty < b < \infty\}$ be the linear space spanned by θ ; for a given convex set $A \subset \mathbb{R}^p$ let the positive part of A in the direction of θ be defined by*

$$A^+ = \{\mathbf{x} \in A : \Pi_V(\mathbf{x}|S_\theta) = b\theta \text{ where } b > 0\}.$$

Let f be an elliptically symmetric and unimodal density function. Suppose that $\int_A f(x)|S_\theta \in A$ for each $\mathbf{x} \in A^+$. Then, $\int_{A^-k\theta} f(x)dx$ is a nonincreasing function of k for $k > 0$. ■

Theorem 3.11.6 (Hu and Wright (1994b)). *Suppose that X has the unimodal and elliptically symmetric density $|V|^{-1/2}h(\|\mathbf{x}\|_V^2)$ where h is nonincreasing on $[0, \infty)$. Let $C \subset \mathbb{R}^p$ be a closed convex cone, \mathcal{M} be a linear subspace such that $\mathcal{M} \subset C$, and $A = \{\mathbf{x} \in \mathbb{R}^p : \|\Pi_V(\mathbf{x})|\mathcal{M}\| - \|\Pi_V(\mathbf{x}|C)\|_V^2 \leq a + b\|\mathbf{x} - \Pi_V(\mathbf{x}|C)\|_V^2\}$, where $a > 0$ and $b > 0$. If $\mu_0 \in \mathcal{M}$ and $\nu_0 \in C$ then $\text{pr}(\mathbf{X} + (\mu_0 + t\nu_0) \in A)$, is nonincreasing in t for $t > 0$.* ■

Now let us consider the case when the covariance matrix V is completely unknown, and the null and alternative hypotheses are $H_0 : \theta = \mathbf{0}$ and $H_1 : \theta \in C$, respectively. The power function of LRT and T^{*2} are monotonic. In particular, the set containment argument can be used to establish that the power functions of LRT and T^{*2} at θ_0 increase in any direction of C (see Sen and Tsai (1999)).

0.12 APPENDIX 1: CONVEX CONES, POLYHEDRALS, AND PROJECTIONS

0.12.1 Introduction

¹This Appendix provides a brief account of the main results on projections onto convex cones in finite dimensional spaces as they relate to the statistical inference problems

discussed in this book. An excellent account of projections onto convex cones is given in Stoer and Witzgall (1970)(SW); other references include Bazaraa et al. (1993) and Hiriart-Urruty and Lemaréchal (1993). Zarantonello (1998) provides an extensive account of projections onto convex sets in Hilbert spaces; while some of the results and ideas therein are relevant, such general results are not required at this stage. In this Appendix, we provide a reasonably self-contained discussion that would be adequate for this book.

Let \mathbb{R}^p denote the p -dimensional Euclidean space, V be a $p \times p$ symmetric positive definite matrix, $\mathbf{x} \in \mathbb{R}^p$, and $\mathbf{y} \in \mathbb{R}^p$. Then, $\langle \mathbf{x}, \mathbf{y} \rangle_V = \mathbf{x}^T V^{-1} \mathbf{y}$ defines an inner product on \mathbb{R}^p . This induces the corresponding norm $\|\mathbf{x}\|_V = \langle \mathbf{x}, \mathbf{x} \rangle_V^{1/2}$; the corresponding distance between \mathbf{x} and \mathbf{y} is $\|\mathbf{x} - \mathbf{y}\|_V$. If $\mathbf{x}^T V^{-1} \mathbf{y} = 0$ then we say that \mathbf{x} and \mathbf{y} are orthogonal with respect to V , which we denote by $\mathbf{x} \perp_V \mathbf{y}$. We may abbreviate ‘orthogonal with respect to V ’, ‘distance with respect to V ’ to V -orthogonal, V -distance, etc. However, if V is obvious we will drop the reference to V .

If $\mathbf{x} \perp_V \mathbf{y}$ then a version of the Pythagoras theorem holds: namely, $\|\mathbf{x} + \mathbf{y}\|_V^2 = \|\mathbf{x}\|_V^2 + \|\mathbf{y}\|_V^2$. A consequence of this is that the shortest distance between a point and a plane is the distance from the point to the plane along a line that is orthogonal to the plane; here, distance and orthogonality are with respect to $\langle \cdot, \cdot \rangle_V$.

Let C be a closed convex set in \mathbb{R}^p and $\mathbf{x} \in \mathbb{R}^p$. Let $\tilde{\mathbf{x}}$ in C be the point in C that is closest to \mathbf{x} with respect to the distance $\|\cdot\|_V$; thus,

$$\|\mathbf{x} - \tilde{\mathbf{x}}\| = \min_{\theta \in C} (\mathbf{x} - \theta)^T V^{-1} (\mathbf{x} - \theta).$$

The vector $\tilde{\mathbf{x}}$ is called the *projection* of \mathbf{x} onto C , and is denoted by $\Pi_V(\mathbf{x} | C)$ (SW p 47); thus

$$\tilde{\mathbf{x}} = \Pi_V(\mathbf{x} | C) = \arg \min_{\theta \in C} (\mathbf{x} - \theta)^T V^{-1} (\mathbf{x} - \theta).$$

Almost all the results on projections onto convex cones discussed here are generalizations of the corresponding results for projections onto linear spaces. Therefore, the basic results on the latter are essential for the rest of this appendix. For detailed discussions of relevant results on projections onto linear spaces, see textbooks on linear models, such as Arnold (1981), Chapter 2, or Rao (1972).

We may assume, without loss of generality, that $V = I$ when considerations are restricted to $\langle \mathbf{x}, \mathbf{y} \rangle_V$ where \mathbf{x} and \mathbf{y} are elements of \mathbb{R}^p . To show this, let $V^{-1} = \mathbf{U}^T \mathbf{U}$ be a factorization of V^{-1} , for example, the Cholesky factorization. Now apply the invertible linear transformation \mathbf{U} on \mathbb{R}^p . Then, any two points \mathbf{x} and \mathbf{y} in \mathbb{R}^p are mapped to \mathbf{p} and \mathbf{q} in \mathbb{R}^p where $\mathbf{p} = \mathbf{U} \mathbf{x}$ and $\mathbf{q} = \mathbf{U} \mathbf{y}$. Since $\langle \mathbf{x}, \mathbf{y} \rangle_V = \mathbf{x}^T V^{-1} \mathbf{y} = (\mathbf{U} \mathbf{x})^T (\mathbf{U} \mathbf{y}) = \mathbf{p}^T \mathbf{q} = \langle \mathbf{p}, \mathbf{q} \rangle_I$, the distance between \mathbf{x} and \mathbf{y} in $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_V\}$ is equal to the distance between their images \mathbf{p} and \mathbf{q} in $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_I\}$. Therefore, for the purposes of studying properties that depend only on distances between points, we may identify $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_V\}$ with $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_I\}$ and, hence, without loss of generality, assume that V is the identity matrix. Instead of writing $\langle \mathbf{p}, \mathbf{q} \rangle_I$ and $\|\mathbf{p} - \mathbf{q}\|_I$ we shall use the simpler notation, $\mathbf{p}^T \mathbf{q}$ and $\|\mathbf{p} - \mathbf{q}\|$, respectively.

Remark: Since U and U^{-1} are one-to-one, onto and, continuous, $U : \{\mathbb{R}^p, \langle \cdot, \cdot \rangle_V\} \rightarrow \{\mathbb{R}^p, \langle \cdot, \cdot \rangle_I\}$ is a *homeomorphism* and hence $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_V\}$ and $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_I\}$ are *homeomorphic*. Thus, these two spaces have the same topological structure. Since $\langle \mathbf{x}, \mathbf{y} \rangle_V = \langle U\mathbf{x}, U\mathbf{y} \rangle_I$, U is an *isomorphism* and $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_V\}$ and $\{\mathbb{R}^p, \langle \cdot, \cdot \rangle_I\}$ are *isomorphic*.

A set S is said to be the *orthogonal sum* of S_1 and S_2 if $S = S_1 + S_2 = \{s_1 + s_2 : s_1 \in S_1, s_2 \in S_2\}$ and every element $\mathbf{x} \in S$ admits a unique *orthogonal decomposition* of the form $\mathbf{x} = \mathbf{y} + \mathbf{z}$ with $\mathbf{y} \in S_1, \mathbf{z} \in S_2$ and $\mathbf{y}^T \mathbf{z} = 0$. Such an orthogonal sum will be denoted by $S_1 \oplus S_2$. As an example, the *direct sum* of two orthogonal linear subspaces is an orthogonal sum.

For any set $S \subset \mathbb{R}^p$, its *orthogonal complement* S^\perp is defined as $\{\mathbf{y} \in \mathbb{R}^p : \mathbf{y}^T \mathbf{x} = 0, \text{ for all } \mathbf{x} \in S\}$. Clearly, S^\perp is a linear space. For a set $S \subset \mathbb{R}^p$, we define the *polar cone* S° as

$$S^\circ = \{\mathbf{x} : \mathbf{x}^T \mathbf{y} \leq 0 \text{ for all } \mathbf{y} \in S\}.$$

The next result states a few important but elementary properties of cones and polar cones; the proofs are easy and hence omitted.

Proposition 3.12.1 *Let C be a closed convex cone and P and Q be subsets of \mathbb{R}^p . Then we have the following:*

- (1) *C is closed under addition: $\mathbf{x} \in C$ and $\mathbf{y} \in C \Rightarrow \mathbf{x} + \mathbf{y} \in C$.*
- (2) *$C + \mathbf{x} \subset C \subset C - \mathbf{x}$, $\forall \mathbf{x} \in C$.*
- (3) *P° is a closed convex cone.*
- (4) *If $P \subset Q$ then $Q^\circ \subset P^\circ$.*
- (5) *If P is a linear subspace then $P^\perp = P^\circ$.*

Clearly, we can also define orthogonal sum, orthogonal complement, polar cone, etc. with respect to any symmetric and positive definite matrix, V ; all that it requires is to replace $\mathbf{x}^T \mathbf{y}$ by $\mathbf{x}^T V^{-1} \mathbf{y}$ in the foregoing definitions. For example, the *polar cone of S with respect to V* is $\{\mathbf{x} : \mathbf{x}^T V^{-1} \mathbf{y} \leq 0 \text{ for all } \mathbf{y} \in S\}$. This can be denoted by S_V° ; however, we will write S° instead if the suffix V is obvious. Similarly, the *orthogonal complement of S with respect to V* is defined as $\{\mathbf{y} \in \mathbb{R}^p : \mathbf{y}^T V^{-1} \mathbf{x} = 0, \text{ for all } \mathbf{x} \in S\}$. This can be denoted by S_V^\perp , although we would typically write S^\perp if the suffix V is obvious. In what follows, the main results are presented for the case $V = I$; however, the corresponding results for a general V would be obvious.

For statistical inference based on the linear model, projections onto linear spaces play an important role. The corresponding results on projections onto convex cones play a similar important role in constrained statistical inference in the exact normal and large sample theory.

3.12.2 Projections onto Convex Cones

Before we consider projections onto convex cones, let us recall some results on projections onto linear spaces. Let P be a $p \times p$ matrix. We say that P is a *projection matrix* if $P = P^T$ and $P^2 = P$.

Proposition 3.12.2 Let P be $p \times p$ projection matrix, $\mathbf{Q} = \mathbf{I} - \mathbf{P}$, $\mathcal{L}_P = \{\mathbf{P}\mathbf{x} : \mathbf{x} \in \mathbb{R}^p\}$ and $\mathcal{L}_Q = \{\mathbf{Q}\mathbf{x} : \mathbf{x} \in \mathbb{R}^p\}$. Then we have the following: (1) \mathbf{P} is positive semi-definite. (2) \mathbf{Q} is a projection matrix. (3) \mathcal{L}_P and \mathcal{L}_Q are orthogonal linear subspaces. (4) Any $\mathbf{x} \in \mathbb{R}$ can be represented uniquely as $\mathbf{x} = \mathbf{p} + \mathbf{q}$ where $\mathbf{p} \in \mathcal{L}_P$ and $\mathbf{q} \in \mathcal{L}_Q$.

Proof: See Bazaraa et al. (1993), page 448, for example. ■

Let \mathcal{L} be a linear subspace of \mathbb{R}^p . Some results on projections onto linear spaces that are important for inference in linear models, include the following:

1. The projection $\Pi(\mathbf{x}|\mathcal{L})$ of \mathbf{x} onto \mathcal{L} exists and is unique.
2. Let $\mathcal{L} = \{\mathbf{X}\beta : \beta \in \mathbb{R}^k\}$, and $\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, where $k \leq p$, \mathbf{X} is $p \times k$ and $\text{rank}(\mathbf{X}) = k$. Then, $\Pi(\mathbf{x}|\mathcal{L}) = \mathbf{P}\mathbf{x}$ for every \mathbf{x} and \mathbf{P} is called the projection matrix onto \mathcal{L} .
3. $\mathbf{x} - \Pi(\mathbf{x}|\mathcal{L})$ is the projection of \mathbf{x} onto \mathcal{L}^\perp , and the projection matrix onto \mathcal{L}^\perp is $(\mathbf{I} - \mathbf{P})$.
4. If \mathbf{x} is projected onto \mathcal{L} and then the projection itself is projected onto a subspace \mathcal{M} of \mathcal{L} , the resulting projection is the projection of \mathbf{x} onto \mathcal{M} [i.e., $\Pi\{\Pi(\mathbf{x}|\mathcal{L})|\mathcal{M}\} = \Pi(\mathbf{x}|\mathcal{M})$].

The essence of many of these results carries over to projections onto convex cones and plays similar important roles in constrained statistical inference. The main difference that causes technical difficulties is that the projection of \mathbf{x} onto a convex cone is not a linear function of \mathbf{x} and hence it cannot be expressed as $\mathbf{Q}\mathbf{x}$ where \mathbf{x} is arbitrary and \mathbf{Q} does not depend on \mathbf{x} . The geometric nature of projections onto linear spaces and their relevance in the context of statistical inference in linear models appear in standard textbooks on linear models. Therefore, we do not discuss such results here, but provide a discussion of the corresponding results for closed convex cones and schematic diagrams to help interpret the results.

Just as in the case for linear spaces, projection of a point onto a closed convex set exists uniquely, and it can be characterized by the angle between the line of projection and the set. These results are given below.

Proposition 3.12.3 Let C be a nonempty closed convex set in \mathbb{R}^p and $\mathbf{a} \in \mathbb{R}^p \setminus C$. Then we have the following:

- (a) There exists a unique $\mathbf{a}^* \in C$ that is closest to \mathbf{a} in the sense that, for $\mathbf{x} \in C$, $\|\mathbf{a} - \mathbf{x}\|$ is a minimum at $\mathbf{x} = \mathbf{a}^*$.
- (b) The point $\mathbf{a}^* \in C$ is the unique point in C closest to \mathbf{a} if and only if $(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{x} - \mathbf{a}^*) \leq 0$ for all $\mathbf{x} \in C$.
- (c) There exists a vector \mathbf{p} and a scalar α such that $\mathbf{p}^T \mathbf{a} > \alpha$ and $\mathbf{p}^T \mathbf{x} \leq \alpha$ for every $\mathbf{x} \in C$.

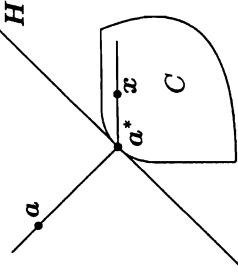


Fig. 3.14 The projection \mathbf{a}^* of \mathbf{a} onto the convex set C ; the plane H , which is orthogonal to $\mathbf{a} - \mathbf{a}^*$, separates C and the point \mathbf{a} . ■

[For (a), see SW 3.3.1 or Bazaraa et al. (1993), page 43; for (b) see Bazaraa et al. (1993), page 43 or Hiriart-Urruty and Lemaréchal (1993), page 117]

Proof of (a): Let $f : C \rightarrow \mathbb{R}$ be defined by $f(\mathbf{x}) = \|\mathbf{a} - \mathbf{x}\|$. Let $\mathbf{x}_0 \in C$ and $S = \{\mathbf{y} \in \mathbb{R}^p : f(\mathbf{y}) \leq f(\mathbf{x}_0)\}$. Then $S \cap C$ is nonempty and compact. Since f is continuous, it has a minimum over $S \cap C$, say at \mathbf{a}^* . Clearly, $f(\mathbf{a}^*) < f(\mathbf{x})$ for $\mathbf{x} \notin S$. Therefore, \mathbf{a}^* is also the global minimum of $f(\mathbf{x})$ over C . Now, to show that \mathbf{a}^* is unique, suppose that $\tilde{\mathbf{a}} \in C$, $\tilde{\mathbf{a}} \neq \mathbf{a}^*$, and $\|\mathbf{a} - \tilde{\mathbf{a}}\| = \|\mathbf{a} - \mathbf{a}^*\|$. Let $\bar{\mathbf{a}} = 0.5(\mathbf{a}^* + \tilde{\mathbf{a}})$. Clearly, $\bar{\mathbf{a}} \in C$ because C is convex. Further, we have that $\|\mathbf{a} - \tilde{\mathbf{a}}\|^2 = 0.5\|\mathbf{a} - \mathbf{a}^*\|^2 + 0.5\|\mathbf{a} - \tilde{\mathbf{a}}\|^2 - 0.25\|\mathbf{a} - \mathbf{a}^*\|^2 < \|\mathbf{a} - \mathbf{a}^*\|^2$. Thus, $\bar{\mathbf{a}}$ is closer to \mathbf{a} than \mathbf{a}^* , this is a contradiction. Therefore, \mathbf{a}^* is unique.

Proof of (b): Assume that $\mathbf{a}^* \in C$ and $(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{x} - \mathbf{a}^*) \leq 0$ for every $\mathbf{x} \in C$. Let $\mathbf{x} \in C$. Then $\|\mathbf{a} - \mathbf{x}\|^2 = \|\mathbf{a} - \mathbf{a}^*\|^2 + \|\mathbf{a}^* - \mathbf{x}\|^2 + 2(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{a}^* - \mathbf{x})$. Since $(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{a}^* - \mathbf{x}) \geq 0$ by assumption, we have that $\|\mathbf{a} - \mathbf{x}\|^2 \geq \|\mathbf{a} - \mathbf{a}^*\|^2$. Therefore, \mathbf{a}^* is the point in C closest to \mathbf{a} . To prove the converse, assume that $\mathbf{a}^* \in C$ is the point in C closest to \mathbf{a} . Let $\mathbf{x} \in C$. Then, $\mathbf{a}^* + \lambda(\mathbf{x} - \mathbf{a}^*) \in C$ for $0 < \lambda \leq 1$ because C is convex. Now, since \mathbf{a}^* is closer to \mathbf{a} than $\mathbf{a}^* + \lambda(\mathbf{x} - \mathbf{a}^*)$, we have $\|\mathbf{a} - \mathbf{a}^*\|^2 \leq \|\mathbf{a} - \{\mathbf{a}^* + \lambda(\mathbf{x} - \mathbf{a}^*)\}\|^2 = \|\mathbf{a} - \mathbf{a}^*\|^2 + \lambda^2\|\mathbf{x} - \mathbf{a}^*\|^2 - 2\lambda(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{x} - \mathbf{a}^*)$. Therefore, $2\lambda(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{x} - \mathbf{a}^*) \leq \lambda^2\|\mathbf{a} - \mathbf{a}^*\|^2$. Now, divide by $\lambda (> 0)$ and take the limit as $\lambda \rightarrow 0^+$. This leads to $(\mathbf{a} - \mathbf{a}^*)^T(\mathbf{x} - \mathbf{a}^*) \leq 0$ for every $\mathbf{x} \in C$.

Proof of (c): Let $\mathbf{a}^* = \Pi(\mathbf{a} | C)$, \mathbf{p} be a unit vector parallel to $(\mathbf{a} - \mathbf{a}^*)$ and $\mathbf{a} = \mathbf{p}^T \mathbf{a}^*$.

Now, the proof follows from part (b). ■

Fig. 3.14 illustrates the nature of the foregoing results. It shows that if $\Pi(\mathbf{a}|C)$ is treated as the origin then the angle between \mathbf{a} and \mathbf{x} must be obtuse for any $\mathbf{x} \in C$ (the angle θ between the vectors \mathbf{x} and \mathbf{y} in \mathbb{R}^p is defined by $\cos \theta = \mathbf{x}^T \mathbf{y} / \{\|\mathbf{x}\| \|\mathbf{y}\|\}^{-1/2}$ and this angle is said to be obtuse if $\cos \theta < 0$). Further, the hyperplane $H = \{\mathbf{x} \in \mathbb{R}^p : \mathbf{x}^T \mathbf{p} = \alpha\}$ is orthogonal to $(\mathbf{a} - \mathbf{a}^*)$, passes through \mathbf{a}^* , and separates \mathbf{a} and C ; here separates means that the point \mathbf{a} and the set C lie on the opposite sides of the hyperplane.

An important result on projections onto linear spaces is that, for any $\mathbf{x} \in \mathbb{R}^p$ and linear space \mathcal{L} , we have a unique orthogonal decomposition of the form $\mathbf{x} =$

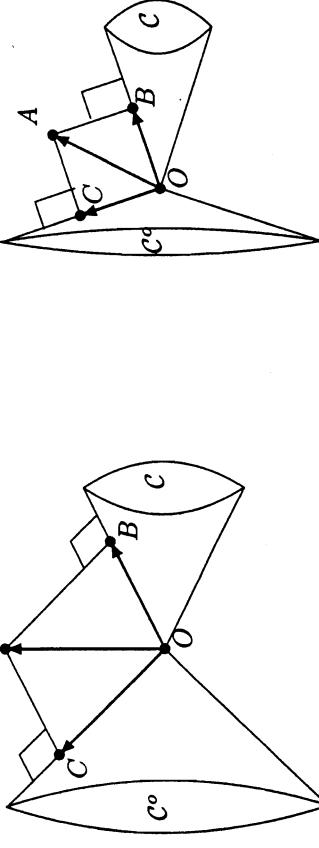


Fig. 3.15 The orthogonal projections of OA onto C and C^o when $V \neq I$.

Proposition 3.12.4 Let C be a closed convex cone and $\mathbf{x} \in \mathbb{R}^p$.

(a) Assume that $\mathbf{x} = \mathbf{y} + \mathbf{z}$ with $\mathbf{y} \in C$, $\mathbf{z} \in C^o$ and $\mathbf{y}^T \mathbf{z} = 0$. Then $\mathbf{y} = \Pi(\mathbf{x}|C)$ and $\mathbf{z} = \Pi(\mathbf{x}|C^o)$.

(b) Conversely, $\mathbf{x} = \Pi(\mathbf{x}|C) + \Pi(\mathbf{x}|C^o)$ and $\Pi(\mathbf{x}|C)^T \Pi(\mathbf{x}|C^o) = 0$.

(c) $\mathbb{R}^p = C \oplus C^o$.

(d) $C^o = C$.

(For (a) and (b) see Hiraiwa-Urruty (1993), pp 120-121 or SW 2.7.5.)

Proof of (a): Let \mathbf{a} be an arbitrary point in C . Then, $\mathbf{a}^T \mathbf{z} \leq 0$, and $\|\mathbf{x} - \mathbf{a}\|^2 = \|\mathbf{x} - \mathbf{y}\|^2 + \|\mathbf{y} - \mathbf{a}\|^2 - 2\mathbf{a}^T \mathbf{z} \geq \|\mathbf{x} - \mathbf{y}\|^2$. Therefore, \mathbf{y} is the point in C closest to \mathbf{x} and hence $\mathbf{y} = \Pi(\mathbf{x}|C)$. Similarly, let $\mathbf{b} \in C^o$. Then $\mathbf{b}^T \mathbf{y} \leq 0$, and $\|\mathbf{x} - \mathbf{b}\|^2 = \|\mathbf{x} - \mathbf{z}\|^2 + \|\mathbf{z} - \mathbf{b}\|^2 - 2\mathbf{b}^T \mathbf{y} \geq \|\mathbf{x} - \mathbf{z}\|^2$. Therefore, \mathbf{z} is the point in C^o closest to \mathbf{x} and hence $\mathbf{z} = \Pi(\mathbf{x}|C^o)$.

Proof of (b): Let \mathbf{x}^* denote $\Pi(\mathbf{x}|C)$ and $\mathbf{y} = (1 + \delta)\mathbf{x}^*$ for $|\delta| < 1$. Then, by Proposition 3.12.3 (on page 114), $(\mathbf{x} - \mathbf{x}^*)^T (\mathbf{y} - \mathbf{x}^*) \leq 0$ for $|\delta| < 1$. Therefore, $\delta(\mathbf{x} - \mathbf{x}^*)^T \mathbf{x}^* \leq 0$ for $|\delta| < 1$. Since δ may be positive or negative, we have that $(\mathbf{x} - \mathbf{x}^*)^T \mathbf{x}^* = 0$. Let \mathbf{u} be an arbitrary point in C and $\mathbf{y} = \mathbf{x}^* + \mathbf{u}$. Then, $\mathbf{y} \in C$ because C is a convex cone. Now, $0 \geq (\mathbf{x} - \mathbf{x}^*)^T (\mathbf{y} - \mathbf{x}^*) = (\mathbf{x} - \mathbf{x}^*)^T \mathbf{u}$. Since \mathbf{u} is arbitrary, we have that $\mathbf{x} - \mathbf{x}^* \in C^o$. Thus, we have shown that $\mathbf{x} = \mathbf{x}^* + (\mathbf{x} - \mathbf{x}^*)$, $\mathbf{x}^* \in C$, $(\mathbf{x} - \mathbf{x}^*) \in C^o$ and $(\mathbf{x} - \mathbf{x}^*)^T \mathbf{x}^* = 0$. Therefore, by part (a), $\mathbf{x} - \mathbf{x}^* = \Pi(\mathbf{x}|C^o)$, and hence $\mathbf{x} = \Pi(\mathbf{x}|C) + \Pi(\mathbf{x}|C^o)$ and $\Pi(\mathbf{x}|C)^T \Pi(\mathbf{x}|C^o) = 0$.

Proof of (c): Follows directly from parts (a) and (b).

Proof of (d): It follows from the definition of polar cones that $C \subset C^{oo}$. To prove the inclusion in the other direction, let $\mathbf{x} \in C^{oo}$ and let $\mathbf{y} + \mathbf{z}$ be the orthogonal decomposition corresponding to $\mathbb{R}^p = C \oplus C^o$. Now, $\mathbf{z}^T \mathbf{z} = \mathbf{z}^T \mathbf{z} + \mathbf{y}^T \mathbf{z} = (\mathbf{y} + \mathbf{z})^T \mathbf{z} = \mathbf{x}^T \mathbf{z} \leq 0$, since $\mathbf{x} \in C^{oo}$ and $\mathbf{z} \in C^o$. It follows that $\mathbf{z}^T \mathbf{z} = 0$ and hence $\mathbf{x} = \mathbf{y} \in C$. Therefore, $C^{oo} \subset C$.

This result is illustrated in Figures 3.15 and 3.16; the first one shows the projections with respect to an arbitrary matrix V and the second shows those for the identity matrix. Let A be an arbitrary point, and \mathbf{x} denote OA . Let $\mathbf{y} = OB$ and $\mathbf{z} = OC$ be the projections of OA onto C and C^o , respectively. Now, the foregoing proposition says that OC and OB are orthogonal, $OBAC$ is a rectangle (with respect to $\langle \cdot, \cdot \rangle_V$), and $OA = OB + OC$. Further, the orthogonal decomposition of OA into OB and OC where B and C lie C and C^o , respectively, is unique.

The following result appears in Rockafellar (1970) on page 146:

Proposition 3.12.5 Let C_1, \dots, C_k be k convex cones in \mathbb{R}^p . Then

$$(C_1 + \dots + C_k)^o = C_1^o + \dots + C_k^o,$$

and $(C_1 \cap \dots \cap C_k)^o = \text{closure of } (C_1^o + \dots + C_k^o)$.

Example: Let $X, P, Q \subset \mathbb{R}^p$, $P \subset Q^o$, and $X = P \oplus Q$. Show that for any $\mathbf{x} \in X$, there exist unique $\mathbf{p} \in P$ and $\mathbf{q} \in Q$ such that $\mathbf{x} = \mathbf{p} + \mathbf{q}$ and $\mathbf{p}^T \mathbf{q} = 0$; further, $\mathbf{p} = \Pi(\mathbf{x}|P)$ and $\mathbf{q} = \Pi(\mathbf{x}|Q)$.

Solution: Let $\mathbf{a} \in P$. Then

$$\begin{aligned} \|\mathbf{x} - \mathbf{a}\|^2 &= \|\mathbf{p} + \mathbf{q} - \mathbf{a}\|^2 = \|\mathbf{p} - \mathbf{a}\|^2 + \|\mathbf{q}\|^2 + 2(\mathbf{p} - \mathbf{a})^T \mathbf{q} \\ &= \|\mathbf{p} - \mathbf{a}\|^2 + \|\mathbf{q}\|^2 - 2\mathbf{a}^T \mathbf{q}, \quad \text{because } \mathbf{p}^T \mathbf{q} = 0 \\ &\geq \|\mathbf{p} - \mathbf{a}\|^2 + \|\mathbf{q}\|^2; \end{aligned} \tag{3.63}$$

the last step follows because $\mathbf{a} \in P$ and $\mathbf{q} \in Q^o$ and hence $\mathbf{a}^T \mathbf{q} \leq 0$. Therefore, $\|\mathbf{x} - \mathbf{a}\|^2$ reaches its minimum over $\mathbf{a} \in P$ when $\mathbf{a} = \mathbf{p}$. Hence, $\mathbf{p} = \Pi(\mathbf{x}|P)$. Similarly, $\mathbf{q} = \Pi(\mathbf{x}|Q)$.

Another important feature of projections onto linear spaces is that if \mathcal{M} and \mathcal{L} are linear subspaces and $\mathcal{M} \subset \mathcal{L}$ then projecting \mathbf{x} onto \mathcal{L} first and then onto \mathcal{M} leads to the same point as projecting \mathbf{x} directly onto \mathcal{M} ; further, this two-stage projection leads to the orthogonal decomposition, $\Pi(\mathbf{x}|\mathcal{L}) = \Pi(\mathbf{x}|\mathcal{M}) + \Pi(\mathbf{x}|\mathcal{L} \cap \mathcal{M}^\perp)$, and hence $\|\Pi(\mathbf{x}|\mathcal{L})\|^2 = \|\Pi(\mathbf{x}|\mathcal{M})\|^2 + \|\Pi(\mathbf{x}|\mathcal{L} \cap \mathcal{M}^\perp)\|^2$. This plays an important role for statistical inference in normal theory linear models. For example, consider the likelihood ratio statistic for testing $\theta \in \mathcal{M}$ against $\theta \in \mathcal{L}$ based on $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{V})$. This LRT statistic is equal to the LRT statistic for testing $\theta = 0$ against $\theta \in \mathcal{L} \cap \mathcal{M}^\perp$, and hence its null distribution is χ_q^2 where $q = \dim(\mathcal{L} \cap \mathcal{M}^\perp) = \dim(\mathcal{L}) - \dim(\mathcal{M})$. Similar results hold for projections onto closed convex cones, C_1 and C_2 , as well provided that $C_1 \subset C_2$ and at least one of them is a linear space. These results are stated below. It is worth noting that if neither C_1 nor C_2 is a linear space then we

cannot obtain orthogonal decompositions in general; it will be seen that this is the root of the difficulties in obtaining least favorable null distributions when the null and alternative parameter spaces are C_1 and C_2 with neither of them being a linear space.

Proposition 3.12.6 Let \mathcal{L} and \mathcal{M} be linear spaces, C be a closed convex cone, $\mathcal{M} \subset C \subset \mathcal{L} \subset \mathbb{R}^p$, and $\mathbf{y} \in \mathbb{R}^p$. Then

$$(a) \Pi\{\Pi(\mathbf{y}|\mathcal{L})|C\} = \Pi(\mathbf{y}|C).$$

$$(b) \Pi\{\Pi(\mathbf{y}|C)|\mathcal{M}\} = \Pi(\mathbf{y}|\mathcal{M}).$$

$$(c) \mathcal{C} = \mathcal{M} \oplus (C \cap \mathcal{M}^\perp).$$

(d-1) Given $\mathbf{x} \in \mathbb{R}^p$, there exist unique $\mathbf{m} \in \mathcal{M}$, $\mathbf{m}_o \in \mathcal{M} \cap C^\perp$, and $\mathbf{c}_o \in C^o$ such that $\mathbf{x} = \mathbf{m} + \mathbf{m}_o + \mathbf{c}_o$ and the components are pairwise orthogonal, namely $\mathbf{m}^T \mathbf{m}_o = \mathbf{m}^T \mathbf{c}_o = \mathbf{m}_o^T \mathbf{c}_o = 0$. Further, the components \mathbf{m} , \mathbf{m}_o and \mathbf{c}_o are the projections of \mathbf{x} onto \mathcal{M} , $\mathcal{C} \cap \mathcal{M}^\perp$ and C^o , respectively.

$$(d-2) \Pi\{\Pi(\mathbf{x} | C) | C \cap \mathcal{M}^\perp\} = \Pi(\mathbf{x} | C \cap \mathcal{M}^\perp).$$

(d-3) For any $\mathbf{x} \in C$, we have

$$\mathbf{x} = \Pi(\mathbf{x} | \mathcal{M}) + \Pi(\mathbf{x} | C \cap \mathcal{M}^\perp) + \Pi(\mathbf{x} | C^o)$$

and the three projections in this decomposition are pairwise orthogonal.

$$(d-4) \mathbb{R}^p = \mathcal{M} \oplus (C \cap \mathcal{M}^\perp) \oplus C^o$$

$$(d-5) (C \cap \mathcal{M}^\perp)^o = C^o \oplus \mathcal{M}.$$

$$(e) \mathcal{L} = \mathcal{C} \oplus (\mathcal{L} \cap \mathcal{C}^o).$$

$$(f) \mathbb{R}^p = \mathcal{C} \oplus (C^o \cap \mathcal{L}) \oplus \mathcal{L}^\perp.$$

$$(g) \mathbf{y} - \Pi(\mathbf{y}|\mathcal{C}) \perp \Pi(\mathbf{y}|\mathcal{M}) - \Pi(\mathbf{y}|C).$$

$$(h) \|\mathbf{y} - \Pi(\mathbf{y}|\mathcal{M})\|^2 - \|\mathbf{y} - \Pi(\mathbf{y}|C)\|^2 = \|\Pi(\mathbf{y}|\mathcal{M}) - \Pi(\mathbf{y}|C)\|^2.$$

$$(i) \Pi(\mathbf{y}|C) - \Pi(\mathbf{y}|\mathcal{M}) = \Pi(\mathbf{y}|C \cap \mathcal{M}^\perp).$$

Proof of (a): Let $\mathcal{L} = \{\mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^k\}$ where \mathbf{X} is a $p \times k$ matrix of rank k , $\mathcal{K} \subset \mathbb{R}^k$ be a closed convex cone such that $C = \mathbf{X}\mathcal{K}$, and $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Then, by using the normal equation, we have that

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 + \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\boldsymbol{\beta}\|^2$$

and $\Pi(\mathbf{X}|\mathcal{L}) = \mathbf{X}\hat{\boldsymbol{\beta}}$. Therefore, the value of $\boldsymbol{\beta} \in \mathcal{K}$, say $\tilde{\boldsymbol{\beta}}$, at which $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ is a minimum is also the same as that at which $\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\boldsymbol{\beta}\|^2$ is a minimum. Since $\Pi(\mathbf{y}|\mathcal{L}) = \mathbf{X}\hat{\boldsymbol{\beta}}$, $\Pi(\mathbf{X}\hat{\boldsymbol{\beta}}|\mathcal{L}) = \mathbf{X}\tilde{\boldsymbol{\beta}}$ and $\Pi(\mathbf{y}|\mathcal{C}) = \mathbf{X}\tilde{\boldsymbol{\beta}}$, the proof follows. One could also show this directly using $\|\mathbf{y} - \mathbf{x}\|^2 = \|\mathbf{y} - \mathbf{P}\mathbf{y}\|^2 + \|\mathbf{P}\mathbf{y} - \mathbf{x}\|^2$ where \mathbf{P} is the projection matrix $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ onto \mathcal{L} , $\mathbf{x} \in \mathcal{L}$, and $\mathbf{y} \in \mathbb{R}^p$.

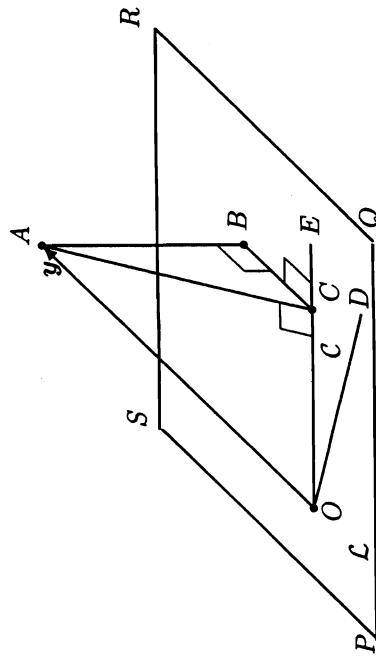


Fig. 3.17 Projections of OA onto the cone C , and the linear space \mathcal{L} , where $C \subset \mathcal{L}$.

Fig. 3.17 illustrates this result; \mathcal{L} is the two-dimensional plane containing PQR ; O is the origin; C is the cone EOD contained in \mathcal{L} ; B is the projection of A onto \mathcal{L} ; and C is the projection of B onto C . Then C is also the projection of A onto C .

Proof of (b) (Rauberias et al. (1986), p 2815): Since $\Pi(\mathbf{y}|C^\perp) = \mathbf{y} - \Pi(\mathbf{y}|C^o)$, we have that $\Pi\{\Pi(\mathbf{y}|C)|\mathcal{M}\} = \Pi\{\mathbf{y} - \Pi(\mathbf{y}|C^o)|\mathcal{M}\} = \Pi\{\Pi(\mathbf{y}|C^o)|\mathcal{M}\}$; the last equality follows since any projection onto any linear space is a linear function. Since $\mathcal{M} \subset C$, we have that $C^o \subset \mathcal{M}^\perp$. Therefore, $\Pi\{\Pi(\mathbf{y}|C^o)|\mathcal{M}\} = 0$, and we conclude that $\Pi\{\Pi(\mathbf{y}|C)|\mathcal{M}\} = \Pi(\mathbf{y}|\mathcal{M})$ (see Fig. 3.10, page 85).

Proof of (c): Let $C^* = \mathcal{M}^\perp \cap C$. First let us show that $C = \{\mathbf{y} + \mathbf{z} : \mathbf{y} \in \mathcal{M}, \mathbf{z} \in C^*\}$. Recall that $\mathbb{R}^p = \mathcal{M} \oplus \mathcal{M}^\perp$. Given $\mathbf{x} \in C$, we can write $\mathbf{x} = \mathbf{y} + \mathbf{z}$, where $\mathbf{y} \in \mathcal{M}$, $\mathbf{z} \in \mathcal{M}^\perp$ and $\mathbf{y}^T \mathbf{z} = 0$. Since $-\mathbf{y} \in \mathcal{M} \subset C$ and C is a closed convex cone, it follows that $\mathbf{z} = \mathbf{x} + (-\mathbf{y}) \in C$. Therefore, $\mathbf{z} \in C^*$ and hence $C \subset \mathcal{M} + C^*$. Since $\mathcal{M} \subset C$, $C^* = \mathcal{M}^\perp \cap C \subset C$ and C is a closed convex cone it follows that $\mathcal{M} + C^* \subset C$. Therefore, $C = \mathcal{M} + C^*$.

The claim that $\mathcal{M} + C^*$ is an orthogonal sum follows directly from $\mathbb{R}^p = \mathcal{M} \oplus \mathcal{M}^\perp$. Let $\mathbf{x} \in C$ and let us start with the decomposition, $\mathbf{x} = \mathbf{y} + \mathbf{z}$ in the previous paragraph where $\mathbf{y} \in \mathcal{M}$, $\mathbf{z} \in \mathcal{M}^\perp$ and $\mathbf{y}^T \mathbf{z} = 0$; the pair $\{\mathbf{y}, \mathbf{z}\}$ is uniquely defined. Since we just proved that \mathbf{z} also belongs to $C \cap \mathcal{M}^\perp$, we have shown the following: $C = \mathcal{M} + C^*$ and given $\mathbf{x} \in C$ there exist unique $\mathbf{y} \in \mathcal{M}$ and $\mathbf{z} \in C^*$ such that $\mathbf{x} = \mathbf{y} + \mathbf{z}$ and $\mathbf{y}^T \mathbf{z} = 0$. Therefore, $C = \mathcal{M} \oplus C^*$.

Remark: Note that \mathbf{z} is the point in \mathcal{M}^\perp closest to \mathbf{x} , and further since $\mathbf{z} \in C^* \subset \mathcal{M}^\perp$ it is also the point in C^* closest to \mathbf{x} as well. Therefore, $\mathbf{z} = \Pi(\mathbf{x} | C^*)$ and $\mathbf{y} = \Pi(\mathbf{x} | \mathcal{M})$.

Proof of (d-I): Let $\mathbf{x} \in \mathbb{R}^p$. Since $\mathbb{R}^p = \mathcal{C} \oplus \mathcal{C}^o$, there exist unique $\mathbf{c} \in \mathcal{C}$ and $\mathbf{c}_0 \in \mathcal{C}^o$ such that $\mathbf{x} = \mathbf{c} + \mathbf{c}_0$ and $\mathbf{c}^T \mathbf{c}_0 = 0$. We also have $\mathbf{c} = \Pi(\mathbf{x}|\mathcal{C})$ and $\mathbf{c}_0 = \Pi(\mathbf{x}|\mathcal{C}^o)$. Since $\mathcal{C} = \mathcal{M} \oplus (\mathcal{C} \cap \mathcal{M}^\perp)$, there exist unique $\mathbf{m} \in \mathcal{M}$ and $\mathbf{m}_0 \in \mathcal{C} \cap \mathcal{M}^\perp$ such that $\mathbf{c} = \mathbf{m} + \mathbf{m}_0$ and $\mathbf{m}^T \mathbf{m}_0 = 0$. We also have

$$\mathbf{m} = \Pi(\mathbf{c} | \mathcal{M}) = \Pi\{\Pi(\mathbf{x} | \mathcal{C}) | \mathcal{M}\} = \Pi(\mathbf{x} | \mathcal{M}) \quad \text{by part (b),} \quad (3.64)$$

$$\mathbf{m}_0 = \Pi(\mathcal{C} \mid \mathcal{C} \cap \mathcal{M}^\perp) = \Pi\{\Pi(\mathbf{x} \mid \mathcal{C}) \mid \mathcal{C} \cap \mathcal{M}^\perp\}. \quad (3.65)$$

Since $\mathcal{M} \subset \mathcal{C}$, we have that $\mathcal{C}^\circ \subset \mathcal{M}^\perp$ and hence $\mathbf{c}_o^T \mathbf{m} = 0$. Now

$$0 = \mathbf{c}_o^T \mathbf{c} = \mathbf{c}_o^T (\mathbf{m} + \mathbf{m}_o) = \mathbf{c}_o^T \mathbf{m} + \mathbf{c}_o^T \mathbf{m}_o = \mathbf{c}_o^T \mathbf{m}_o.$$

Therefore, given $\mathbf{x} \in \mathbb{R}^p$, we have shown that there exist $\mathbf{m} \in \mathcal{M}$, $\mathbf{m}_o \in \mathcal{M} \cap \mathcal{C}^\perp$, and $\mathbf{c}_o \in \mathcal{C}^\circ$ such that $\mathbf{x} = \mathbf{m} + \mathbf{m}_o + \mathbf{c}_o$ and the three components are pairwise orthogonal, (i.e., $\mathbf{m}^T \mathbf{m}_o = \mathbf{m}_o^T \mathbf{c}_o = \mathbf{m}_o^T \mathbf{c}_o = 0$).

To prove uniqueness of the decomposition let $\mathbf{x} \in \mathbb{R}^p$. Suppose that $\mathbf{x} = \mathbf{m}' + \mathbf{m}'_o + \mathbf{c}'_o$ where $\mathbf{m}' \in \mathcal{M}$, $\mathbf{m}'_o \in \mathcal{C} \cap \mathcal{M}^\perp$, $\mathbf{c}'_o \in \mathcal{C}^\circ$, and $\{\mathbf{m}', \mathbf{m}'_o, \mathbf{c}'_o\}$ are pairwise orthogonal. Then, since $\mathbf{m}' \in \mathcal{C}$ and $\mathbf{m}'_o \in \mathcal{C}$ it follows that $\mathbf{m}' + \mathbf{m}'_o \in \mathcal{C}$. It also follows from the pairwise orthogonality of the components that $(\mathbf{m}' + \mathbf{m}'_o)^T \mathbf{c}'_o = 0$. Since $\mathbb{R}^p = \mathcal{C} \oplus \mathcal{C}^\circ$ it follows from the orthogonal decomposition $\mathbf{x} = (\mathbf{m}' + \mathbf{m}'_o) + \mathbf{c}'_o$ that $\mathbf{c}'_o = \Pi(\mathbf{x} \mid \mathcal{C}^\circ) = \mathbf{c}_0$ and $(\mathbf{m}' + \mathbf{m}'_o) = \Pi(\mathbf{x} \mid \mathcal{C})$. Since $\mathcal{C} = \mathcal{M} \oplus (\mathcal{C} \cap \mathcal{M}^\perp)$, $\mathbf{m}' + \mathbf{m}'_o \in \mathcal{C}$, $\mathbf{m}' \in \mathcal{M}$, $\mathbf{m}'_o \in \mathcal{C} \cap \mathcal{M}^\perp$, and $\mathbf{m}^T \mathbf{m}'_o = 0$ it follows that $\mathbf{m}' = \Pi\{\mathbf{m}' + \mathbf{m}'_o \mid \mathcal{M}\} = \Pi\{\Pi(\mathbf{x} \mid \mathcal{C}) \mid \mathcal{M}\} = \Pi(\mathbf{x} \mid \mathcal{M})$. Therefore, to complete the proof of (d-1) it suffices to prove that $\mathbf{m}_o = \Pi(\mathbf{x} \mid \mathcal{C} \cap \mathcal{M}^\perp)$. To this end, let $\mathbf{b} \in \mathcal{C} \cap \mathcal{M}^\perp$. Now,

$$\begin{aligned} & (\mathbf{m} + \mathbf{c}_o)^T (\mathbf{m}_o - \mathbf{b}) = \mathbf{m}^T (\mathbf{m}_o - \mathbf{b}) + \mathbf{c}_o^T (\mathbf{m}_o - \mathbf{b}) \\ & = 0 + \mathbf{c}_o^T (\mathbf{m}_o - \mathbf{b}), \quad \text{because } \mathbf{m}_o - \mathbf{b} \in \mathcal{M}^\perp \text{ and } \mathbf{m} \in \mathcal{M}. \\ & = \mathbf{c}_o^T \mathbf{m} - \mathbf{c}_o^T \mathbf{b} \\ & = 0 - \mathbf{c}_o^T \mathbf{b} \quad \text{because } \mathbf{c}_o \in \mathcal{C}^\circ \subset \mathcal{M}^\perp \text{ and } \mathbf{m} \in \mathcal{M} \end{aligned}$$

The last term is nonnegative because $\mathbf{c}_o \in \mathcal{C}^\circ$ and $\mathbf{b} \in \mathcal{C}$. Therefore, $(\mathbf{m} + \mathbf{c}_o)^T (\mathbf{m}_o - \mathbf{b})$ is nonnegative and reaches its minimum of zero when $\mathbf{b} = \mathbf{m}_0$. Now,

$$\begin{aligned} & \|\mathbf{x} - \mathbf{b}\|^2 = \|\mathbf{m} + \mathbf{m}_o + \mathbf{c}_o - \mathbf{b}\|^2 \\ & = \|\mathbf{m} + \mathbf{c}_o\|^2 + \|\mathbf{m}_o - \mathbf{b}\|^2 + 2(\mathbf{m} + \mathbf{c}_o)^T (\mathbf{m}_o - \mathbf{b}) \\ & \geq \|\mathbf{m} + \mathbf{c}_o\|^2 + \|\mathbf{m} - \mathbf{b}\|^2 \quad \forall \mathbf{b} \in \mathcal{C} \cap \mathcal{M}^\perp, \end{aligned}$$

equality holding when $\mathbf{b} = \mathbf{m}_0$. Hence, $\|\mathbf{x} - \mathbf{b}\|^2$ reaches its minimum over $\mathbf{b} \in \mathcal{C} \cap \mathcal{M}^\perp$ when $\mathbf{b} = \mathbf{m}_o$. Therefore, $\Pi(\mathbf{x} \mid \mathcal{C} \cap \mathcal{M}^\perp) = \mathbf{m}_o = \Pi\{\Pi(\mathbf{x} \mid \mathcal{C}) \mid \mathcal{C} \cap \mathcal{M}^\perp\}$.

Proofs of part (d-2), (d-3), and (d-4): Contained in the proof of (d-1).

Proof of part (d-5): Since $\mathcal{C} \cap \mathcal{M}^\perp \subset \mathcal{C}$ we have $\mathcal{C}^\circ \subset (\mathcal{C} \cap \mathcal{M}^\perp)^\circ$. Since $\mathcal{C} \cap \mathcal{M}^\perp \subset \mathcal{M}^\perp$ it follows that $\mathcal{M} = (\mathcal{M}^\perp)^\perp = (\mathcal{M}^\perp)^\circ \subset (\mathcal{C} \cap \mathcal{M}^\perp)^\circ$. Since $(\mathcal{C} \cap \mathcal{M}^\perp)^\circ$ is a closed convex cone, it is closed under addition, and hence $\mathcal{C}^\circ + \mathcal{M} \subset (\mathcal{C} \cap \mathcal{M}^\perp)^\circ$.

To prove the set containment in the reverse order, first note that $\mathcal{C}^\circ \subset \mathcal{C} + \mathcal{M}$, and hence $(\mathcal{C}^\circ + \mathcal{M})^\circ \subset \mathcal{C}$. Since $\mathcal{M} \subset \mathcal{C} + \mathcal{M}$, it follows that $(\mathcal{C}^\circ + \mathcal{M})^\circ \subset \mathcal{C} + \mathcal{M}$. Therefore, $(\mathcal{C}^\circ + \mathcal{M})^\circ \subset \mathcal{C} \cap \mathcal{M}^\perp$ and $(\mathcal{C} \cap \mathcal{M}^\perp)^\circ \subset \mathcal{C} + \mathcal{M}$.

Therefore, we have proved $(\mathcal{C} \cap \mathcal{M}^\perp)^\circ = \mathcal{C}^\circ + \mathcal{M}$. The proof that this is an orthogonal sum follows from part (d-4).

Proof of (e): The proof is similar to that for part (c) and is based on $\mathbb{R}^p = \mathcal{C} \oplus \mathcal{C}^\circ$.

Proof of (f): Given $\mathbf{y} \in \mathbb{R}^p$ we have the unique orthogonal decomposition, $\mathbf{y} = \Pi(\mathbf{y} | \mathcal{L}) + \Pi(\mathbf{y} | \mathcal{L}^\perp)$. Since $\mathcal{L} = \mathcal{C} \oplus (\mathcal{C}^\circ \cap \mathcal{L})$, we also have the unique orthogonal decomposition

$$\begin{aligned} \Pi(\mathbf{y} | \mathcal{L}) &= \Pi\{\Pi(\mathbf{y} | \mathcal{L}) | \mathcal{C}\} + \Pi\{\Pi(\mathbf{y} | \mathcal{L}) | \mathcal{C}^\circ \cap \mathcal{L}\} \\ &= \Pi(\mathbf{y} | \mathcal{C}) + \Pi(\mathbf{y} | \mathcal{C}^\circ \cap \mathcal{L}) \quad \text{by part (a).} \end{aligned} \quad (3.66)$$

Therefore, we have the orthogonal decomposition,

$$\begin{aligned} \mathbf{y} &= \Pi(\mathbf{y} | \mathcal{C}) + \Pi(\mathbf{y} | \mathcal{C}^\circ \cap \mathcal{L}) + \Pi(\mathbf{y} | \mathcal{L}^\perp). \\ \mathbf{y} &= \Pi(\mathbf{y} | \mathcal{C}) + \Pi(\mathbf{y} | \mathcal{C} \cap \mathcal{L}) \oplus \mathcal{L}^\perp. \end{aligned}$$

Uniqueness of the decomposition can also be established as for part (d-1). Hence, $\mathbb{R}^p = \mathcal{C} \oplus (\mathcal{C}^\circ \cap \mathcal{L}) \oplus \mathcal{L}^\perp$.
Proofs of (g) and (h): It is easily verified that $\mathcal{M}^\perp = \mathcal{M}^\circ$, and if $\mathcal{C}_1 \subset \mathcal{C}_2$ then $\mathcal{C}_2^\circ \subset \mathcal{C}_1^\circ$. Since $\mathcal{M} \subset \mathcal{C}$ we have that $\mathcal{C}^\circ \subset \mathcal{M}^\perp$ and $\Pi(\mathbf{y} | \mathcal{C}^\circ)^T \Pi(\mathbf{y} | \mathcal{M}) = 0$. Note that $\Pi(\mathbf{y} | \mathcal{C}^\circ) \in \mathcal{C}^\circ \subset \mathcal{M}^\perp$. Therefore, $\Pi(\mathbf{y} | \mathcal{C}^\circ) \perp \mathcal{M}$. Now,

$$\begin{aligned} & \{\mathbf{y} - \Pi(\mathbf{y} | \mathcal{C})\}^T \{\Pi(\mathbf{y} | \mathcal{C}) - \Pi(\mathbf{y} | \mathcal{M})\} \\ & = \Pi(\mathbf{y} | \mathcal{C}^\circ)^T \Pi(\mathbf{y} | \mathcal{C}) - \Pi(\mathbf{y} | \mathcal{C}^\circ)^T \Pi(\mathbf{y} | \mathcal{M}) = 0. \end{aligned} \quad (3.67)$$

Therefore, $\|\mathbf{y} - \Pi(\mathbf{y} | \mathcal{M})\|^2 = \|\mathbf{y} - \Pi(\mathbf{y} | \mathcal{C})\|^2 + \|\Pi(\mathbf{y} | \mathcal{C}) - \Pi(\mathbf{y} | \mathcal{M})\|^2$.

Proof of (i):

$$\begin{aligned} \Pi(\mathbf{x} \mid \mathcal{C}) &= \Pi\{\Pi(\mathbf{x} \mid \mathcal{C}) \mid \mathcal{M}\} + \Pi\{\Pi(\mathbf{x} \mid \mathcal{C}) \mid \mathcal{C} \cap \mathcal{M}^\perp\} \\ &= \Pi(\mathbf{x} \mid \mathcal{M}) + \Pi(\mathbf{x} \mid \mathcal{C} \cap \mathcal{M}^\perp) \quad \text{by parts (b) and (c)} \end{aligned} \quad (3.68)$$

The Fig. 3.17 illustrates the main ideas of this proof in three-dimensions. ■

A consequence of the orthogonal sum $\mathcal{C} = \mathcal{M} \oplus (\mathcal{C} \cap \mathcal{M}^\perp)$ in part(c) of the foregoing proposition is that any $\mathbf{y} \in \mathcal{C}$ has a unique decomposition of the form $\mathbf{y} = \mathbf{a} + \mathbf{b}$, where $\mathbf{a} \in \mathcal{M}$, $\mathbf{b} \in \mathcal{C} \cap \mathcal{M}^\perp$, $\mathbf{a}^T \mathbf{b} = 0$; further for such a decomposition, $\mathbf{a} = \Pi(\mathbf{y} \mid \mathcal{M})$, and $\mathbf{b} = \Pi(\mathbf{y} \mid \mathcal{C} \cap \mathcal{M}^\perp)$. Similarly, a consequence of the orthogonal sum $\mathbb{R}^p = \mathcal{M} \oplus (\mathcal{C} \cap \mathcal{M}^\perp) \oplus \mathcal{C}^\circ$ is that any $\mathbf{y} \in \mathbb{R}^p$ has a unique decomposition of the form $\mathbf{y} = \mathbf{a} + \mathbf{b} + \mathbf{c}$, where $\mathbf{a} \in \mathcal{M}$, $\mathbf{b} \in \mathcal{C} \cap \mathcal{M}^\perp$, $\mathbf{c} \in \mathcal{C}^\circ$, $\mathbf{a}^T \mathbf{b} = 0$, $\mathbf{a}^T \mathbf{c} = 0$, $\mathbf{b}^T \mathbf{c} = 0$; and, further, for such a decomposition we have that $\mathbf{a} = \Pi(\mathbf{y} \mid \mathcal{M})$, $\mathbf{b} = \Pi(\mathbf{y} \mid \mathcal{C} \cap \mathcal{M}^\perp)$ and $\mathbf{c} = \Pi(\mathbf{y} \mid \mathcal{C}^\circ)$.

In contrast to the decompositions in linear subspaces, the decomposition $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$, with $\mathbf{x}_1 \in \mathcal{C}$ and $\mathbf{x}_2 \in \mathcal{C}^\circ$ is not unique when \mathcal{C} is a closed convex cone unless we impose the requirement $\mathbf{x}_1^T \mathbf{x}_2 = 0$; however, the decomposition is optimal in the sense stated below (Hiriart-Urruty and Lemaréchal (1993), pp 118, 121).

Proposition 3.12.7 Let C be a closed convex cone $\mathbf{x} \in \mathbb{R}^p$. Then we have the following:

- (a) $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2, \mathbf{x}_1 \in C, \mathbf{x}_2 \in C^\circ \Rightarrow \|\mathbf{x}_1\| \geq \|\Pi(\mathbf{x}|C)\| \text{ and } \|\mathbf{x}_2\| \geq \|\Pi(\mathbf{x}|C^\circ)\|$.
- (b) $\|\Pi(\mathbf{x}_1|C) - \Pi(\mathbf{x}_2|C)\|^2 \leq \{\Pi(\mathbf{x}_1|C) - \Pi(\mathbf{x}_2|C)\}^T (\mathbf{x}_1 - \mathbf{x}_2), \text{ for } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^p$.
- (c) $\|\Pi(\mathbf{x}_1|C) - \Pi(\mathbf{x}_2|C)\| \leq \|(\mathbf{x}_1 - \mathbf{x}_2)\|, \text{ for } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^p$. ■

3.12.3 Polyhedral Cones

Let $\mathbf{a}_1, \dots, \mathbf{a}_q$ be q points in \mathbb{R}^p and $\mathcal{P} = \{\mathbf{x} \in \mathbb{R}^p : \mathbf{a}_i^T \mathbf{x} \geq 0 \text{ for } i = 1, \dots, q\}$. Then \mathcal{P} is a closed convex cone and it is called a *polyhedral cone*. Note that \mathcal{P} is the intersection of the half-spaces, $\{\mathbf{x} : \mathbf{a}_1^T \mathbf{x} \geq 0\}, \dots, \{\mathbf{x} : \mathbf{a}_q^T \mathbf{x} \geq 0\}$. With the $p \times q$ matrix \mathbf{A} defined as $[\mathbf{a}_1, \dots, \mathbf{a}_q]$, we may express \mathcal{P} as $\{\mathbf{x} \in \mathbb{R}^p : \mathbf{A}^T \mathbf{x} \geq 0\}$. Because the equality constraint $\mathbf{a}^T \mathbf{x} = 0$ is equivalent to the two inequality constraints $\mathbf{a}^T \mathbf{x} \geq 0$ and $(-\mathbf{a})^T \mathbf{x} \geq 0$, the set of constraints $\mathbf{A}^T \mathbf{x} \geq 0$ may contain equality constraints as well. We will say that the set of constraints $\mathbf{A}^T \mathbf{x} \geq 0$ defining \mathcal{P} is *tight* if \mathcal{P} cannot be defined using a submatrix of \mathbf{A} with fewer columns than \mathbf{A} ; in other words, \mathbf{A} has no redundant columns, and hence no redundant constraints. Without loss of generality, we shall always assume that the set of constraints defining \mathcal{P} is tight.

Let S be a subset of \mathbb{R}^p . Then the cone generated by S consists of elements of the form $w_1 \mathbf{b}_1 + \dots + w_m \mathbf{b}_m$ where $\{\mathbf{b}_1, \dots, \mathbf{b}_m\} \subset S$, $w_i \geq 0$, for $i = 1, \dots, m$; it is also called the *conical hull* or *positive hull* of S . The elements of S are called *generators* of the cone generated by S . The conical hull is the smallest cone that contains S . The set consisting of elements of the form $w_1 \mathbf{b}_1 + \dots + w_m \mathbf{b}_m$ where $\{\mathbf{b}_1, \dots, \mathbf{b}_m\} \subset S$, $\sum w_i = 1$, and $w_i \geq 0$, for $i = 1, \dots, m$, is called the *convex hull* of S . The convex hull is the smallest convex set containing S . An explicit formula for the relationship between a polyhedral cone and its polar cone is given in the next result.

Proposition 3.12.8 Let $\mathbf{a}_1, \dots, \mathbf{a}_q$ be q points in \mathbb{R}^p and let $C = \{w_1 \mathbf{a}_1 + \dots + w_q \mathbf{a}_q : w_i \geq 0, i = 1, \dots, q\}$, the cone generated by $\{\mathbf{a}_1, \dots, \mathbf{a}_q\}$. Then $C^\circ = \{\mathbf{x} \in \mathbb{R}^p : \mathbf{a}_i^T \mathbf{x} \leq 0 \text{ for } i = 1, \dots, q\}$. Consequently, if $\mathcal{P} = \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} \geq 0, \text{ for } i = 1, \dots, q\}$ then $\mathcal{P}^\circ = \{w_1 \mathbf{a}_1 + \dots + w_q \mathbf{a}_q : w_i \leq 0, i = 1, \dots, q\}$. [For example see, Hiriart-Urruty and Lemaréchal (1993), p 119].

Proof. Let $\mathbf{x} \in C^\circ$. Then $\mathbf{x}^T(w_1 \mathbf{a}_1 + \dots + w_q \mathbf{a}_q) \leq 0$ for $(w_1, \dots, w_q) \geq 0$. Now, by choosing $w_i > 0$ and the $w_j = 0$ for $j \neq i$, we have that $\mathbf{a}_i^T \mathbf{x} \leq 0$ for $i = 1, \dots, q$. Therefore, $\mathbf{x} \in \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} \leq 0, \text{ for } i = 1, \dots, q\}$. To prove the set containment in the opposite direction, let $\mathbf{x} \in \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} \leq 0, \text{ for } i = 1, \dots, q\}$ and $\mathbf{z} \in C$. Then $\mathbf{z} = (w_1 \mathbf{a}_1 + \dots + w_q \mathbf{a}_q)$ for some $w_i \geq 0$ for $i = 1, \dots, q$. Now, $\mathbf{x}^T \mathbf{z} = w_1 \mathbf{a}_1^T \mathbf{x} + \dots + w_q \mathbf{a}_q^T \mathbf{x} \leq 0$. Therefore, $\mathbf{x} \in C^\circ$. The second part of the proposition follows from the first by substituting $\mathcal{P} = -C^\circ$ and noting that $\mathcal{P}^\circ = -C^{\circ\circ} = -C$. ■

The following restatement of the foregoing proposition in matrix notation is useful.

Corollary 3.12.9 Let $\mathbf{a}_1, \dots, \mathbf{a}_q$ be q points in \mathbb{R}^p and let $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_q]$.

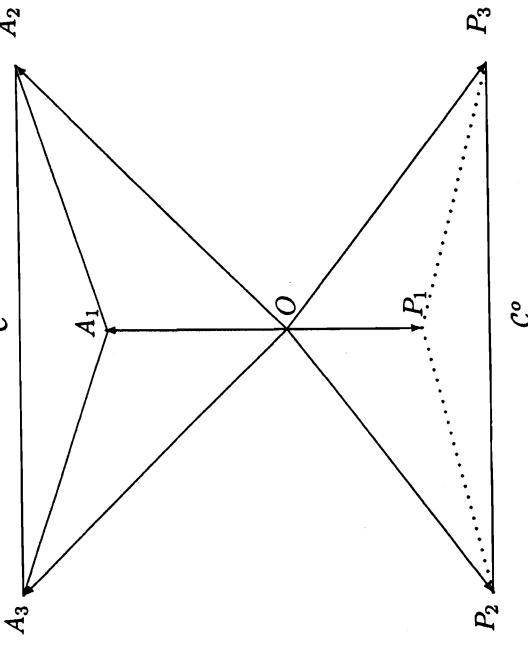


Fig. 3.18 A polyhedral C , and its polar cone C°

1. If $\mathcal{P} = \{\mathbf{x} : \mathbf{A}^T \mathbf{x} \geq 0\}$ then $\mathcal{P}^\circ = \{\mathbf{y} : \mathbf{y} \leq 0\}$.

2. If $C = \{\mathbf{Ax} : \mathbf{x} \geq 0\}$ then $C^\circ = \{\mathbf{y} : \mathbf{A}^T \mathbf{y} \leq 0\}$. ■

Fig. 3.18 provides a representation of this relationship between a polyhedral and its polar in three-dimensions. Let OA_1, OA_2 , and OA_3 correspond to $\mathbf{a}_1, \mathbf{a}_2$, and \mathbf{a}_3 , respectively, and $C = \{\mathbf{Ax} : \mathbf{x} \geq 0\}$. Let $OA_1 \perp OP_2 P_3, OA_2 \perp OP_1 P_3$, and $OA_3 \perp OP_1 P_2$. Then $OP_1 \perp OA_2 A_3, OP_2 \perp OA_1 A_3$, and $OP_3 \perp OA_1 A_2$. The polyhedrals $OA_1 A_2 A_3$ and $OP_1 P_2 P_3$ are C and C° , respectively. Note that every boundary plane of C corresponds to an edge of C° , the two being orthogonal to each other; similarly, every edge of C corresponds to a boundary plane of C° , again the two being orthogonal to each other. This correspondence defines the function Φ in Proposition 3.12.11.

Now we state two useful results on polyhedral cones; for example see SW 2.8.6 and 2.8.8.

Proposition 3.12.10 (a) Minkowski's theorem: For every polyhedral cone \mathcal{P} there exist a_1, \dots, a_q in \mathbb{R}^p such that $\mathcal{P} = \{w_1 a_1 + \dots + w_q a_q : w_i \geq 0 \text{ for } i = 1, \dots, q\}$.

(b) Weyl's theorem: Every closed convex cone with finitely many generators is a polyhedral cone. ■

Let $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_q)$ be a $p \times q$ matrix, and $\mathcal{P} = \{\mathbf{x} : \mathbf{A}^T \mathbf{x} \geq 0\}$. Let $J = \{1, \dots, q\} \setminus I$ be a subset of $\{1, \dots, q\}$; J may be empty. Let $I = \{1, \dots, q\} \setminus J$. Let \mathbf{A}_J and \mathbf{A}_I denote the matrices with their columns being $\{\mathbf{a}_j : j \in J\}$ and $\{\mathbf{a}_i : i \in I\}$ respectively. The set $F_J = \{\mathbf{x} : \mathbf{A}_J^T \mathbf{x} = 0, \mathbf{A}_I^T \mathbf{x} \geq 0\}$ is called a *face* of \mathcal{P} ; without loss of generality, it is assumed that the set of constraints defining F_J

There is an important relationship between the faces of a polyhedral cone and those of its polar cone. This was useful for deriving the chi-bar-square distribution for a general closed convex cone; however, these relationships are not required to understand the proofs in the form they are presented in this monograph. For completeness, we state these results below; for example, see SW 2.13.2, 2.13.3.

Proposition 3.12.11 *Let $\mathcal{F}(\mathcal{P}) = \{F : F \text{ is a face of } \mathcal{P} \text{ and } F \text{ is nonempty}\}$; it is the collection of faces of \mathcal{P} . Then*

1. $F^\perp \cap \mathcal{P}^\circ$ is a face of \mathcal{P}° .

2. For any given face F of \mathcal{P} , let Φ denote the operation of choosing $F^\perp \cap \mathcal{P}^\circ$. Then $\mathcal{F}(\mathcal{P}) \xrightarrow{\Phi} \mathcal{F}(\mathcal{P}^\circ) \xrightarrow{\Phi} \mathcal{F}(\mathcal{P}^{\circ\circ}) = \mathcal{F}(\mathcal{P})$ and $\Phi \Phi(F) = F$.

3. $\langle F^\perp \cap \mathcal{P}^\circ \rangle = \langle F \rangle^\perp$ where $\langle A \rangle^\perp$ denotes the linear space spanned by A . Further, every face of \mathcal{P} is uniquely associated with the face $F^\perp \cap \mathcal{P}^\circ$ of \mathcal{P}° , and any face of \mathcal{P}° is of the form $F^\perp \cap \mathcal{P}^\circ$ for some face F of \mathcal{P} . ■

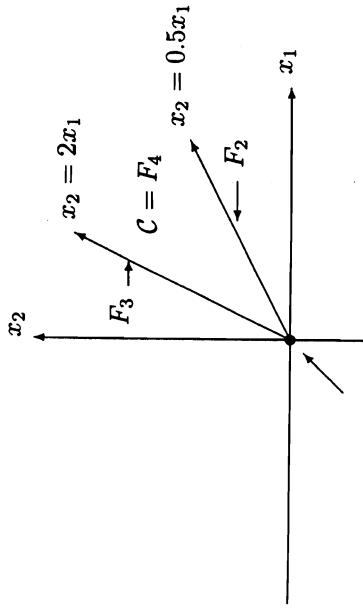


Fig. 3.19 The faces of the polyhedral $\{x : 2x_1 - \frac{1}{2}x_2 \geq 0, x_2 - \frac{1}{2}x_1 \geq 0\}$.

is tight. Now, the *relative interior* of F_J , denoted, $ri(F_J)$, is defined as the interior of F_J with respect to the relative topology induced in the linear space spanned by F_J ; it may be verified that the relative interior of a face is obtained by replacing the inequality constraints defining F_J by their corresponding strict inequalities; thus, $ri(F_J) = \{x : A_J^T x = 0, A_J^T x > 0\}$. The next example illustrates these definitions.

Example: Let $C = \{x : 2x_1 - x_2 \geq 0, x_2 - \frac{1}{2}x_1 \geq 0\}$.

Then C is a polyhedral cone and its faces are (see Fig. 3.19)

$$\begin{aligned} F_1 &= \{x : 2x_1 - x_2 = 0, x_2 - 2^{-1}x_1 = 0\}, \\ F_2 &= \{x : 2x_1 - x_2 \geq 0, x_2 - 2^{-1}x_1 = 0\}, \\ F_3 &= \{x : 2x_1 - x_2 = 0, x_2 - 2^{-1}x_1 \geq 0\}, \\ F_4 &= \{x : 2x_1 - x_2 \geq 0, x_2 - 2^{-1}x_1 \geq 0\}, \end{aligned}$$

Descriptions of the faces of C and their relative interiors are:

$$\begin{aligned} ri(F_1) &= \text{the vertex,} \\ F_2 &= \text{the lower boundary of } C \quad ri(F_2) = F_2 \text{ excluding the vertex} \\ F_3 &= \text{the upper boundary of } C \quad ri(F_3) = F_3 \text{ excluding the vertex} \\ F_4 &= C \text{ excluding its boundary, } F_2 \cup F_3. \end{aligned}$$

Note that the polyhedral \mathcal{P} is the union of sets of the form $ri(F_J)$, each of which is an open set in a linear subspace, and hence does not include its boundary. A main feature of statistical inference problems under inequality constraints is that the parameter space includes the boundary. An advantage of $ri(F_J)$ is that the parameter space \mathcal{P} that includes its boundary is partitioned into sets of the form $ri(F_J)$, each of which is open with respect to the linear space spanned by $ri(F_J)$ and has a linear structure. Consequently, we can apply well-known results for open parameter spaces to a polyhedral \mathcal{P} as well by considering each member of the partition of \mathcal{P} separately.

3.13 APPENDIX 2: PROOFS

Proof of Theorem 3.4.2: To prove this theorem it is convenient to establish several lemmas that are of independent interest.

Lemma 3.13.1 .

1. Suppose that $Z \sim N(0, I)$ and let $\mathcal{P} \subset \mathbb{R}^p$ be a cone. Then the direction and length of Z are independent. Consequently, $\|Z\|$ and $Z/\|Z\|$ are independent random variables and $\{\|Z\| \geq c\}$ and $\{Z \in \mathcal{P}\}$ are independent events.
2. More generally, suppose that the distribution of X is orthogonally invariant (i.e., X and PX have the same distribution for any orthogonal matrix P). Then $\{\|Z\| \geq c\}$ and $\{Z \in \mathcal{P}\}$ are independent events and $Z/\|Z\|$ is distributed uniformly on the unit sphere.

Proof: First, let us transform Z into polar coordinates: $Z_1 = R \sin \theta_1, Z_2 = R \cos \theta_1 \sin \theta_2, Z_3 = R \cos \theta_1 \cos \theta_2 \sin \theta_3, \dots, Z_{p-1} = R \cos \theta_1 \cos \theta_2 \dots \cos \theta_{p-2} \sin \theta_{p-1}, Z_p = R \cos \theta_1 \cos \theta_2 \dots \cos \theta_{p-2} \cos \theta_{p-1}$, where $-(\pi/2) < \theta_i \leq (\pi/2)$ for $i = 1, \dots, p-2$ and $-\pi < \theta_{p-1} \leq \pi$. Then $R^2 = \|Z\|^2$ and the Jacobian of the transformation is $R^{p-1} \cos^{p-2} \theta_1 \cos^{p-3} \theta_2 \dots \cos \theta_{p-2}$ (see Anderson (1984), page 280). Now, since the density of Z at z is a function of $\|z\|$ only, it follows that the density of $(R, \theta_1, \dots, \theta_{p-1})$ factors into two components, one is a function of R only and the other is a function of $(\theta_1, \dots, \theta_{p-1})$. Now the proof follows because $\{\|Z\| \geq c\}$ depends on $(\theta_1, \dots, \theta_{p-1})$ only. A similar result also appeared in Kudo (1963). The second result is stated in Perlman (1969) and is essentially the same as the first. ■

Lemma 3.13.2 *Let $\mathcal{P} \subset \mathbb{R}^p$ be a polyhedral cone, $y \in \mathbb{R}^p$, F be a face of \mathcal{P} , L denote the linear space spanned by F , and P denote the projection matrix onto L . Let the statements (A), (B) and (C) be defined as follows:*

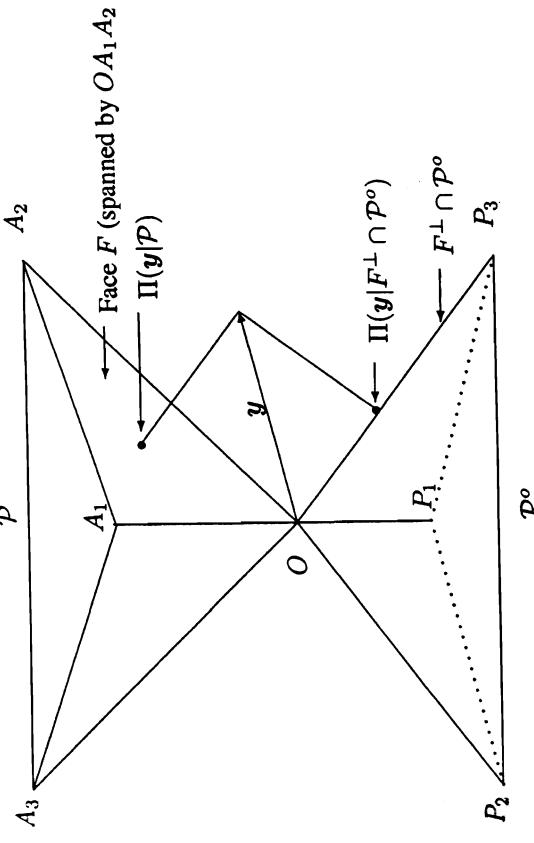


Fig. 3.20 Projections onto F and $F^\perp \cap P^\circ$. (3.69)

Proof of (B) \Rightarrow (C): Assume that (B) is true. Since $\Pi(\mathbf{y}|\mathcal{L}) = \mathbf{P}\mathbf{y} \in ri(F) \subset \mathcal{P}$, $(\mathbf{I} - \mathbf{P})\mathbf{y} = \Pi(\mathbf{y}|\mathcal{L}^\perp) = F^\perp \cap P^\circ \subset \mathcal{P}^\circ$ and $\Pi(\mathbf{y}|\mathcal{L})^T \Pi(\mathbf{y}|\mathcal{L}^\perp) = 0$ it follows that the two orthogonal decompositions in (3.69) are identical, and hence $\mathbf{P}\mathbf{y} = \Pi(\mathbf{y}|\mathcal{L}) = \Pi(\mathbf{y}|\mathcal{P})$ and $(\mathbf{I} - \mathbf{P})\mathbf{y} = \Pi(\mathbf{y}|\mathcal{L}^\perp) = \Pi(\mathbf{y}|\mathcal{P}^\circ)$; this completes the proof of (C).

Proof of (B) \Rightarrow (A): Follows since $\Pi(\mathbf{y}|\mathcal{P}) = \mathbf{P}\mathbf{y} \in ri(F)$. (see Fig. 3.20).

Proof of (A) \Rightarrow (B): Assume that (A) is true. Let \mathbf{y}^* denote $\Pi(\mathbf{y}|\mathcal{P})$. Then $(\mathbf{y} - \mathbf{y}^*)^T(\mathbf{x} - \mathbf{y}^*) \leq 0$ for every $\mathbf{x} \in \mathcal{P}$ (by Proposition 3.12.3 on page 114). Suppose that $(\mathbf{y} - \mathbf{y}^*)^T(\mathbf{x} - \mathbf{y}^*)$ is negative for some $\mathbf{x} \in ri(F)$. Then $g(t) = (\mathbf{y} - \mathbf{y}^*)^T(\mathbf{t} - \mathbf{y}^*)$ is negative for values of t in a small neighborhood of x contained in $ri(F)$. Therefore, there is an $\mathbf{x}_1 \in ri(F)$ such that $(\mathbf{y} - \mathbf{y}^*)^T(\mathbf{x}_1 - \mathbf{y}^*) < 0$, and hence \mathbf{y} is closer to \mathbf{x}_1 in \mathcal{P} than to \mathbf{y}^* ; this is a contradiction because \mathbf{y}^* is the closest point in \mathcal{P} . Therefore, $(\mathbf{y} - \mathbf{y}^*)^T(\mathbf{x} - \mathbf{y}^*) = 0$ for every \mathbf{x} in $ri(F)$ and $\mathbf{y} - \mathbf{y}^*$ is orthogonal to \mathcal{L} . Since $\mathbf{y}^* \in ri(F)$ it follows that $\mathbf{y}^* \in \mathcal{L}$. Therefore, $\mathbf{P}\mathbf{y} = \Pi(\mathbf{y}|\mathcal{L}) = \mathbf{y}^* \in ri(F)$. Further, $\mathbf{y} - \mathbf{y}^* = \Pi(\mathbf{y}|\mathcal{P}^\circ) \in \mathcal{P}^\circ$ by (3.69) and $\mathbf{y} - \mathbf{y}^* = (\mathbf{I} - \mathbf{P})\mathbf{y} = \Pi(\mathbf{y}|\mathcal{L}^\perp) \in \mathcal{L}^\perp \subset F^\perp$. Therefore, $\mathbf{y} - \mathbf{y}^* \in F^\perp \cap P^\circ$. This completes the proof of (A) \Rightarrow (B). ■

(A) $\Pi(\mathbf{y}|\mathcal{P}) \in ri(F)$.

(B) $\mathbf{P}\mathbf{y} \in ri(F)$ and $(\mathbf{I} - \mathbf{P})\mathbf{y} \in F^\perp \cap P^\circ$.

(C) $\mathbf{P}\mathbf{y} = \Pi(\mathbf{y}|\mathcal{P})$ and $(\mathbf{I} - \mathbf{P})\mathbf{y} = \Pi(\mathbf{y}|\mathcal{P}^\circ)$.

Then, (A) is true if and only if (B) is true. If (B) is true then (C) is true.

Proof: Let us recall that $\mathbb{R}^p = \mathcal{L} \oplus \mathcal{L}^\perp$ where \oplus is an orthogonal sum. Consequently, $\mathbf{y} \in \mathbb{R}^p$ has a unique orthogonal decomposition of the form $\mathbf{y} = \mathbf{a} + \mathbf{b}$, where $\mathbf{a} \in \mathcal{L}$, $\mathbf{b} \in \mathcal{L}^\perp$ and $\mathbf{a}^T \mathbf{b} = 0$; further, $\mathbf{a} = \Pi(\mathbf{y}|\mathcal{L}) = \mathbf{P}\mathbf{y}$ and $\mathbf{b} = \Pi(\mathbf{y}|\mathcal{L}^\perp) = (\mathbf{I} - \mathbf{P})\mathbf{y}$. A similar result holds for closed convex cones as well: $\mathbb{R}^p = \mathcal{P} = \mathcal{P} \oplus \mathcal{P}^\circ$ is an orthogonal sum. Consequently, $\mathbf{y} \in \mathbb{R}^p$ has a unique orthogonal decomposition of the form $\mathbf{y} = \mathbf{p} + \mathbf{q}$, where $\mathbf{p} \in \mathcal{P}$, $\mathbf{q} \in \mathcal{P}^\circ$, and $\mathbf{p}^T \mathbf{q} = 0$; further, $\mathbf{p} = \Pi(\mathbf{y}|\mathcal{P})$ and $\mathbf{q} = \Pi(\mathbf{y}|\mathcal{P}^\circ)$. Consequently, if $\mathbf{y} = \mathbf{p} + \mathbf{q}$ where $\mathbf{p} \in \mathcal{P}$, $\mathbf{q} \in \mathcal{P}^\circ$, and $\mathbf{p}^T \mathbf{q} = 0$, then $\mathbf{p} = \Pi(\mathbf{y}|\mathcal{P})$ and $\mathbf{q} = \Pi(\mathbf{y}|\mathcal{P}^\circ)$. Let us write down the two unique orthogonal decompositions of \mathbf{y} :

$$\mathbf{y} = \Pi(\mathbf{y}|\mathcal{L}) + \Pi(\mathbf{y}|\mathcal{L}^\perp) = \Pi(\mathbf{y}|\mathcal{P}) + \Pi(\mathbf{y}|\mathcal{P}^\circ) \quad (3.69)$$

Lemma 3.13.3 Let $\mathbf{Y} \sim N(\mathbf{0}, \mathbf{I})$, \mathbf{P} be a projection matrix, F be a face of \mathcal{P} , \mathbf{P} be the projection matrix onto the linear space spanned by F , and $r = rank(\mathbf{P})$. Then we have the following.

- (I) $pr\{\|\Pi(\mathbf{Y}|\mathcal{P})\|^2 \geq c\} \geq pr\{\Pi(\mathbf{Y}|F)\} = pr\{\chi_r^2 \geq c\}$.
- (II) Conditional on $\Pi(\mathbf{Y}|\mathcal{P}) \in ri(F)$, $\|\Pi(\mathbf{y}|\mathcal{P})\|^2$ and $\|\Pi(\mathbf{y}|\mathcal{P}^\circ)\|^2$ are independent

and are distributed as χ_r^2 and χ_{p-r}^2 , respectively. Hence

$$\begin{aligned} pr\{\|\Pi(\mathbf{Y}|\mathcal{P})\|^2 \geq c_1, \|\Pi(\mathbf{Y}|\mathcal{P}^\circ)\|^2 \geq c_2 | \Pi(\mathbf{Y}|F) \in ri(F)\} \\ = pr\{\chi_r^2 \geq c_1\} pr\{\chi_{p-r}^2 \geq c_2\}. \end{aligned} \quad (3.70)$$

Lemma 3.13.3 Let $\mathbf{Y} \sim N(\mathbf{0}, \mathbf{I})$, \mathbf{P} be a projection matrix, $r = rank(\mathbf{P})$ and $\mathcal{C} \subset \mathbb{R}^p$ be a nonempty cone. Then the distribution of $\|\mathbf{PY}\|^2$, conditional on $\mathbf{PY} \in \mathcal{C}$, is χ_r^2 .

Proof: Let $\mathbf{P} = UDU^T$ where U is an orthogonal matrix and D is a diagonal matrix that we assume, without loss of generality, to be $diag(1, \dots, 1, 0, \dots, 0)$ with each of the first r diagonal elements equal to 1 and the rest being zeros. Let $\mathbf{Z} = \mathbf{U}^T \mathbf{Y}$, and $\mathbf{Z}_a = (Z_1, \dots, Z_r)^T$. Then $\mathbf{Z} \sim N_p(\mathbf{0}, \mathbf{I})$ and $\mathbf{Z}_a \sim N_r(\mathbf{0}, \mathbf{I})$. Further,

Proof:

- (I) Since \mathbf{PY} and $(\mathbf{I} - \mathbf{P})\mathbf{Y}$ are independent and $ri(F)$ is a cone, it follows from the previous two lemmas that $pr\{\|\Pi(\mathbf{Y}|\mathcal{P})\|^2 \geq c | \Pi(\mathbf{Y}|F) \in ri(F)\} = pr\{\|\mathbf{PY}\|^2 \geq c | \Pi(\mathbf{Y}|\mathcal{P}) \in ri(F)\} = pr\{\|\mathbf{PY}\|^2 \geq c | \mathbf{PY} \in ri(F), (\mathbf{I} - \mathbf{P})\mathbf{Y} \in F^\perp \cap \mathcal{P}^\circ\} = pr\{\|\mathbf{PY}\|^2 \geq c | \mathbf{PY} \in ri(F)\} = pr\{\chi_r^2 \geq c\}$.
- (II) By arguments similar to those for the previous part and the fact that $F^\perp \cap \mathcal{P}^\circ$ is

a cone, we have the following:

$$\begin{aligned} \text{pr}\{\|\Pi(Y|\mathcal{P})\|^2 \geq c_1, \|\Pi(Y|\mathcal{P}^\circ)\|^2 \geq c_2 | \Pi(Y|\mathcal{P}) \in ri(F)\} \\ = \text{pr}\{\|\mathbf{P}Y\|^2 \geq c_1, \|\mathbf{(I-P)Y}\|^2 \geq c_2 | \mathbf{PY} \in ri(F), \mathbf{(I-P)Y} \in F^\perp \cap \mathcal{P}^\circ\} \\ = \text{pr}\{\|\mathbf{PY}\|^2 \geq c_1 | \mathbf{PY} \in ri(F)\} \times \\ \text{pr}\{\|\mathbf{(I-P)Y}\|^2 \geq c_2 | (\mathbf{I-P)Y} \in F^\perp \cap \mathcal{P}^\circ\} \\ = \text{pr}(\chi_r^2 \geq c_1) \text{pr}(\chi_{p-r}^2 \geq c_2). \blacksquare \end{aligned}$$

Lemma 3.13.5 Let \mathcal{P} be a polyhedral cone in \mathbb{R}^p . Then there exist a collection of faces of \mathcal{P} , say $\{F_1, \dots, F_K\}$, such that the collection of their relative interiors, $\{ri(F_1), \dots, ri(F_K)\}$, forms a partition of \mathcal{P} . Further,

$$\|\Pi(\mathbf{y}|\mathcal{P})\|^2 = \sum_{i=0}^K I\{\Pi(\mathbf{y}|\mathcal{P}) \in ri(F_i)\} \|P_i \mathbf{y}\|^2$$

where $I(\cdot)$ is the indicator function and P_i is the projection matrix onto the linear space spanned by F_i .

Proof: Since \mathcal{P} is a closed convex set, $\Pi(Y|\mathcal{P})$ exists and is unique (see Proposition 3.12.3 on page 114), or Stoer and Witzgall (1970), p 51). If F and G are two different faces of \mathcal{P} , then $ri(F) \cap ri(G)$ is the empty set (see Stoer and Witzgall (1970), p 43, 2.4.14). The polyhedral cone \mathcal{P} , a typical face F , and its relative interior $ri(F)$ can be expressed as follows :

$$\begin{aligned} \mathcal{P} &= \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} \leq 0 \text{ for } i = 1, \dots, m\}, \\ F &= \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} = 0 \text{ for } i \in I, \mathbf{a}_j^T \mathbf{x} \leq 0 \text{ for } j \in J\}, \\ ri(F) &= \{\mathbf{x} : \mathbf{a}_i^T \mathbf{x} = 0 \text{ for } i \in I, \mathbf{a}_j^T \mathbf{x} < 0 \text{ for } j \in J\} \end{aligned} \quad (3.71)$$

where $\{I, J\}$ is a partition of $\{1, \dots, m\}$; it is assumed that the equality constraints defining F are contained in $\{\mathbf{a}_i^T \mathbf{x} = 0, i \in I\}$ and hence the set of constraints $\{\mathbf{a}_j^T \mathbf{x} \leq 0, j \in J\}$ does not include a constraint of the form $\mathbf{a}_j^T \mathbf{x} = 0$ for some $\mathbf{a} \neq 0$. Thus, it is clear that the collection of all relative interiors of faces of \mathcal{P} covers \mathcal{P} .

Let $\{F_1, \dots, F_K\}$ denote a collection of faces of \mathcal{P} such that $\{ri(F_1), \dots, ri(F_K)\}$ forms a partition of \mathcal{P} . Let P_i denote the projection matrix onto the linear space spanned by F_i , $i = 1, \dots, K$. Now, if $\Pi(\mathbf{y}|\mathcal{P}) \in ri(F_i)$ then $\Pi(\mathbf{y}|\mathcal{P}) = P_i \mathbf{y}$, by Lemma 3.13.2. This completes the proof. \blacksquare

The next lemma is really a restatement of the results in the previous lemmas, but it is stated here because this alternative statement is helpful.

Lemma 3.13.6 Let \mathcal{P} be a convex polyhedral in \mathbb{R}^p , and $\mathbf{X}_{p \times 1} \sim N(\mathbf{0}, \mathbf{I})$. Then there exists a collection of faces of \mathcal{P} , say $\{F_1, \dots, F_k\}$, such that their relative interiors, $\{ri(F_1), \dots, ri(F_k)\}$, are a partition of \mathcal{P} . Let $S_i = \{\mathbf{x} : \Pi(\mathbf{x}|\mathcal{P}) \in ri(F_i)\}$, P_i denote the projection matrix onto the linear space spanned by F_i , and $\tilde{\chi}^2(I, \mathcal{P}) = \|\mathbf{X}\|^2 - \min\{\|\mathbf{X} - \mathbf{a}\|^2 : \mathbf{a} \in \mathcal{P}\}$. Then, we have the following :

$$(i) \tilde{\chi}^2(I, \mathcal{P}) = \sum_{i=1}^k I(\mathbf{X} \in S_i) \|P_i \mathbf{X}\|^2.$$

- (ii) If $\mathbf{X} \in S_i$, then $P_i \mathbf{X}$ and $(I - P_i) \mathbf{X}$ are the projections of \mathbf{X} onto \mathcal{P} and \mathcal{P}° , respectively; hence $\tilde{\chi}^2(I, \mathcal{P}) = \|P_i \mathbf{X}\|^2$ and $\tilde{\chi}^2(I, \mathcal{P}^\circ) = \|(I - P_i) \mathbf{X}\|^2$.
- (iii) If $\mathbf{X} \in S_i$, then $\|P_i \mathbf{X}\|^2$ and $\|(I - P_i) \mathbf{X}\|^2$ are independent and are distributed as χ_ν^2 and $\chi_{p-\nu}^2$, respectively, where $\nu = \text{rank}(P_i)$. \blacksquare

Lemma 3.13.7 Let $\mathbf{Y} \sim N_p(\mathbf{0}, \mathbf{I})$, \mathcal{P} be a polyhedral cone, and let $\tilde{\chi}^2$ denote $\|\Pi(\mathbf{Y}|\mathcal{P})\|^2$. Then

$$\text{pr}(\tilde{\chi}^2 \geq c) = \sum_{i=0}^p w_i \text{pr}(\chi_i^2 \geq c), \text{ where } w_i = \sum_{\text{rank}(P_j)=i} \text{pr}\{\Pi(\mathbf{Y}|\mathcal{P}) \in ri(F_j)\}.$$

Thus, $w_i = \text{pr}\{\Pi(\mathbf{Y}|\mathcal{P}) \text{ lies in the relative interior of a face of } \mathcal{P} \text{ of dimension } i\}$, for $i = 0, \dots, p$. Further, $w_0 + \dots + w_p = 1$.

Proof: In view of the previous lemma, we have the following:

$$\begin{aligned} \text{pr}(\tilde{\chi}^2 \geq c) &= \sum_{j=0}^k \text{pr}\{\|\Pi(\mathbf{Y}|\mathcal{P})\|^2 \geq c | \Pi(\mathbf{Y}|\mathcal{P}) \in ri(F_j)\} \times \text{pr}\{\Pi(\mathbf{Y}|\mathcal{P}) \in ri(F_j)\} \\ &= \sum_{i=0}^p w_i \text{pr}(\chi_i^2 \geq c), \text{ where } w_i = \sum_{\text{rank}(P_j)=i} \text{pr}\{\Pi(\mathbf{Y}|\mathcal{P}) \in ri(F_j)\}. \end{aligned} \quad (3.72)$$

Thus, $w_i = \text{pr}\{\Pi(\mathbf{Y}|\mathcal{P}) \text{ lies on the relative interior of a face of dimension } i\}$, the dimension of a face being defined as the dimension of the linear space spanned by the face.

Proof of Theorem 3.4.2 continued: We will prove the theorem for a polyhedral cone first, and then deduce the result for a general cone. Assume that \mathcal{C} is a polyhedral cone.

Let $\mathbf{V}^{-1} = \mathbf{A}^T \mathbf{A}$ be a factorization of \mathbf{V}^{-1} , $\mathbf{X} = \mathbf{A}\mathbf{Y}$, and \mathcal{P} be the polyhedral cone $\mathcal{AC} = \{\mathbf{A}\mathbf{x} : \mathbf{x} \in \mathcal{C}\}$. Then the faces of \mathcal{C} take the form \mathbf{AF} for some face F of \mathcal{P} and $\dim(\mathbf{AF}) = \dim(F) = \dim(\mathcal{P})$ = dimension of the linear space spanned by F . Now,

$$\begin{aligned} LRT &= \mathbf{Y}^T \mathbf{V}^{-1} \mathbf{Y} - \min\{(\mathbf{Y} - \boldsymbol{\alpha})^T \mathbf{V}^{-1} (\mathbf{Y} - \boldsymbol{\alpha}) : \boldsymbol{\alpha} \in \mathcal{P}\} \\ &= \mathbf{X}^T \mathbf{X} - \min\{(\mathbf{X} - \boldsymbol{\beta})^T (\mathbf{X} - \boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathcal{P}\}. \end{aligned} \quad (3.73)$$

It follows from Lemma 3.13.7 that $\text{pr}(LRT \leq c | H_0) = \sum_{i=0}^p w_i(p, \mathbf{V}, \mathcal{P}) \text{pr}(\chi_i^2 \leq c)$, where $w_i(p, \mathbf{V}, \mathcal{P}) = \text{pr}\{\Pi(\mathbf{Y}|\mathcal{P}) \text{ lies on the relative interior of a face } G \text{ of } \mathcal{P} \text{ with } \dim(G) = i\} = \text{pr}\{\Pi(\mathbf{Y}|G) \text{ lies on the relative interior of a face } F \text{ of } \mathcal{C} \text{ with } \dim(F) = i\}$.

Now let us relax the assumption that \mathcal{C} is a polyhedral and assume only that it is a closed convex cone. Let $C_1 \subset C_2 \subset \dots$ be a sequence of polyhedral cones $\{ri(F_i)\}, F_i$ denote the projection matrix onto the linear space spanned by F_i , and $\tilde{\chi}^2(I, \mathcal{P}) = \|\mathbf{X}\|^2 - \min\{\|\mathbf{X} - \mathbf{a}\|^2 : \mathbf{a} \in \mathcal{P}\}$. Then, for any fixed (i, p, \mathbf{V}) , the sequence

$\{w_i(p, V, C_n)\}$ is bounded, and hence has limit points. Therefore, $\text{pr}(\|\Pi_V(Y, C_n)\|_V^2 \leq c)$ has limit points of the form $\sum_{i=0}^p w_i \text{pr}(X_i^2 \leq c)$ for some $\{w_i\}$. Since all such limit points are equal to $\text{pr}(\|\Pi_V(Y, C)\|_V^2 \leq c)$, for any c , it follows that limit points of $\{w_i(p, V, C_n)\}$ are unique. Therefore, we define $w_i(p, V, C) = \lim_{n \rightarrow \infty} w_i(p, V, C_n)$. ■

Proof of Proposition 3.6.1: The proofs of virtually all of these are based on the details of the proof of the previous theorem. The proof of part 1 is contained in the proof of Theorem 3.4.2 just given. Part 2 is a corollary to part 1. It was also proved by Nuesch (1964); see also Perlman (1969) and Wolak (1987). To our knowledge, part 3 has been a known result but it is not clear whether or not a proof has been published. This was posed as a conjecture (see Shapiro (1987)); Professor Alexander Shapiro informed us that he has received more than one proof of the conjecture, one by Professor J. Kinoses and another by Professor J. Lawrence; Professor Tim Robertson also informed us that he is aware of other proofs. Part 4 follows from part 3. Part 5 is shown in Shapiro (1985); see also Shapiro (1988, equation 3.4). Part 6 can be established easily by applying the transformation $Y = RZ$ where $Z \sim N(0, V)$; see also Shapiro (1988, equation 5.3). See equations (5.4), (5.5), (5.8), and (5.10) in Shapiro (1988) for parts 7, 8, 9, and 10, respectively. ■

Proof of Theorem 3.9.1: (a) Let $Z = E/\sigma$ and assume that the null hypothesis holds. Then, we have that

$$\bar{E}_A^2 = \left\{ \min_{\mathbf{a} \in \mathcal{M}} Q^*(\mathbf{a}) - \min_{\mathbf{a} \in \mathcal{C}} Q^*(\mathbf{a}) \right\} / \min_{\mathbf{a} \in \mathcal{M}} Q^*(\mathbf{a}) \quad (3.74)$$

where $Q^*(\mathbf{a}) = (\mathbf{A}Z - \mathbf{a})^T W^{-1}(\mathbf{A}Z - \mathbf{a})$ and $\mathbf{A} = W^{-1}\mathbf{X}^T U^{-1}$. The distribution of Z and hence that of (3.74) does not depend on any unknown parameters. Consequently, the null distribution of \bar{E}_A^2 is the same for any value of θ in \mathcal{M} and any value of $\sigma > 0$; by similar arguments, the corresponding result holds for \bar{F}_A as well.

(b) Let $\theta \in \mathcal{C}$. Let

$$\begin{aligned} g_1(\theta, E) &= \min\{\|\mathbf{X}\theta + E - \mathbf{X}\mathbf{a}\|_U^2 : \mathbf{a} \in \mathcal{C}\} \\ \text{and } g_2(\theta, E) &= \min\{\|\mathbf{X}\theta + E - \mathbf{X}\mathbf{a}\|_U^2 : \mathbf{a} \in \mathbb{R}^p\}. \end{aligned} \quad (3.75)$$

To show the dependence of \bar{E}_B^2 on θ , let us write

$$\bar{E}_B^2(\theta, E) = \{g_1(\theta, E) - g_2(\theta, E)\}/g_1(\theta, E).$$

Now,

$$\begin{aligned} \min\{\|\mathbf{Y} - \mathbf{X}\mathbf{a}\|_U^2 : \mathbf{a} \in \mathcal{C}\} &= \min\{\|\mathbf{X}\theta + E - \mathbf{X}\mathbf{a}\|_U^2 : \mathbf{a} \in \mathcal{C}\} \\ &= \min\{\|\mathbf{E} - \mathbf{X}(\mathbf{a} - \theta)\|_U^2 : \mathbf{a} \in \mathcal{C}\} = \min\{\|\mathbf{E} - \mathbf{X}\mathbf{b}\|_U^2 : \mathbf{b} \in \mathcal{C} - \theta\} \\ &\leq \min\{\|\mathbf{E} - \mathbf{X}\mathbf{b}\|_U^2 : \mathbf{b} \in \mathcal{C}\} \end{aligned}$$

the last step follows since $\mathcal{C} \subset \mathcal{C} - \theta$. Similarly,

$$\min\{\|\mathbf{Y} - \mathbf{X}\mathbf{a}\|_U^2 : \mathbf{a} \in \mathbb{R}^p\} = \min\{\|\mathbf{E} - \mathbf{X}\mathbf{b}\|_U^2 : \mathbf{b} \in \mathbb{R}^p\}.$$

Thus, $g_1(\theta, E) \leq g_1(0, E)$ and $g_2(\theta, E) = g_2(0, E)$. Now,

$$\begin{aligned} \bar{E}_B^2(\theta, E) &= 1 - \{g_2(\theta, E)/g_1(\theta, E)\} = 1 - \{g_2(0, E)/g_1(0, E)\} \\ &\leq 1 - \{g_2(0, E)/g_1(0, E)\} = \bar{E}_B^2(0, E). \end{aligned}$$

Therefore, $\bar{E}_B^2(\theta, E) \leq \bar{E}_B^2(0, E)$ and

$$\text{pr}\{\bar{E}_B^2(\theta, E) \geq c\} \leq \text{pr}\{\bar{E}_B^2(0, E) \geq c\}, \quad \theta \in \mathcal{C}.$$

Hence the least favorable null value for \bar{E}_B^2 is 0.

To prove the corresponding result for \bar{F}_B , note that $S^2(\theta, E) = S^2(0, E)$. Now,

$$\begin{aligned} \bar{F}(\theta, E) &= \{g_1(\theta, E) - g_2(\theta, E)\}/S^2(\theta, E) \\ &= \{g_1(\theta, E) - g_2(0, E)\}/S^2(0, E) \\ &\leq \{g_1(0, E) - g_2(0, E)\}/S^2(0, E) = \bar{F}(0, E). \end{aligned}$$

Therefore, $\text{pr}\{\bar{F}_B(\theta, E) \geq c\} \leq \text{pr}\{\bar{F}_B(0, E) \geq c\}$ for $\theta \in \mathcal{C}$. ■

Proof of Theorem 3.9.2 : (a) Without loss of generality assume that $\sigma = 1$, U is the identity matrix, and the null hypothesis is $R\theta = \mathbf{0}$ where R is a $r \times p$ matrix of rank r . Now

$$\bar{E}_B^2 = \{q(\theta^0) - q(\theta^1)\}/Q(\theta^0).$$

Note that

$$Q(\theta^0) = [q(\theta^0) - q(\theta^1)] + \|\mathbf{Y} - \mathbf{X}\hat{\theta}\|^2 + q(\theta^1).$$

Since $\hat{\theta}$ is independent of $(\mathbf{Y} - \mathbf{X}\hat{\theta})$, we have that $\mathbf{Y} - \mathbf{X}\hat{\theta}$ is also independent of $\{q(\theta^0) - q(\theta^1), q(\theta^1)\}$. Let $\{F_1, \dots, F_k\}$ be a collection of faces of $C \cap \mathcal{M}^\perp$ as in Lemma 3.13.4.

Now, suppose that $\Pi_W(\hat{\theta}|C \cap \mathcal{M}^\perp) \in ri(F)$ for some $F \in \{F_1, \dots, F_k\}$. Let $\dim\{ri(F)\} = j$. Then $q(\theta^0) - q(\theta^1) = \|\mathbf{P}\hat{\theta}\|^2$ where \mathbf{P} is the projection matrix onto the linear space spanned by F . Now, since $q(\theta^0) \sim \chi_r^2, q(\theta^0) - q(\theta^1) = \|\mathbf{P}\hat{\theta}\|^2 \sim \chi_j^2$, and $q(\theta^0) - q(\theta^1)$ and $q(\theta^1)$ are independent by Proposition 3.12.6 part (g), we have that $q(\theta^1) \sim \chi_{r-j}^2$. Therefore,

$$\begin{aligned} \bar{E}_A^2 &= \{q(\theta^0) - q(\theta^1)\}/[\{q(\theta^0) - q(\theta^1)\} + \|\mathbf{Y} - \mathbf{X}\hat{\theta}\|^2 + q(\theta^1)] \\ &\sim \chi_j^2/(\chi_j^2 + \chi_{N-p+r-j}^2), \end{aligned} \quad (3.76)$$

where the two terms in the denominator are independent. Therefore,

Let B_i denote the event $\{\hat{\theta}|C \cap M^\perp\} \in ri(F_i)\}$. Now, by adopting an argument similar to the derivation of $\bar{\chi}^2$ we have the following:

$$\begin{aligned}\text{pr}(\bar{E}_A^2 \leq c|H_0) &= \sum_i \text{pr}\{\bar{E}_A^2 \leq c \mid B_i\} \text{pr}\{B_i\}. \\ &= \sum_{j=0}^p \sum_{\dim(F_i)=j} \text{pr}\{\bar{E}_A^2 \leq c \mid B_i\} \text{pr}(B_i) \\ &= \sum_{j=0}^p w_j(p, \mathbf{W}, \mathcal{C} \cap M^\perp) \text{pr}[\beta(2^{-1}j, 2^{-1}(N-p+r-j)) \leq c].\end{aligned}$$

The distribution of \bar{F} is obtained similarly.

$$\begin{aligned}\text{pr}(\bar{F}_A \leq c|H_0) &= \sum_{i=1}^K \text{pr}\{\|\Pi(\hat{\theta}|C \cap M^\perp)\|^2 S^{-2} \leq c \mid B_i\} \text{pr}(B_i) \\ &= \sum_{j=0}^p \sum_{\dim(F_i)=j} \text{pr}\{\|\mathbf{P}_i \hat{\theta}\|^2 S^{-2} \leq c \mid B_i\} \text{pr}(B_i) \\ &= \sum_{j=0}^p w_j(p, \mathbf{W}, \mathcal{C} \cap M^\perp) \text{pr}(j F_{j,\nu} \leq c).\end{aligned}$$

Problems

Section 3.3

3.1 The length and direction of $N(\mathbf{0}, \mathbf{I})$ are independent; two-dimensional case:
Let $\mathbf{X} = (X_1, X_2)^T \sim N(\mathbf{0}, \mathbf{I})$, where \mathbf{I} is the 2×2 identity matrix. Let (r, α) be the polar coordinates of \mathbf{X} given by $(X_1, X_2) = r(\sin \alpha, \cos \alpha)$. Obtain the joint density of (r, α) , and deduce that r and α are independent and that $\|\mathbf{X}\|$ and $\mathbf{X}/\|\mathbf{X}\|$ are independent [see Anderson (1984), p 279].

3.2 Derivation of a $\bar{\chi}^2$ distribution in three dimensions: Let $\mathbf{X} = (X_1, X_2, X_3)^T \sim N(\boldsymbol{\theta}, \mathbf{I})$ and let the null and alternative hypotheses be $H_0 : \boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \boldsymbol{\theta} \geqslant \mathbf{0}$. Let $\tilde{\boldsymbol{\theta}}$ be the mle of $\boldsymbol{\theta}$ under H_1 based on a single observation of \mathbf{X} . Show that of H_0 against H_1 is also to be based on a single observation of \mathbf{X} . Show that $\tilde{\boldsymbol{\theta}} = (X_1 I\{X_1 > 0\}, X_2 I\{X_2 > 0\}, X_3 I\{X_3 > 0\})$, where I is the indicator function and $LRT = X_1^2 I(X_1 > 0) + X_2^2 I(X_2 > 0) + X_3^2 I(X_3 > 0)$. Show that \mathbf{X} has exactly 3, 2, 1, or 0 positive components with probabilities (1/8), (3/8), (3/8), or (1/8), respectively. Deduce that, for $c > 0$, the null distribution of LRT is given by, $\text{pr}(LRT \leq c) = \sum_{i=0}^3 w_i \text{pr}(X_i^2 \leq c)$, where $(w_0, w_1, w_2, w_3) = (1/8, 3/8, 3/8, 1/8)$. Verify that $w_i = \text{pr}(\text{exactly } i \text{ components of } \mathbf{Z} \text{ are positive})$, where $\mathbf{Z}_{3 \times 1} \sim N(\mathbf{0}, \mathbf{I})$, and $i = 0, \dots, 3$.

3.3 Derivation of a $\bar{\chi}^2$ distribution in three dimensions: Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{I})$, where \mathbf{X} is 3×1 , and let $\mathcal{C} = \{\boldsymbol{\theta} \in \mathbb{R}^3 : \theta_2 \geqslant 0, \theta_3 \geqslant 0\}$. Show that the mle of $\boldsymbol{\theta}$ subject to $\boldsymbol{\theta} \in \mathcal{C}$ and the LRT of $H_0 : \boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \boldsymbol{\theta} \in \mathcal{C}$ are given by

$$(\tilde{\theta}_1, \tilde{\theta}_2, \tilde{\theta}_3) = \begin{cases} (X_1, X_2, X_3) & \text{if } X_2 \geqslant 0 \text{ and } X_3 \geqslant 0 \\ (X_1, X_2, 0) & \text{if } X_3 < 0 \text{ and } X_2 > 0 \\ (X_1, 0, X_3) & \text{if } X_2 < 0 \text{ and } X_3 > 0 \\ (X_1, 0, 0) & \text{if } X_2 < 0 \text{ and } X_3 < 0 \end{cases}$$

and

$$LRT = \begin{cases} X_1^2 + X_2^2 + X_3^2 & \sim \chi^2_3 \text{ if } X_2 \geqslant 0 \text{ and } X_3 \geqslant 0 \\ X_1^2 + X_2^2 & \sim \chi^2_2 \text{ if } X_3 < 0 \text{ and } X_2 > 0 \\ X_1^2 + X_3^2 & \sim \chi^2_2 \text{ if } X_2 < 0 \text{ and } X_3 > 0 \\ X_1^2 & \sim \chi^2_1 \text{ if } X_2 < 0 \text{ and } X_3 < 0, \end{cases}$$

respectively. Deduce that the null distribution of LRT is given by

$$\text{pr}(LRT \leq c|H_0) = (1/4)\text{pr}(\chi_1^2 \leq c) + (1/2)\text{pr}(\chi_2^2 \leq c) + (1/4)\text{pr}(\chi_3^2 \leq c).$$

3.4 Derivation of a $\bar{\chi}^2$ distribution in four dimensions: Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{I})$ where \mathbf{X} is 4×1 . Let the null and alternative hypothesis be $H_0 : \boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \theta_3 \geqslant 0, \theta_4 \geqslant 0$. Show that

$$\begin{aligned}\text{pr}(\bar{F}_B \leq c|\boldsymbol{\theta} = \mathbf{0}) &= \sum_{i=0}^K \text{pr}[q(\boldsymbol{\theta}^1)/\{q(\boldsymbol{\theta}^1) + \|Y - \mathbf{X}\hat{\theta}\|_U^2\} \leq c \mid A_i] \text{pr}(A_i) \\ &= \sum_{j=0}^p \sum_{\dim(F_i)=j} \text{pr}[\chi_{p-j}^2 / \{\chi_{p-j}^2 + \chi_{N-p}^2\} \leq c] \text{pr}(A_i) \\ &= \sum_{j=0}^p w_{p-j}(p, \mathbf{W}, C) \text{pr}[\beta(2^{-1}j, 2^{-1}(N-p)) \leq c].\end{aligned}$$

The proof for the \bar{F} is similar.

$$\begin{aligned}\text{pr}(\bar{F}_B \leq c|\boldsymbol{\theta} = \mathbf{0}) &= \sum_{i=0}^K \text{pr}[q(\boldsymbol{\theta}^1)/S^2 \leq c] \text{pr}(A_i) = \sum_{i=0}^p w_{p-j}(p, \mathbf{W}, C) \text{pr}[F_{j,\nu} \leq c].\end{aligned}$$

$$= \sum_{i=0}^k \sum_{j=0}^p \text{pr}[\chi_{p-j}^2 / S^2 \leq c] \text{pr}(A_i) = \sum_{j=0}^p w_{p-j}(p, \mathbf{W}, C) \text{pr}[F_{j,\nu} \leq c].$$

where θ_2 is a subvector of θ consisting of k elements. Show that $\text{pr}(LRT \leq c | H_0) = w_{p-k} \text{pr}(\chi^2_{p-k} \leq c) + \dots + w_p \text{pr}(\chi^2_p \leq c)$ where $w_{p-k+i} = 2^{-k} k! / \{(k-i)! i!\}$.

3.6 Application of LRT in two dimensions: Let $X = (X_1, X_2)^T \sim N(\theta, I)$. The mean of a sample of five observations on X is $(-1, 2)$. Compute the *mle* under the constraint $\theta_1 \geq 0, \theta_2 \geq 0$ and the *p*-value for testing $H_0 : \theta_1 = \theta_2 = 0$ against $H_1 : \theta_1 \geq 0, \theta_2 \geq 0$.

3.7 Derivation of a $\bar{\chi}^2$ distribution in two dimensions: Let $X \sim N(\theta, I)$, $H_0 : \theta_1 = \theta_2 = 0$, $H_1 : \theta_1 \geq 0, \theta_2 - m\theta_1 \geq 0$, where m is known. Show from first principles that the distribution of the LRT under H_0 is given by $\text{pr}(LRT \leq c | H_0) = \sum_{i=2}^{i=2} w_i \text{pr}(\chi^2_i \leq c)$ where $w_0 = q, w_1 = 0.5, w_2 = 0.5 - q$, and $q = (2\pi)^{-1} \cos^{-1} \{-m/(1+m^2)^{1/2}\}$. Indicate the critical region corresponding to the critical value 4. [The details of the derivations are also given in (Gourioux et al., 1982, section 4.2).]

3.8 Derivation of the critical region of LRT in two dimensions: For Example 3.3.5, prove the following: (i) $(1, \rho) \perp_V (0, 1)$ and $(1, \rho^{-1}) \perp_V (1, 0)$. (ii) The boundaries of the polar cone of the nonnegative orthant are $\theta_2 = \rho^{-1}\theta_1$ and $\theta_2 = \rho\theta_1$. (iii) For $\rho = 0.9$ and $\bar{x} = (-3, 2)$, the *mle* of θ subject to $\theta \geq 0$ is the point of intersection of the θ_2 axis and the line of slope ρ through \bar{x} . (iv) For $\rho = 0.9$ and $\bar{x} = (3, -2)$, the *mle* of θ subject to $\theta \geq 0$ is the point of intersection of the θ_1 axis and the line of slope ρ^{-1} through \bar{x} . Indicate the critical region $\{LRT \geq 6\}$.

3.9 Let X_i be an observation from the i th treatment and assume that $X_i \sim N(\mu_i, 1)$ for $i = 1, 2, 3$. Let $H_0 : \mu_1 = \mu_2 = \mu_3$ and $H_1 : \mu_1 \geq \mu_2$ and $\mu_1 \geq \mu_3$. (i) Verify that the LRT of H_0 against H_1 rejects H_0 at 5% level when $(X_1, X_2, X_3) = (50, 150, 70)$. (ii) Verify that the phenomenon in part (i) is essentially the same as that in Example 3.3.6.

[Hint: Let $Y_1 = X_1 - X_2, Y_2 = X_1 - X_3, Y = (Y_1, Y_2)^T$ and $X = (X_1, X_2)^T$. Then $Y = AX, Y \sim N(\theta, V)$, and $\theta = A\mu$ where $V = AA^T, A = (1, -1, 0; 1, 0, -1)$, and $V = 2(1, 0.5; 0.5, 1)$. Verify that the LRT of H_0 against H_1 is equivalent to the LRT of $H_0^* : \theta = 0$ against $H_1^* : \theta \geq 0$ based on a single observation of Y . Verify (by using the result in Problem 8) that $Y = (-100, -20)^T$, which corresponds to $X = (50, 150, 70)^T$, is in the critical region.]

Section 3.4

3.10 Pythagoras theorem in \mathbb{R}^p : Let V be a $p \times p$ positive definite matrix. Define the inner product $\langle \mathbf{x}, \mathbf{y} \rangle_V = \mathbf{x}^T V^{-1} \mathbf{y}$, and the distance function $\| \mathbf{x} - \mathbf{y} \|_V = \{ \mathbf{x}^T V^{-1} \mathbf{y} \}^{1/2}$. We say that \mathbf{x} is orthogonal to \mathbf{y} with respect to V , denoted $\mathbf{x} \perp_V \mathbf{y}$, if $\mathbf{x} \neq \mathbf{0}, \mathbf{y} \neq \mathbf{0}$, and $\langle \mathbf{x}, \mathbf{y} \rangle_V = 0$. If A, B, C are three points in \mathbb{R}^p such that $AB \perp_V AC$, show that $\| AB \|_V^2 + \| AC \|_V^2 = \| BC \|_V^2$.

3.11 Let \mathcal{L} be a subspace of \mathbb{R}^p , $r = \dim(\mathcal{L})$, and $\mathbf{x} \in \mathbb{R}^p$. Let O be the origin and A be the point such $\mathbf{x} = OA$. Let V be a $p \times p$ positive definite matrix. Show that the point in \mathcal{L} that is V -closest to A is obtained by dropping a V -perpendicular to \mathcal{L} from A . [Hint: Use the Pythagoras theorem.]

3.12 Let $X \sim N(\theta, V)$, $H_0 : \theta = 0$, and $H_1 : \theta \in C$. Suppose that $V^{-1} = A^T A = B^T B$ for some A and B . Show that BA^{-1} is an orthogonal matrix. Now, transform the setting from $\{X, \theta, V, C\}$ by A as in Example 3.3.4 on page 68; also, consider the transformation by B . Show that the geometry of the image under transformation A is the same as that under transformation B except for a rotation defined by the orthogonal matrix BA^{-1} .

3.13 Let Y_1, \dots, Y_n be independent standard normal rv's, and let $\mathbf{Y} = (Y_1, \dots, Y_n)^T$. Let B be a $k \times n$ matrix of given constants such that $\{y : By \geq 0\}$ is nonempty. Show that the conditional distribution of $\| \mathbf{Y} \|^2$ given $BY \geq 0$ is χ^2_n .

3.14 A function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is said to be Lipschitz continuous if $|f(\mathbf{x}) - f(\mathbf{y})| \leq k \|\mathbf{x} - \mathbf{y}\|$ for some $k > 0$. Let $\mathcal{A} \subset \mathbb{R}^p$ and $d_{\mathcal{A}}$ be the distance function $d_{\mathcal{A}} : \mathbb{R}^p \rightarrow \mathbb{R}$ defined by $d_{\mathcal{A}}(\mathbf{x}) = \|\mathbf{x} - \mathcal{A}\|$. Show that $d_{\mathcal{A}}$ is Lipschitz continuous with $k = 1$.

3.15 Let A be an $m \times n$ matrix and $\text{rank}(A) = m$. Let $P = I - A^T(AA^T)^{-1}A$. Then, show that P is the projection matrix onto the null space of A .

3.16 Let V be a positive definite matrix, and let $V^{-1} = P^T D^{-1} P$ for some D and P , where P is orthogonal. Consider the problem of testing $H_0 : \theta = 0$ against $H_1 : \theta \in C$ based on a single observation of X where $X \sim N(\theta, V)$. Show that the LRT of H_0 against H_1 based on a single observation of X is equivalent to testing $H_0^* : \theta = 0$ against $H_1^* : \theta \in \mathcal{P}$ based on a single observation of Y , where $Y = PX$ and $\mathcal{P} = PC = \{Px : x \in C\}$ which is also a closed convex cone. Deduce that $\bar{\chi}^2(V, C) \sim \bar{\chi}^2(D, PC)$.

Similarly, let $V^{-1} = A^T A$; thus A can be the Cholesky factor or $D^{-1/2}P$. Let $Z = AX$ and $Q = AC = \{Ax : x \in C\}$. Show that Q is also a closed convex cone, the testing problem is equivalent to test of $H_0' : \gamma = 0$ vs $H_1' : \gamma \in Q$ based on a single observation of $Z \sim N(\gamma, I)$ and $\bar{\chi}^2(V, C) \sim \bar{\chi}^2(I, Q)$.

3.17 The $\bar{\chi}^2(V, C)$ distribution and the weights depend on V only through the correlation matrix corresponding to V : Let $V = (v_{ij})$ be a positive definite matrix of order $p \times p$. Consider the problem of testing $H_0 : \theta = 0$ against $H_1 : \theta \in C$ based on a single observation of $X \sim N(\theta, V)$. Let Ω denote the correlation matrix corresponding to V . Let A be the diagonal matrix $\text{diag}(a_1, \dots, a_p)$ where $a_i = v_{ii}^{1/2}$. Verify that $V = A\Omega A$. Show that $AC = C$. Deduce that $\bar{\chi}^2(V, C) \sim \bar{\chi}^2(\Omega, C)$ and hence the $\bar{\chi}^2$ -weights depend on V only through its correlation matrix.

3.18 Let R be a $k \times p$ matrix of rank $k, k \leq p$, and $\mathbf{y} = Ra$. Verify that $\min\{(\mathbf{x}-\mathbf{a})^T V^{-1}(\mathbf{x}-\mathbf{a}) : Ra \geq 0\} = \min\{(\mathbf{y}-\mathbf{b})^T (RV R^T)^{-1}(\mathbf{y}-\mathbf{b}) : b \geq 0\}$. Dede that

$$\bar{\chi}^2(V, C) \sim \bar{\chi}^2(RVR^T, \mathbb{R}^{+k})$$

where $C = \{ \theta : R\theta \geq 0 \}$.

[Hint: Find a matrix Q such that $RVQ^T = 0$ and $C^T = [R^T, Q^T]$ is nonsingular; then show that $\min\{(\mathbf{x}-\mathbf{a})^T V^{-1}(\mathbf{x}-\mathbf{a}) : Ra \geq 0\} = \min\{(\mathbf{z}-c)^T (CV C^T)^{-1}(\mathbf{z}-c) : c \geq 0\}$. Now use the fact that $CV C^T$ is block diagonal.]

Section 3.8

3.19 Derivation of the null distribution of LRT for a Type B problem in two dimensions: Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{I})$, $H_1 : \theta_1 \geq 0, \theta_2 - m\theta_1 \geq 0$ and $H_2 : \text{not } H_1$, where m is known. Verify the following: (a) $LRT = \min\{(\mathbf{X} - \mathbf{a})^T(\mathbf{X} - \mathbf{a}) : a_1 \geq 0, a_2 - ma_1 \geq 0\}$. (b) $\text{pr}(LRT \leq c|\boldsymbol{\theta} = \mathbf{0}) = (0.5 - q)\text{pr}(\chi_0^2 \leq c) + 0.5\text{pr}(\chi_1^2 \leq c) + q\text{pr}(\chi_2^2 \leq c)$, where $q = (2\pi)^{-1} \cos^{-1}(-m(1 + m^2)^{-1/2})$.

3.20 Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{I})$, where \mathbf{X} is 2×1 , $H_0 : \boldsymbol{\theta} \in \mathcal{C}$ and $H_2 : \boldsymbol{\theta} \notin \mathcal{C}$. Let \mathcal{C}° denote the polar cone of \mathcal{C} with respect to \mathbf{I} . Let S_1, \dots, S_4 denote the four cones as in Example 3 in Section 3.3. (a) Show that, conditional on $\mathbf{X} \in S_i$, $\bar{\boldsymbol{\theta}}$ and $\mathbf{X} - \bar{\boldsymbol{\theta}}$ are independent, $i = 1, \dots, 4$. (b) Let L_A and L_B denote the LRT for the Type A testing problem, $H_0 : \boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \boldsymbol{\theta} \geq \mathbf{0}$, and for the Type B testing problem, $H_1 : \boldsymbol{\theta} \in \mathcal{C}$ against $H_2 : \boldsymbol{\theta} \notin \mathcal{C}$, respectively. Show that the joint distribution of L_A and L_B at $\boldsymbol{\theta} = \mathbf{0}$ is given by

$$\begin{aligned} \text{pr}(L_A \geq c_1, L_B \geq c_2 | \boldsymbol{\theta} = \mathbf{0}) &= (0.5 - q)\text{pr}(\chi_0^2 \geq c_1)\text{pr}(\chi_2^2 \geq c_2) + \\ &0.5\text{pr}(\chi_1^2 \geq c_1)\text{pr}(\chi_1^2 \geq c_2) + q\text{pr}(\chi_2^2 \geq c_1)\text{pr}(\chi_0^2 \geq c_2). \end{aligned}$$

3.21

3.21.1. Let $\mathbf{X} \sim N(\boldsymbol{\theta}, V)$. Let $\bar{\boldsymbol{\theta}}$ denote the mle subject to $\boldsymbol{\theta} \in \mathcal{C}$, L_A and L_B denote the LRT's for $H_0 : \boldsymbol{\theta} = \mathbf{0}$ against $H_1 : \boldsymbol{\theta} \in \mathcal{C}$ and $H_1 : \boldsymbol{\theta} \in \mathcal{C}$ against $H_2 : \boldsymbol{\theta} \notin \mathcal{C}$, respectively. Then the joint distribution of L_A and L_B at $\boldsymbol{\theta} = \mathbf{0}$ is given by the following:

$$\text{pr}(L_A \leq c_1, L_B \leq c_2 | \boldsymbol{\theta} = \mathbf{0}) = \sum_{i=0}^p w_i(p, V, C) \text{pr}(\chi_i^2 \leq c_1) \text{pr}(\chi_{p-i}^2 \leq c_2).$$

3.21.2. Let the hypotheses H_0, H_1 , and H_2 be defined as $H_0 : \boldsymbol{\theta} \in \mathcal{M}, H_1 : \boldsymbol{\theta} \in \mathcal{C}$, and $H_2 : \boldsymbol{\theta} \in \mathcal{L}$ where \mathcal{M} and \mathcal{L} are linear spaces, \mathcal{C} is a closed convex cone, $\mathcal{M} \subset \mathcal{C} \subset \mathcal{L}$, and $\text{dim}(\mathcal{L}) = k \leq p$. Obtain an expression for $\text{pr}(L_A \leq c_1, L_B \leq c_2)$.

3.21.3. Suppose that $\mathbf{X} \sim N(\boldsymbol{\theta}, \sigma^2 U)$ where σ is unknown and U is known. Using similar arguments obtain expressions for $\text{pr}_0(\bar{E}_A^2 \geq c, \bar{E}_B^2 \geq d)$ and $\text{pr}_0(\bar{F}_A \geq c, \bar{F}_B \geq d)$.

3.21.4. Consider the linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}$ where $\mathbf{E} \sim N(0, \sigma^2 U)$, $\boldsymbol{\sigma}$ is unknown and \mathbf{U} is known. Using similar arguments obtain expressions for $\text{pr}_0(\bar{E}_A^2 \geq c, \bar{E}_B^2 \geq d)$ and $\text{pr}_0(\bar{F}_A \geq c, \bar{F}_B \geq d)$. [Hint:

$$\begin{aligned} \text{pr}_0(L_A \leq c_1, L_B \leq c_2) &= \sum \text{pr}_0(L_A \leq c_1, L_B \leq c_2 | X \in S_i) \text{pr}(X \in S_i) \\ &= \sum \text{pr}_0(\|P_i X\|^2 \leq c_1, \|P_i X\|^2 \leq c_2 | (I - P)_i X \|^2 \leq c_2 | X \in S_i) \text{pr}(X \in S_i). \end{aligned}$$

Now use Lemma 3.13.6.]

Section 3.9

3.22 For Example 1.2.2 in Chapter 1 on El Niño, recall that the hypothesis of interest is that

the warm phase of El Niño tends to suppress hurricanes whereas cold phase encourages hurricanes.

- (a) Test the foregoing hypothesis as a null hypothesis.
- (b) Test the same hypothesis as the alternative.

3.23 Connections among \bar{E}^2, \bar{F} and LRT: Let $\mathbf{X} = (X_1, X_2)^T$ and assume that $\mathbf{X} \sim N(\boldsymbol{\theta}, \sigma^2 \mathbf{U})$. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be n independently and identically distributed observations on \mathbf{X} . Show that this set up is a special case of the linear model in 3.40. Let the null and alternative hypotheses be $H_0 : \boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \boldsymbol{\theta} \geq \mathbf{0}$. Let $Q_1(\boldsymbol{\theta}) = \sum (\mathbf{X}_i - \boldsymbol{\theta})^T \mathbf{U}^{-1} (\mathbf{X}_i - \boldsymbol{\theta})$ and $q_1(\boldsymbol{\theta}) = n(\bar{\mathbf{X}} - \boldsymbol{\theta})^T \mathbf{U}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\theta})$. Show that $Q_1(\boldsymbol{\theta}) = \sum (\mathbf{X}_i - \bar{\mathbf{X}})^T \mathbf{U}^{-1} (\mathbf{X}_i - \bar{\mathbf{X}}) + q_1(\boldsymbol{\theta})$ and that the loglikelihood $L(\boldsymbol{\theta}, \sigma^2)$ takes the form $L(\boldsymbol{\theta}, \sigma^2) = -(1/2)\sigma^{-2} Q_1(\boldsymbol{\theta}) - \frac{1}{2} np \log(\sigma^2) + Const$. Let the null and alternative hypotheses be $H_a : \boldsymbol{\theta} \in \mathcal{C}_a$ and $H_b : \boldsymbol{\theta} \in \mathcal{C}_b$, respectively, where $\mathcal{C}_a \subset \mathcal{C}_b$. Let $\boldsymbol{\theta}^a$ and $\boldsymbol{\theta}^b$ be the points at which $Q_1(\boldsymbol{\theta})$ is minimized over \mathcal{C}_a and \mathcal{C}_b , respectively. Show that, the likelihood ratio statistic for testing H_a against H_b is given by $LRT = np \log\{\frac{Q_1(\boldsymbol{\theta}^a)}{Q_1(\boldsymbol{\theta}^b)} / Q_1(\boldsymbol{\theta}^b)\}$. Let $\bar{E}^2 = \{Q_1(\boldsymbol{\theta}^a) - Q_1(\boldsymbol{\theta}^b)\} / Q_1(\boldsymbol{\theta}^a)$ and $\bar{F} = \{Q_1(\boldsymbol{\theta}^a) - Q_1(\boldsymbol{\theta}^b)\} / S^2$, where $S^2 = \nu^{-1} \sum (\mathbf{X}_i - \bar{\mathbf{X}})^T \mathbf{U}^{-1} (\mathbf{X}_i - \bar{\mathbf{X}})$ and $\nu = 2(n-1)$. Show that $\bar{E}^2 = [1 - \exp\{-LRT/(np)\}]$. Deduce that the \bar{E}^2 -test that rejects H_a for large values of \bar{E}^2 , is equivalent to the likelihood ratio test.

(a) Now let the null and alternative hypotheses be $H_0 : \boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \boldsymbol{\theta} \geq \mathbf{0}$. Let the four quadrants Q_1, \dots, Q_4 of \mathbb{R}^2 be defined as in Examples 1 and 2 of Section 3.3. Show that $\bar{X}_1, \bar{X}_2, \sum \|\mathbf{X}_i - \bar{\mathbf{X}}\|^2$ are mutually independent. By arguments similar to those for (3.7) show that

$$\bar{E}_A^2 = \begin{cases} T / \{T + \sum \|\mathbf{X}_i - \bar{\mathbf{X}}\|^2\} & \sim \beta(\frac{n}{2}, \frac{n-2}{2}) \quad \text{given } \bar{\mathbf{X}} \in Q_1 \\ (n\bar{X}_1^2) / \{T + \sum \|\mathbf{X}_i - \bar{\mathbf{X}}\|^2\} & \sim \beta(\frac{1}{2}, \frac{n-1}{2}) \quad \text{given } \bar{\mathbf{X}} \in Q_4 \\ (n\bar{X}_2^2) / \{T + \sum \|\mathbf{X}_i - \bar{\mathbf{X}}\|^2\} & \sim \beta(\frac{1}{2}, \frac{n-1}{2}) \quad \text{given } \bar{\mathbf{X}} \in Q_2 \\ 0 / \{0 + T + \sum \|\mathbf{X}_i - \bar{\mathbf{X}}\|^2\} & \sim \beta(\frac{1}{2}, \frac{n-0}{2}) \quad \text{given } \bar{\mathbf{X}} \in Q_3 \end{cases}$$

where $T = n\bar{X}_1^2 + n\bar{X}_2^2$.

By applying $\text{pr}(\bar{E}_A^2 \leq c | H_0) = \sum_1^4 \text{pr}(\bar{E}_A^2 \leq c | \bar{\mathbf{X}} \in Q_i) \text{pr}(\bar{\mathbf{X}} \in Q_i)$ or otherwise, deduce that $\text{pr}(\bar{E}_A^2 \leq c | H_0)$ is equal to

$$\frac{1}{4} \text{pr}\{\beta(\frac{0}{2}, \frac{n-0}{2}) \leq c\} + \frac{1}{2} \text{pr}\{\beta(\frac{1}{2}, \frac{n-1}{2}) \leq c\} + \frac{1}{4} \text{pr}\{\beta(\frac{2}{2}, \frac{n-2}{2}) \leq c\}.$$

- (b) Show that

$$\bar{F}_A = \begin{cases} (n\bar{X}_1^2 + n\bar{X}_2^2) / S^2 & \sim 2F_{2,\nu} \quad \text{conditional on } \bar{\mathbf{X}} \in Q_1 \\ n\bar{X}_1^2 / S^2 & \sim 1F_{1,\nu} \quad \text{conditional on } \bar{\mathbf{X}} \in Q_4 \\ n\bar{X}_2^2 / S^2 & \sim 1F_{1,\nu} \quad \text{conditional on } \bar{\mathbf{X}} \in Q_2 \\ 0 / S^2 & \sim 0F_{0,\nu} \quad \text{conditional on } \bar{\mathbf{X}} \in Q_3. \end{cases}$$

Show that, under the null hypothesis,

$$\text{pr}(\bar{F}_A \leq c) = (1/4)\text{pr}(2F_{2,\nu} \leq c) + (1/2)\text{pr}(1F_{1,\nu} \leq c) + (1/4)\text{pr}(0F_{0,\nu} \leq c). \blacksquare$$

3.24 Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \sigma^2 U)$, where σ is unknown. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be independently and identically distributed as \mathbf{X} . Let the null and alternative hypotheses be $H_1 : \mathbf{R}_1 \boldsymbol{\theta} \geq \mathbf{0}$, $\mathbf{R}_2 \boldsymbol{\theta} = \mathbf{0}$ and $H_2 : \boldsymbol{\theta}$ is not restricted, where \mathbf{R}_1^T is $s \times m$ and \mathbf{R}_2^T is $t \times m$, and $[\mathbf{R}_1^T, \mathbf{R}_2^T]$ has full rank. Show that the least favorable null value of \bar{E}^2 and of \bar{F} is $\boldsymbol{\theta} = \mathbf{0}$, and their least favorable null distributions are

$$\text{pr}_0(\bar{E}^2 \leq c | \boldsymbol{\theta} = \mathbf{0}) = \sum_{i=0}^s w_{s-i}(s, A) \text{pr}\{\beta\{\frac{1}{2}(t+i), \frac{1}{2}(np-p)\}\} \leq c.$$

$$\text{pr}_0(\bar{F} \leq c | \boldsymbol{\theta} = \mathbf{0}) = \sum_{i=0}^s w_{s-i}(s, A) \text{pr}\{(t+i)F_{t+i, \nu} \leq c\}$$

where $A = \mathbf{R}_1 \mathbf{V} \mathbf{R}_1^T - (\mathbf{R}_1 \mathbf{V} \mathbf{R}_2^T)(\mathbf{R}_2 \mathbf{V} \mathbf{R}_2^T)^{-1}(\mathbf{R}_2 \mathbf{V} \mathbf{R}_1^T)$.

3.25 Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be independently and identically distributed as \mathbf{X} where $\mathbf{X} \sim N(\boldsymbol{\theta}, \sigma^2 U)$. Let the null and alternative hypotheses be $H_0 : \mathbf{R}\boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \mathbf{R}_1 \boldsymbol{\theta} \geq \mathbf{0}$, where \mathbf{R} is a full row-rank matrix of order $r \times p$ and \mathbf{R}_1 is a submatrix of \mathbf{R} of order $q \times p$. Show that \bar{E}^2 and \bar{F} tests are similar and that their null distributions are given by

$$\text{pr}(\bar{E}^2 \leq c | H_0) = \sum_{j=0}^q w_j(q, \mathbf{R}_1 \mathbf{V} \mathbf{R}_1^T) \text{pr}\{\beta(2^{-1}(r-q+j), 2^{-1}(np-p+q-j)) \leq c\},$$

$$\text{pr}(\bar{F} \leq c) = \sum_{j=0}^q w_j(q, \mathbf{R}_1 \mathbf{V} \mathbf{R}_1^T) \text{pr}\{(r-q+j)F_{r-q+j, \nu} \leq c\}.$$

Section 3.10

3.26 For the hypertensive patients in the sodium reduction example (see Example 1.2.17 in Chapter 1), obtain an overall p -value for testing the hypothesis that a reduction in sodium intake leads to a reduction in blood pressure. State the assumptions that you made and discuss their limitations.

3.27 Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{V})$ where \mathbf{V} is unknown. Let $\boldsymbol{\theta}_2$ and \mathbf{S} be as in Section 3.10. Consider the test of $H_0 : \boldsymbol{\theta}_2 = \mathbf{0}$ against $H_1 : \boldsymbol{\theta}_2 \in \mathcal{P}$, where \mathcal{P} is a one-sided closed cone. Show that the LRT takes the simpler form (see Perlman (1969), equation (6.13))

$$LRT = \{\|\mathbf{Y}\|_F^2 - \|\mathbf{Y} - \Pi_T(\mathbf{Y}|\mathcal{P})\|_F^2\} \{1 + \|\mathbf{Y} - \Pi_T(\mathbf{Y}|\mathcal{P})\|_F^2\}^{-1}$$

where $\mathbf{Y} = \bar{\mathbf{X}}_{22}$ and $T = \mathbf{S}_{22}$.

3.28 Let \mathbf{A} be the $k \times (k-1)$ matrix such that $\mathbf{A}\boldsymbol{\mu} \geq \mathbf{0}$ is equivalent to $\mu_1 \leq \dots \leq \mu_k$. Let \mathbf{y}^* be the point at which $(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{W}(\mathbf{y} - \boldsymbol{\mu})$ reaches its minimum over $\mathbf{A}\boldsymbol{\mu} \geq \mathbf{0}$, where \mathbf{W} is a diagonal matrix with positive diagonal elements. Let $\mathcal{Q} = \{\boldsymbol{\mu} : \mathbf{A}\boldsymbol{\mu} \geq \mathbf{0}\}$. Let \mathcal{M} be the linear space spanned by $(1, \dots, 1)$. Deduce the following [Hint: Use the results relating to Fig. 3.10, in particular those in Section 3.12.2. Proofs of these from first principles are provided in Barlow et al. (1972); see also Robertson et al. (1988)].

1. \mathcal{M} is the largest linear space contained in \mathcal{Q} .
2. $(\mathbf{y} - \mathbf{y}^*) \perp \mathbf{w}_{-1}, \mathcal{M}$.

3.28 Let $(\mathbf{y} - \mathbf{y}^*)^T \mathbf{W} \mathbf{1} = 0$. (Hint: $\mathbf{1} \in \mathcal{M}$).

4. $(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{W}(\mathbf{y} - \boldsymbol{\mu}) \geq (\mathbf{y} - \mathbf{y}^*)^T \mathbf{W}(\mathbf{y} - \mathbf{y}^*) + (\mathbf{y}^* - \boldsymbol{\mu})^T \mathbf{W}(\mathbf{y}^* - \boldsymbol{\mu})$ for every $\boldsymbol{\mu} \in \mathcal{Q}$.

$$5. (\mathbf{y} - \mathbf{y}^*)^T \mathbf{W} \mathbf{y}^* = 0.$$

$$6. (\mathbf{y} - \mathbf{y}^*)^T \mathbf{W} \boldsymbol{\mu} \leq 0 \text{ for every } \boldsymbol{\mu} \in \mathcal{Q}.$$

$$7. \text{If } \tilde{\mathbf{y}} \in \mathcal{Q} \text{ and } (\mathbf{y} - \tilde{\mathbf{y}})^T \mathbf{W}(\tilde{\mathbf{y}} - \boldsymbol{\mu}) \geq 0 \text{ for every } \boldsymbol{\mu} \in \mathcal{Q} \text{ then } \tilde{\mathbf{y}} = \mathbf{y}^*.$$

$$8. \text{Let } \tilde{\mathbf{y}} \in \mathcal{Q}. \text{ Then } \tilde{\mathbf{y}} = \mathbf{y}^* \text{ if and only if } (\mathbf{y} - \tilde{\mathbf{y}})^T \mathbf{W} \tilde{\mathbf{y}} = 0 \text{ and } (\mathbf{y} - \tilde{\mathbf{y}})^T \mathbf{W} \boldsymbol{\mu} \leq 0 \text{ for every } \boldsymbol{\mu} \in \mathcal{Q}.$$

3.29 Let $\hat{\boldsymbol{\theta}}$ be an estimator of $\boldsymbol{\theta}$ where $\boldsymbol{\theta} \in \mathbb{R}^p$. Suppose that $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$ converges in distribution to $N(\mathbf{0}, \mathbf{V})$ where \mathbf{V} is positive definite. Let $\hat{\mathbf{V}}$ be a consistent estimator of \mathbf{V} (i.e., $\hat{\mathbf{V}} \xrightarrow{p} \mathbf{V}$), $H_0 : \boldsymbol{\theta} \in \mathcal{M}$, $H_1 : \boldsymbol{\theta} \notin \mathcal{C}$. Let $L_A = n\{\|\hat{\boldsymbol{\theta}} - \mathcal{M}\|^2 - \|\hat{\boldsymbol{\theta}} - \mathcal{C}\|^2\}$ and $L_B = n\{\|\hat{\boldsymbol{\theta}} - \mathcal{C}\|^2\}$, where $\|\cdot\|$ is defined with respect to $\langle \cdot \rangle_{\hat{\mathbf{V}}}$. Following the notation introduced in this chapter, it is implicitly assumed that \mathcal{C} is a closed convex cone, \mathcal{M} is a linear space, and $\mathcal{M} \subset \mathcal{C}$. Show that L_A is consistent at any $\boldsymbol{\theta} \notin (\mathcal{C} \cap \mathcal{M}^\perp)^\circ$ and that L_B is consistent at any $\boldsymbol{\theta} \notin \mathcal{C}$. Now assume that \mathcal{C} is a polyhedral. Explain why the limiting probability of rejecting the null hypothesis may depend on the particular face of the polyhedral on which the true value lies.

3.30

3.30.1. Let $H_0 : \boldsymbol{\theta} = \mathbf{0}$ and $H_1 : \boldsymbol{\theta} \in \mathcal{C}$ be the null and alternative hypotheses where $\boldsymbol{\theta} \in \mathbb{R}^p$. Assume that $\mathbf{X} = \boldsymbol{\theta} + \mathbf{E}$, where the distribution of \mathbf{E} does not depend on any unknown parameters. Let V be a positive definite matrix, the inner product be defined with respect to $\langle \cdot \rangle_V$, and let

$$L_A = \|\mathbf{X}\|^2 - \|\mathbf{X} - \mathcal{C}\|^2.$$

Then prove the following: (a) $L_A = \|\mathbf{X} - \mathcal{C}^\circ\|^2 = \|\mathbf{E} - \mathcal{P}_{\boldsymbol{\theta}}\|^2$, where $\mathcal{P}_{\boldsymbol{\theta}} = \mathcal{C}^\circ - \boldsymbol{\theta}$.

(b) $\|\boldsymbol{\theta} - \mathcal{C}^\circ\| = \|\mathbf{0} - \mathcal{P}_{\boldsymbol{\theta}}\|$; in other words, the distance between $\boldsymbol{\theta}$ and \mathcal{C}° is equal to the distance between $\mathbf{0}$ and $\mathcal{P}_{\boldsymbol{\theta}}$. (c) For any given $K > 0$, if the distance between $\boldsymbol{\theta}$ and \mathcal{C}° is more than K then $B_K(\mathbf{0}) \cap \mathcal{P}_{\boldsymbol{\theta}}$ is the empty set, where $B_K(\mathbf{0})$ is the ball of radius K centered at $\mathbf{0}$. (d) $L_A \xrightarrow{p} \infty$ as $\|\boldsymbol{\theta} - \mathcal{C}^\circ\| \rightarrow \infty$. (e) For the test that rejects the H_0 for large values of L_A , show that the power tends to 1 as $\|\boldsymbol{\theta} - \mathcal{C}^\circ\| \rightarrow \infty$.

3.30.2. Let $\mathbf{X}, \boldsymbol{\theta}, \mathbf{E}, \mathbf{V}$ and the inner product be as in the previous part. Let $H_0 : \boldsymbol{\theta} \in \mathcal{M}$ and $H_1 : \boldsymbol{\theta} \in \mathcal{C}$, and let

$$L_A = \|\mathbf{X} - \mathcal{M}\|^2 - \|\mathbf{X} - \mathcal{C}\|^2;$$

\mathcal{M} is a linear space, \mathcal{C} is a closed convex cone, and $\mathcal{M} \subset \mathcal{C}$. Consider the test that rejects the H_0 for large values of L_A . Prove that the power of L_A tends to 1 as $\|\boldsymbol{\theta} - (\mathcal{C} \cap \mathcal{M}^\perp)^\circ\| \rightarrow \infty$.

[Hint: Recall that $L_A = \|\mathbf{X} - (\mathcal{C} \cap \mathcal{M}^\perp)\|_0^2$ and hence L_A has the same numerical value as the L_A in the previous part with \mathcal{C} replaced by $\mathcal{C} \cap \mathcal{M}^\perp$. Now apply the previous part.]

3.30.3. Let $\mathbf{X}, \boldsymbol{\theta}, E, V$ and the inner product be as in the previous part. Let $H_1 : \boldsymbol{\theta} \in \mathcal{C}$ and $H_2 : \boldsymbol{\theta} \notin \mathcal{C}$, and let $L_B = \|\mathbf{X} - \boldsymbol{\theta}\|^2$. Show that $L_B \xrightarrow{P} \infty$ and that the power of the test that rejects H_1 for large values of L_B tends to 1 as $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \rightarrow \infty$. [Hint: Apply a method similar to that for the first part].

3.31 Let \mathbf{A} be a nonsingular matrix of order $p \times p$, $\mathbf{x} \in \mathbb{R}^p$ and \mathcal{P} be a cone. Then prove the following:

1. $A\Pi(\mathbf{x}|\mathcal{P}) = \Pi_{AV\mathbf{A}^T}(\mathbf{Ax}|\mathcal{AP})$.
2. $\|\Pi_V(\mathbf{x}|\mathcal{P})\|_V = \|\Pi_{AV\mathbf{A}^T}(\mathbf{Ax}|\mathcal{AP})\|_{AV\mathbf{A}^T}$.
3. $\|\mathbf{x} - \Pi_V(\mathbf{x}|\mathcal{P})\|_V = \|\mathbf{Ax} - \Pi_{AV\mathbf{A}^T}(\mathbf{Ax}|\mathcal{AP})\|_{AV\mathbf{A}^T}$.
4. Let $\mathcal{P}_1 \supset \mathcal{P}_2 \dots \supset \mathcal{P}_n$ be a sequence of decreasing cones in \mathbb{R}^p . Then $\|\mathbf{x} - \Pi_V(\mathbf{x}|\mathcal{P}_n)\|_V$ is a nondecreasing sequence and converges to $\|\mathbf{x} - \Pi_V(\mathbf{x}|\mathcal{P})\|_V$ where $\mathcal{P} = \lim \mathcal{P}_n$; a similar result holds for increasing sequences.

[Perlman (1969)]

3.32 (a) Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$, \mathcal{C} be a closed convex cone, V be a positive definite matrix of order $p \times p$, Π denote Π_V , $\|\cdot\|$ denote $\|\cdot\|_V$, and \mathcal{C}° denote the polar cone of \mathcal{C} with respect to $\langle \cdot, \cdot \rangle_V$. Show that [see Lemma 2.2 in Mukerjee et al. (1986)]

$$\Pi(\mathbf{x} + \mathbf{y} | \mathcal{C}) \leq \|\Pi(\mathbf{x} | \mathcal{C}) + \Pi(\mathbf{y} | \mathcal{C})\| \leq \|\Pi(\mathbf{x} | \mathcal{C})\| + \|\Pi(\mathbf{y} | \mathcal{C})\|.$$

- (b) $\|\Pi(\mathbf{y} + \mathbf{x} | \mathcal{C})\| \leq \|\Pi(\mathbf{y} | \mathcal{C})\|$, $\mathbf{x} \in \mathcal{C}^\circ$.
(c) $\|\mathbf{y} - \mathbf{x}\|^2 \geq \|\Pi(\mathbf{y} | \mathcal{C}) - \Pi(\mathbf{x} | \mathcal{C})\|^2$, $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$.

3.33 Let $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ where Y_1, \dots, Y_n are independent standard normal random variables. Let \mathcal{L} be a linear subspace of \mathbb{R}^n , $\dim(\mathcal{L}) = r$, \mathbf{A} be an $m \times n$ matrix such that each row of \mathbf{A} is orthogonal to \mathcal{L} , and $\tilde{\mathbf{Y}}$ denote the projection of \mathbf{Y} onto \mathcal{L} . Assume that $\{\mathbf{y} : \mathbf{Ay} \geq 0\}$ is nonempty. Then, show that the conditional distribution of $\|\tilde{\mathbf{Y}}\|^2$, given $\mathbf{AY} \geq 0$, is χ_r^2 . [Hint: Use Lemma 3.13.3. Let \mathbf{P} be the symmetric projection matrix onto \mathcal{L} . Then $\mathcal{L} = \{\mathbf{Py} : \mathbf{y} \in \mathbb{R}^n\}$. Verify that if a_i^T is the i^{th} row of \mathbf{A} then $a_i = \mathbf{Py}_i$ for some y_i , $(i = 1, \dots, m)$. Deduce that $\mathbf{A} = \mathbf{BP}$ where $B^T = [\mathbf{y}_1, \dots, \mathbf{y}_m]$, and hence $\mathbf{AY} \geq 0$ is equivalent to $\mathbf{PY} \in \mathcal{C} = \{\mathbf{x} : \mathbf{Bx} \geq 0\}$. For a different proof, see Meyer (2003b).]

3.34 Let $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ where Y_1, \dots, Y_n are independent standard normal random variables. Let \mathcal{L} be a linear subspace of \mathbb{R}^n , $\dim(\mathcal{L}) = r$. Let \mathbf{A} be an $m \times n$ matrix such that each row of \mathbf{A} is a vector in \mathcal{L} . Let \mathbf{Y} denote the projection of \mathbf{Y} onto \mathcal{L} . Assume that $\{\mathbf{Ay} \geq 0\}$ is non-empty. Then, show that the conditional distribution of $\|\tilde{\mathbf{Y}}\|^2$, given $\mathbf{AY} \geq 0$, is χ_r^2 . [Hint: Use Lemma 3.13.3. Let \mathbf{P} be the symmetric projection matrix onto \mathcal{L} . Then $\mathcal{L} = \{\mathbf{Py} : \mathbf{y} \in \mathbb{R}^n\}$. Verify that if a_i^T is the i^{th} row of \mathbf{A} then $a_i = \mathbf{Py}_i$ for some y_i , $(i = 1, \dots, m)$. Deduce that $\mathbf{A} = \mathbf{BP}$ where $B^T = [\mathbf{y}_1, \dots, \mathbf{y}_m]$, and hence $\mathbf{AY} \geq 0$ is equivalent to $\mathbf{PY} \in \mathcal{C} = \{\mathbf{x} : \mathbf{Bx} \geq 0\}$. For a different proof, see Meyer (2003b).]

[Hint: Use Lemma 3.13.3. Let \mathbf{P} be a projection matrix onto \mathcal{L} . Let the i^{th} row of \mathbf{A} be a_i^T . Then $a_i = (\mathbf{I} - \mathbf{P})y_i$ for some y_i , $(i = 1, \dots, m)$. Deduce that $\mathbf{A} = \mathbf{Y}_0^T(\mathbf{I} - \mathbf{P})$, where $\mathbf{Y}_0^T = [\mathbf{y}_1, \dots, \mathbf{y}_m]$, and hence $\mathbf{AZ} \geq 0$ is equivalent to $(\mathbf{I} - \mathbf{P})\mathbf{Z} \in \mathcal{P}$ where $\mathcal{P} = \{\mathbf{x} : \mathbf{Y}_0^T \mathbf{x} \geq 0\}$. Now, use the fact that \mathbf{PZ} and $(\mathbf{I} - \mathbf{P})\mathbf{Z}$ are independent to verify that $\text{pr}(\|\mathbf{PZ}\|^2 \leq c \mid \mathbf{BZ} \geq 0) = \text{pr}(\|\mathbf{PZ}\|^2 \leq c \mid (\mathbf{I} - \mathbf{P})\mathbf{Z} \in \mathcal{P}) = \text{pr}(\|\mathbf{PZ}\|^2 \leq c)$

4

Tests in General Parametric Models

4.1 INTRODUCTION

In Chapter 3, tests of hypotheses concerning the mean of a multivariate normal were considered. Often, the parameter of interest in a hypothesis is more than just the mean of a multivariate normal. Some examples of this type were discussed in Chapter 1. In this chapter, the relevant theory is developed in a general context. To illustrate the nature of the problem, consider Example 1.2.8. It describes the results of a study to establish that a new treatment is better than an old one when the response is an ordinal variable. In this example, the inference problem can be formulated as test of $g(\boldsymbol{\pi}) = \mathbf{0}$ against $g(\boldsymbol{\pi}) \geq \mathbf{0}$ where $\boldsymbol{\pi}$ is the vector consisting of all the $\{\pi_{ij}\}$ and π_{ij} is the probability that the response for Treatment i is the j^{th} ordinal category of the response variable, and g is a nonlinear function; in fact, every component of g is an odds ratio. This example does not fit into the setting in Chapter 3. The theory developed in this chapter is applicable to this and much more general settings. A sample of models to which the theory of this chapter is applicable include the following:

- (1) nonlinear regression
- (2) generalized linear models (logistic regression, log-linear models, etc.)
- (3) proportional hazards model, Cox's regression model
- (4) time series model (for example, ARCH models used in finance)
- (5) mixed effects models (positive semi-definite covariance matrices)
- (6) quasi-likelihood
- (7) generalized estimating equations (GEE)

(8) generalized method of moments (GMM).

It is thus clear that the results of this chapter are applicable to a broad range of contexts.

Let $\boldsymbol{\theta}$ denote the parameter of a statistical model and Θ denote the parameter space. A common feature of most of the inference problems with no inequality constraints on $\boldsymbol{\theta}$ is that Θ is open and the likelihood function is smooth over Θ . This feature plays a crucial role in ensuring that (i) $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \approx N\{\mathbf{0}, \mathcal{I}(\boldsymbol{\theta})^{-1}\}$, and (ii) the asymptotic null distribution of the likelihood ratio statistic for testing $\boldsymbol{\theta}_a = \mathbf{0}$ against $\boldsymbol{\theta}_a \neq \mathbf{0}$ is a chi-square where $\boldsymbol{\theta}_a$ is a subvector of $\boldsymbol{\theta}$. These two results form the basis of statistical inference in these models when there are no inequality constraints on $\boldsymbol{\theta}$.

In this chapter, we will build on these and obtain corresponding results when there are inequality constraints on $\boldsymbol{\theta}$ or the parameter space for $\boldsymbol{\theta}$ is constrained so that the constrained space is not open (topologically). These results will form the basis of inference in this setting. Most of the details in this chapter are given for the case when the observations are independently and identically distributed. However, corresponding results hold for non-*iid* cases, which include regression type models, time series models, and stochastic processes. As a general guide, with $\boldsymbol{\theta}$ denoting an unconstrained estimator of $\boldsymbol{\theta}$, if $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$ is asymptotically normal and the asymptotic null distribution of the likelihood ratio statistic for testing $\boldsymbol{\theta}_a = \mathbf{0}$ against $\boldsymbol{\theta}_a \neq \mathbf{0}$ is a chi-square, where $\boldsymbol{\theta}_a$ is a subvector of $\boldsymbol{\theta}$, then it is very likely that most of the results in this chapter for likelihood ratio tests would be applicable for constrained statistical inference on $\boldsymbol{\theta}$ as well.

The plan of this chapter is as follows. Section 4.2 provides a brief discussion of regularity conditions relating to the likelihood function that would be relevant to the rest of the chapter. Section 4.3 considers statistical models, which are usually referred to as *regular models*, and develops the likelihood ratio, Wald and score-type tests of $R\boldsymbol{\theta} = \mathbf{0}$ against $R\boldsymbol{\theta} \geq \mathbf{0}$ where R is a given matrix of constants. These results are of particular importance because a large number of real-life examples with inequality constraints on $\boldsymbol{\theta}$ are of this type. In any case, a detailed study of the case when the constraints are linear is essential before we consider nonlinear constraints. Section 4.4 extends the results of the previous section to the case when the constraints are nonlinear in $\boldsymbol{\theta}$. In particular, tests of $\boldsymbol{h}(\boldsymbol{\theta}) = \mathbf{0}$ against $\boldsymbol{h}(\boldsymbol{\theta}) \geq \mathbf{0}$ where \boldsymbol{h} is a vector function will be discussed; it will be seen that there are several tests that are asymptotically equivalent to the LRT. Subsection 4.4.2 provides a statistic for testing $\boldsymbol{h}(\boldsymbol{\theta}) = \mathbf{0}$ against $\boldsymbol{h}_2(\boldsymbol{\theta}) \geq \mathbf{0}$, where $\boldsymbol{h}(\boldsymbol{\theta})$ is a smooth vector function of $\boldsymbol{\theta}$ and $\boldsymbol{h}_2(\boldsymbol{\theta})$ is a subvector of $\boldsymbol{h}(\boldsymbol{\theta})$. The main requirement for this is that $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$ be asymptotically normal; thus $\hat{\boldsymbol{\theta}}$ may not be the *mle*. In Section 4.5, a brief overview of likelihood ratio, Wald and score tests is provided with particular reference to the results that would be required for developing a class of score tests against inequality constraints. In Section 4.6, the Rao's score and Neyman's C(α) tests are generalized to tests against inequality constraints. These are further extended to develop tests against inequality constraints when the model is estimated by solving a set of estimating equations; for example, the class of *Generalized Estimating Equation*(GEE) methods and the *Generalized Method of Moments* (GMM) that is widely used in economics, fall into this category.

The results in Sections 4.3–4.6 provide a good coverage of the standard results required for testing $R\boldsymbol{\theta} = \mathbf{0}$ against $R\boldsymbol{\theta} \geq \mathbf{0}$. These results are extended further in the next two sections where the null and alternative hypotheses are stated in the general form, $H_0 : \boldsymbol{\theta} \in \Theta_0$ and $H_1 : \boldsymbol{\theta} \in \Theta_1$, respectively, where $\Theta_0 \subset \Theta_1 \subset \Theta$. In this case, the asymptotic null distributions of test statistics such as the likelihood ratio/Wald/ score, depend on the local shapes of Θ_0 and Θ_1 at the assumed true value in the null parameter space. The relevant local shapes of Θ_0 and Θ_1 are characterized by cones that approximate them at $\boldsymbol{\theta}_0$; these are called *approximating cones*. These will be discussed in detail later in this chapter.

General results for testing $H_0 : \boldsymbol{\theta} \in \Theta_0$ against $H_1 : \boldsymbol{\theta} \in \Theta_1$ are discussed in Section 4.8. A topic of particular interest in economics involves test of $\boldsymbol{h}(\boldsymbol{\theta}) \geq \mathbf{0}$ against $\boldsymbol{h}(\boldsymbol{\theta}) \geq \mathbf{0}$ where \boldsymbol{h} is a given nonlinear function. An example of this type is test of the hypothesis that a certain function (for example, a production function or a cost function) is concave against that it is not concave. The general theory for this is discussed in Section 4.8. The asymptotic distribution of the constrained local and global *mle* are studied in the last section. Some of the proofs are relegated to the Appendix at the end of this chapter.

4.2 PRELIMINARIES

Let Y_1, \dots, Y_n be independently and identically distributed with common density function $f(y; \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p$ and y can be univariate or multivariate. Let

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= \sum \log f(Y_i; \boldsymbol{\theta}), \\ \text{and } \mathcal{I}_{\boldsymbol{\theta}} = \mathcal{I}(\boldsymbol{\theta}) &= E_{\boldsymbol{\theta}}\{(\partial/\partial\boldsymbol{\theta}) \log f(Y; \boldsymbol{\theta})(\partial/\partial\boldsymbol{\theta}^T) \log f(Y; \boldsymbol{\theta})\} \end{aligned} \quad (4.1)$$

denote the loglikelihood, score function, and the information matrix for one observation, respectively. Let $\boldsymbol{\theta}_0$ denote the *true value* and $\hat{\boldsymbol{\theta}}$ denote the *global mle* over Θ .

Let A be a subset of Θ containing the true value of $\boldsymbol{\theta}$. Let $\hat{\boldsymbol{\theta}}_A$ denote the *mle* over A defined by

$$\hat{\boldsymbol{\theta}}_A = \arg \max_{\boldsymbol{\theta} \in A} \ell(\boldsymbol{\theta});$$

equivalently, $\ell(\hat{\boldsymbol{\theta}}_A) = \max_{\boldsymbol{\theta} \in A} \ell(\boldsymbol{\theta})$. Strictly speaking, it is better to define the *mle* over A by $\ell(\hat{\boldsymbol{\theta}}_A) = \sup_{\boldsymbol{\theta} \in A} \ell(\boldsymbol{\theta}) + o_p(1)$; if we do so, $\hat{\boldsymbol{\theta}}_A$ would be well defined even if $\sup_{\boldsymbol{\theta} \in A} \ell(\boldsymbol{\theta})$ is reached at a point on the boundary of A that is not in A . However, we shall not be concerned with this technical detail. Every result in this chapter requires $\hat{\boldsymbol{\theta}}_A$ to be consistent, where A may the natural parameter space or the ones defined by null/alternative hypotheses. For the case when the true value of $\boldsymbol{\theta}$ is an interior point of A , consistency of the *mle* has been discussed in great detail in the literature; for example, see Lehmann (1983, section 6.3). For a discussion of the case when the true parameter may be on the boundary of the parameter space, see Self and Liang (1987, Theorem 1) and Andrews (1999, Section 3.1); van der Vaart and Wellner (1996, Section 3.2) consider the more general case, which includes M-estimation.

Consistency of $\hat{\theta}_A$ would follow from the following two, albeit strong, conditions (see Andrews (1999)): There exists a function $\ell_0(\theta)$ such that

- (i) $\sup_{\theta \in A} |n^{-1}\ell(\theta) - \ell_0(\theta)| \xrightarrow{P} 0$, and
- (ii) $\sup_{\theta \in A, \|\theta - \theta_0\| > \epsilon} \ell_0(\theta) < \ell_0(\theta_0)$, $\epsilon > 0$.

The first condition says that $n^{-1}\ell(\theta)$ converges uniformly over the parameter space, and the second says that θ_0 is the global maximum of $\ell_0(\theta)$ and it is separated from all the other maxima. The uniform law of large numbers of Andrews (1992) may be helpful in establishing the uniform convergence of $n^{-1}\ell(\theta)$. In any case, let us note that if the estimator is the global maximizer of a function, such as the likelihood, then constraints on θ do not generally introduce additional difficulties for establishing consistency. Throughout this chapter, we shall assume the following:

Assumption A1: The mle, $\hat{\theta}_A$, is consistent where $\theta_0 \in A \subset \Theta$. ■

Further, we shall also typically assume that a set of regularity conditions are satisfied to ensure that the asymptotic null distribution of the LRT of $\theta = \theta_0$ against $\theta \neq \theta_0$ is a chi-square, and $\sqrt{n}(\hat{\theta} - \theta_0)$ is asymptotically normal. If this were not true, it is quite unlikely that the results stated in this section about the asymptotic null distribution of LRT for test against inequality restrictions would hold. It is also of interest to note that such a set of regularity conditions is sufficient to derive the asymptotic distributions of LRT under inequality constraints. With this in mind, let us define the following set of regularity conditions:

Condition Q:

1. Distinct values of θ correspond to distinct distributions.
2. The first three partial derivatives of $\log f(y; \theta)$ with respect to θ exist almost everywhere.
3. There exists a $G(y)$ such that $\int G(y)dy < \infty$ and the absolute values of the first three partial derivatives of $\log f(y; \theta)$ with respect to θ are bounded by $G(y)$ in a neighborhood of θ_0 .
4. The Fisher information matrix, $\mathcal{I}(\theta)$, is finite and positive definite.

Remark: In most cases, it would be possible to replace the condition involving the third derivative of $\log f(y; \theta)$ by one that involves only the second derivative. However, we shall not concern ourselves with this technical detail.

In what follows we shall assume that Condition Q is satisfied. Under this condition, we have the following standard results; for example, see Cox and Hinkley (1974), Sen and Singer (1993), Lehmann (1991), and Ferguson (1996).

Proposition 4.2.1 Suppose that Condition Q is satisfied. Then we have the following:

1. $n^{-1/2}S(\theta_0) \xrightarrow{d} N\{\mathbf{0}, \mathcal{I}(\theta_0)\}$.
2. $n^{-1/2}\mathcal{I}(\theta_0)^{-1}S(\theta_0) = n^{1/2}(\hat{\theta} - \theta_0) + o_p(1)$.
3. $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{\mathbf{0}, \mathcal{I}(\theta_0)^{-1}\}$.

4. The LRT for testing $R\theta = 0$ against $R\theta \neq 0$ is asymptotically χ^2 under H_0 where $r = \text{rank}(R)$. ■

This proposition states important first order results that are useful for inference when there are no inequality constraints on θ . Typically, these results are proved using a quadratic approximation of $\ell(\theta)$ and a linear approximation of $\nabla \ell(\theta)$ near θ_0 . Since we will be making use of such quadratic approximations, they are stated in the next proposition.

Proposition 4.2.2 Assume that Condition Q is satisfied. Let $\mathbf{u} = \sqrt{n}(\hat{\theta} - \theta_0)$ and $K > 0$ be given. Then we have the following:

$$(A) \quad \ell(\theta) = \ell(\theta_0) + n^{-1/2}\mathbf{u}^T S(\theta_0) - (1/2)\mathbf{u}^T \mathcal{I}_{\theta_0} \mathbf{u} + r_n(\mathbf{u}), \quad (4.2)$$

where $\sup_{\|\mathbf{u}\| < K} |r_n(\mathbf{u})| = o_p(1)$.

$$(B) \quad \ell(\theta) = \ell(\theta_0) + (1/2)n^{-1}S(\theta_0)^T \mathcal{I}_{\theta_0}^{-1}S(\theta_0) \\ - (1/2)(Z_n - \mathbf{u})^T \mathcal{I}_{\theta_0} (Z_n - \mathbf{u}) + \delta_n(\mathbf{u}), \quad (4.3)$$

where $Z_n = n^{-1/2}\mathcal{I}_{\theta_0}^{-1}S(\theta_0)$ and $\sup_{\|\mathbf{u}\| < K} |\delta_n(\mathbf{u})| = o_p(1)$.

$$(C) \quad \ell(\theta) = \ell(\hat{\theta}) - (1/2)(Z_n - \mathbf{u})^T \mathcal{I}_{\theta_0} (Z_n - \mathbf{u}) + \epsilon_n(\mathbf{u}), \quad (4.4)$$

where $Z_n = \sqrt{n}(\hat{\theta} - \theta_0)$, $\hat{\theta}$ is \sqrt{n} -consistent, and

$$\sup_{\|\mathbf{u}\| < K} |\epsilon_n(\mathbf{u})| = o_p(1).$$

Further, $\delta_n(\mathbf{u}) = n^{-1/2}\|\mathbf{u}\|^3 O_p(1)$. ■

These three quadratic approximations are very important in constrained inference, just as in unconstrained inference. A detailed proof of the first approximation is given in Sen and Singer (1993, Section 5.3); the second approximation is used in Self and Liang (1987), and the third one is used in Silvapulle (1994). For a detailed account see Andrews (1999).

Practically all the asymptotic results of this chapter are based on the assumption that the mle is \sqrt{n} -consistent-i.e., $\sqrt{n}(\hat{\theta}_A - \theta_0) = O_p(1)$, where $\hat{\theta}_A \in A \subset \Theta$. Fortunately, \sqrt{n} -consistency of mle follows from the consistency of the mle and the quadratic approximations in the previous proposition; this is stated in the next result.

Lemma 4.2.3 Suppose that Condition Q is satisfied and $\theta_0 \in A \subset \Theta$. Then, the consistency of $\hat{\theta}_A$ ensures that $\sqrt{n}(\hat{\theta}_A - \theta_0) = O_p(1)$.

Proof: For the case when the observations are iid, see Chernoff (1954). For the non-iid case, see Andrews (1999). For a proof of a similar result, see page 141 in Pollard (1984); see also Lemma 4.3 in Geyer (1994).

The \sqrt{n} -consistency of $\hat{\theta}_A$ enables us to restrict our attention to neighborhoods of the form $\{\theta : \sqrt{n}\|\theta - \theta_0\| < K\}$, where K is arbitrary but fixed, for the purposes of evaluating the global maximum of $\ell(\theta)$. Once we restrict attention to such

a neighborhood, the remainder terms in Proposition 4.2.2 become negligible because they would converge to zero uniformly over the neighborhood. This is the approach that will be used throughout.

A result similar to Lemma 4.2.3 is also available for local mle. Except for the very last section of this chapter, the results in this chapter relate to test of hypotheses and hence local mle would not play a role. However, for completeness let us state the result; the following is Lemma 1 in Self and Liang (1987)

$$D(z) = \min_{R\theta=0} (z - \theta)^T \mathcal{I}_{\theta_0}(z - \theta) - \min_{R\theta \geq 0} (z - \theta)^T \mathcal{I}_{\theta_0}(z - \theta). \quad (4.8)$$

Proposition 4.2.4 Suppose that Condition Q is satisfied and that, for some $\delta_0 > 0$, $B = A \cap \{\theta : \|\theta - \theta_0\| \leq \delta\}$ is a closed set for $0 < \delta < \delta_0$. Then, there exists a sequence of points θ^\dagger in B at which $\ell(\theta)$ has a local maximum and $\theta^\dagger \xrightarrow{p} \theta_0$. Further, $\sqrt{n}(\theta^\dagger - \theta_0) = O_p(1)$. ■

4.3 TESTS OF $R\theta = 0$ AGAINST $R\theta \geq 0$

Let the setting be as in the previous section; in particular, $\theta \in \mathbb{R}^p$ denotes the parameter associated with a statistical model, Θ denotes the parameter space for θ , θ_0 denotes the true parameter, and $\{\ell(\theta), S(\theta), \mathcal{I}(\theta)\}$ be as in (4.1). Let R be a known full-rank matrix of constants of order $r \times p$ ($r \leq p$). In this section, we consider test of

$$H_0 : R\theta = 0 \quad \text{vs} \quad H_1 : R\theta \geq 0. \quad (4.5)$$

In this section, we will assume that Θ is open. The results for the case when Θ is open also hold for many cases when it is not open, but the technical details are more involved. The main advantage of considering the case when Θ is open is that the classical arguments based on Taylor series, for example, as in Cox and Hinkley (1974) and Lehmann (1991), can also be used. For most of this section, the observations will be assumed to be independent and identically distributed.

For the time being, let us note that practically every result in this section for testing (4.5) also holds, with minor modifications, for testing

$$H_0^* : \theta \in \mathcal{M} \quad \text{vs} \quad H_1^* : \theta \in C \quad (4.6)$$

as well where \mathcal{M} is a linear space, C is a closed convex cone, and $\mathcal{M} \subset C$. These results also have generalizations when the C in H_1^* is replaced by a set that may be neither concave nor conical; such generalizations will be discussed in a later section.

Notation: Unless the contrary is clear, the following notation will be adopted: $\hat{\theta}$ is the unconstrained estimator over Θ , and $\hat{\theta}$ and $\bar{\theta}$ denote the estimators under H_1 and H_0 , respectively.

4.3.1 Likelihood Ratio Test with iid Observations

For testing $H_0 : R\theta = 0$ against $H_1 : R\theta \geq 0$, we have

$$LRT = 2[\sup_{R\theta \geq 0} \ell(\theta) - \sup_{R\theta=0} \ell(\theta)]; \quad (4.7)$$

its asymptotic null distribution is given in the next result.

Proposition 4.3.1 Let Y_1, \dots, Y_n be independently and identically distributed with common density function $f(y; \theta)$. Assume that Condition Q is satisfied. Let H_0 and H_1 be as in (4.5). Suppose that H_0 is satisfied and let θ_0 denote the true value in H_0 . Let the function $D(z)$, $z \in \mathbb{R}^p$, be defined by

$$\lim_{n \rightarrow \infty} pr_{\theta_0}(LRT \geq c | H_0) = \sum_{i=0}^r w_i(r, R\mathcal{I}_{\theta_0}^{-1} R^T) pr(\chi_i^2 \geq c). \quad (4.9)$$

Let $Z_n = n^{-1/2} \mathcal{I}_{\theta_0}^{-1} S(\theta_0)$ and $Z \sim N(\mathbf{0}, \mathcal{I}_{\theta_0}^{-1})$. Then, we have the following:

1. $LRT = D(Z_n) + o_p(1)$, and hence $LRT \xrightarrow{d} D(Z)$.

2. The asymptotic null distribution of the LRT is the chi-bar square given by

$$\lim_{n \rightarrow \infty} pr_{\theta_0}(LRT \geq c | H_0) = \sum_{i=0}^r w_i(r, R\mathcal{I}_{\theta_0}^{-1} R^T) pr(\chi_i^2 \geq c). \quad (4.9)$$

Proof: The large sample arguments used here are similar to those used for the case when there are no inequality constraints. Let the notation be as in (4.3). By the same result, we have that

$$\ell(\theta) = A_n - (1/2)q(Z_n - u) + \delta_n(u), \quad (4.10)$$

where $Z_n = n^{-1/2} \mathcal{I}_{\theta_0}^{-1} S(\theta_0)$, $q(Z_n - u) = (Z_n - u)^T \mathcal{I}_{\theta_0}(Z_n - u)$, A_n does not depend on u , and $\sup\{\delta_n(u) : \|u\| < K\} = o_p(1)$. Since $R\theta_0 = 0$, it follows that $R\theta \geq 0$ if and only if $Ru \geq 0$, and $R\theta = 0$ if and only if $Ru = 0$. Now, it follows from the foregoing quadratic approximation that

$$\begin{aligned} LRT &= 2[\sup\{\ell(\theta) : R\theta \geq 0\} - \sup\{\ell(\theta) : R\theta = 0\}] \\ &= \inf\{q(Z_n - u) : Ru = 0\} - \inf\{q(Z_n - u) : Ru \geq 0\} + o_p(1) \\ &= D(Z_n) + o_p(1). \end{aligned}$$

Note that $Z_n \xrightarrow{d} Z$, where $Z \sim N(\mathbf{0}, \mathcal{I}_{\theta_0}^{-1})$, and $D(z)$ is continuous in z (in fact, it can be shown, using the triangular inequality, that $|D(z) - D(x)| \leq \|z - x\|$ and hence $D(z)$ is Lipschitz-continuous). It follows that $D(Z_n)$ and the LRT have a common limiting distribution and it is equal to the distribution of $D(Z)$. By Corollary 3.7.2 (on page 85), the distribution of $D(Z)$ is the chi-bar-square in (4.9). ■

This result is a direct generalization of the classical result that the asymptotic distribution of the likelihood ratio statistic for testing $R\theta = 0$ against $R\theta \neq 0$ is χ_r^2 . In fact, for testing (4.5) and (4.6) the family of chi-bar-square distributions plays a role similar to that of the chi-squared distributions in inference under no inequality constraints.

It would be helpful to note that the main steps in the proof of Proposition 4.3.1 are very similar to those for proving that the *LRT* of $R\theta = 0$ against $R\theta \neq 0$ is asymptotically χ_r^2 under H_0 . In fact, the main idea of the proof carries over to much more general contexts, although some of the technical details would be different.

Let t_{obs} denote the observed sample value of *LRT*. If

$$\sum_{i=0}^r w_i(r, R\mathbf{I}_{\theta_0}^{-1}R^T) \text{pr}(\chi_i^2 \geq t_{obs}) \quad (4.11)$$

does not depend on the assumed true value θ_0 in the null parameter space, then it would be the large sample *p*-value. However, in general, it depends on θ_0 through $\mathcal{I}_{\theta_0}^{-1}$ in $w_i(q, R\mathbf{I}_{\theta_0}^{-1}R^T)$. Thus, θ_0 is a nuisance parameter and hence (4.11) is not a usable definition of a *p*-value. Some methods of dealing with this are discussed in the next subsection.

4.3.2 Tests in the Presence of Nuisance Parameters

Tests of statistical hypotheses in the presence of nuisance parameters are encountered frequently in statistical analysis. The presence of such parameters usually causes difficulties for statistical inference. Let θ be partitioned as $(\lambda : \psi)$ where $\psi \in \mathbb{R}^r$; here, and in what follows, we shall use the following notation for partitioning a vector such as θ to avoid writing the superscript "T" several times:

$$(\lambda : \psi) \text{ denotes } (\lambda^T, \psi^T)^T.$$

For simplicity, we shall consider the special case,

$$H_0 : \psi = 0 \text{ and } H_1 : \psi \geq 0; \quad (4.12)$$

it will be clear that the main procedures discussed below are applicable for testing $\psi = 0$ against $\psi \in C$ where $C \subset \mathbb{R}^r$ and for more general testing problems as well.

Let T be a test statistic such that large values of T provide evidence against H_0 , t denote the sample value of T , and λ_0 denote the true value of λ . Since the discussion of this subsection is centered around *p*-value and probability of Type I error, let us suppose that $H_0 : \psi = 0$ holds; therefore, we let $\theta_0 = (\lambda_0 : 0)$ denote the true value. Suppose that $\text{pr}(T \geq t | \psi = 0, \lambda)$ depends on λ . If λ_0 were known, then the *p*-value would be $\text{pr}(T \geq t | \psi = 0, \lambda = \lambda_0)$. Because λ_0 is unknown, it is not usable as a *p*-value.

Under the null hypothesis, $H_0 : \psi = 0$, the only unknown parameter is λ , and therefore the notation $\text{pr}_\lambda(\cdot | H_0)$ will be used for the probability evaluated under H_0 at $\theta = (\lambda : 0)$. Let

$$p_n(c; \lambda) = \text{pr}_\lambda(LRT \geq c | H_0) \quad \text{and} \quad p_\infty(c; \lambda) = \lim_{n \rightarrow \infty} p_n(c; \lambda).$$

In many cases, closed-form expressions are available for $p_\infty(t; \lambda)$, but not for $p_n(t; \lambda)$.

In this case, our main interest would be on large sample tests based on $p_\infty(t; \lambda)$.

By Proposition 4.3.1,

$$p_\infty(t; \lambda) = \sum_{i=0}^r w_i(r, \mathcal{T}^{\psi\psi}(\lambda)) \text{pr}(\chi_i^2 \geq t) \quad (4.13)$$

where $\mathcal{T}^{\psi\psi}(\lambda)$ is the (ψ, ψ) block of $\mathcal{I}^{-1}(\theta)$ evaluated at $\theta = (\lambda : 0)$. If $p_\infty(t; \lambda)$ does not depend on λ then the large sample *p*-value corresponding to *LRT* = t is $p_\infty(t; \lambda)$. Therefore, for the following discussions, we shall consider only the case when $p_\infty(c; \lambda)$ does depend on λ .

Bounds test:

Let $p_{n,L}(t)$ and $p_{n,U}(t)$ be positive numbers such that

$$p_{n,L}(t) \leq p_n(t; \lambda) \leq p_{n,U}(t) \quad \text{for any } \lambda.$$

Then, the *bounds test* at level α is the following: reject H_0 if $p_{n,U}(t) \leq \alpha$, do not reject H_0 if $p_{n,L}(t) > \alpha$ and the test is inconclusive if $p_{n,L}(t) \leq \alpha < p_{n,U}(t)$. Usually, it is difficult to obtain sharp bounds for $p_n(t; \lambda)$. If the sample size is large then $p_n(t; \lambda)$ can be approximated by the tail probability of the chi-bar-square distribution in (4.13), and hence the bounds in (3.23) can be used for a large sample bounds test.

Supremum of pr(Type I Error) over the null parameter space (p_{sup} -test):

A test, which we call the p_{sup} -test, is based on

$$p\text{-value} = \sup_{\lambda \in \Lambda} p_n(t; \lambda) \quad (4.14)$$

where t is the sample value of *LRT* and Λ is the parameter space for λ , for example, see Lehmann (1983). This is the generally accepted procedure for exact inference.

The argument underlying this procedure is that if the value of t for *LRT* is consistent with at least one point in the null parameter space then we should not reject H_0 . A large sample approximation of (4.14) is

$$\text{reject } H_0 \text{ if } \sup_{\lambda \in \Lambda} p_\infty(t; \lambda) \leq \alpha. \quad (4.15)$$

Ideally, it would be better to approximate $\lim_n \sup_{\lambda \in \Lambda} p_n(t; \lambda)$ than $\sup_{\lambda \in \Lambda} \lim_n p_n(t; \lambda)$.

It does not appear that this has been investigated for tests against inequality constraints. In some cases, $p_\infty(t; \lambda)$ may vary substantially as a function of λ in Λ and hence this test may be too conservative. Let us illustrate this using an example. Consider the Perlman's test of $\psi = 0$ against $\psi \geq 0$ based on a random sample from $N(\psi, \lambda)$, where λ is an unknown $q \times q$ covariance matrix. For Perlman's statistic, \mathcal{U} , which is equivalent to the *LRT*, the finite sample null distribution is the distribution given by (see Proposition 3.10.3 on page 103)

$$p_n(t; \lambda) = \text{pr}_\lambda(\mathcal{U} \geq t | H_0) = \sum_{i=1}^r w_i(r, \lambda) \text{pr}(\chi_i^2 / \chi_{n-r}^2 \geq t); \quad (4.16)$$

it was shown in (3.56) that

$$p_{sup} = \sup_{\lambda \in \Lambda} p_n(t; \lambda) = 0.5 \{ \text{pr}(\chi_{r-1}^2 / \chi_{n-r}^2 \geq t) + \text{pr}(\chi_r^2 / \chi_{n-r}^2 \geq t) \}. \quad (4.17)$$

It can be shown that the supremum of $p_n(t; \lambda)$ over λ is achieved when the correlation matrix corresponding to λ is equal to J , where J is the matrix in which every element is equal to 1 [Silvapulle (1996b), Perlman (1969)]. Clearly, J is an extreme form for a correlation matrix. The true correlation matrix in most practical situations is unlikely to be close to J . If we restrict the parameter space for λ away from a neighborhood of J , then the supremum of $p_n(t; \lambda)$ would be smaller than the expression on the right-hand side of (4.17). Consequently, the test in (4.14) with p_{sup} as in (4.17) appears to be conservative for most practical applications. For a discussion of this example, see Silvapulle (1996b).

There are some values of λ that one would be able to eliminate without looking at the data. After looking at the data, one may be tempted to eliminate values of λ that are far from an estimate of λ ; however, this needs to be done carefully to ensure that the probability of Type I error does not exceed the nominal level. One approach that may overcome the conservative nature of p_{sup} -test is to approximate the p -value by the supremum of $\text{pr}(LRT \geq t|\lambda)$ over a confidence region for λ . The p^* -test discussed below provides such an approach.

Supremum of $\text{pr}(\text{Type I Error})$ over a confidence region (p^* -test):

In this method, the test is carried out in two stages. For illustrative purposes, suppose that we wish to test H_0 vs H_1 at level 0.05. In the first stage, we construct a confidence region for λ , say a 99% confidence region; let this be denoted by \mathcal{A} . Let $p^* = \{(1 - 0.99) + \sup_{\lambda \in \mathcal{A}} p_n(t; \lambda)\}$. Now, in the second stage, a test at level 0.05 is

$$\text{reject } H_0 \text{ if } p^* \leq 0.05.$$

An interpretation of this is the following: In the first stage, when we replace the full parameter space for λ by a 99% confidence region for it, the “probability” of an error is $(1 - 0.99)$. In the second stage, we reject H_0 if $\sup_{\lambda \in \mathcal{A}} p_n(t; \lambda) \leq 0.04$, so that the combined probability of Type I error is not more than $(1 - 0.99) + 0.04 = 0.05$. We call this, the p^* -test, and the foregoing p^* is an upper bound for the p -value of this test. Obviously, it is also possible to use a 98% confidence region for Stage 1 and carry out the test in Stage 2 with $p^* = \{(1 - 0.98) + \sup_{\lambda \in \mathcal{A}} p_n(t; \lambda)\}$. A rigorous statement of this is given below in Proposition 4.3.2.

Suppose that we wish to test

$$H_0 : \psi = 0 \text{ against } H_1 : \psi \in C$$

for some C . Let T_λ be a test statistic for a given λ , and let t_λ denote the sample value of T_λ ; note that the test statistic may depend on λ . For simplicity, assume that large values of T_λ favor H_1 . Let α_1 be a given small number in the range $0 \leq \alpha_1 < 1$. Let \mathcal{A} be a random subset of Λ such that

$$\text{pr}(\lambda_0 \in \mathcal{A} | H_0) \geq 1 - \alpha_1.$$

Thus, once the sample values have been substituted, \mathcal{A} becomes a $100(1 - \alpha_1)\%$ confidence region for λ_0 under H_0 . Let

$$p^* = \alpha_1 + \sup_{\lambda \in \mathcal{A}} \text{pr}\{T_\lambda \geq t_\lambda | \theta = (\lambda : \mathbf{0})\}.$$

Now, a level- α p^* -test is

$$\text{reject } H_0 \text{ if } p^* \leq \alpha. \quad (4.18)$$

This is stated in the next result, and the proof is given in the appendix to this chapter (see page 215).

Proposition 4.3.2 *The size of the p^* -test in (4.18) does not exceed α . Consequently, an upper bound for the p -value of this test is p^* .*

Although the p^* -test was motivated by the need for a test with better power than the p_{sup} -test, there is no guarantee that the former would have such a power advantage. However, it is worth noting that the power of p^* -test at level α is no less than that of the p_{sup} -test at level $(\alpha - \alpha_1)$. This suggests that the power of the p^* -test cannot be much worse than that of the p_{sup} test provided α_1 is small. Simulations in Silvapulle (1996b) suggest that the p^* -test has better power than the p_{sup} -test in the specific cases considered; for illustrative examples of this method and further discussions, see Silvapulle (1996b) and Berger and Boos (1994).

Point estimate of the p -value (\hat{p} -test):

Because λ_0 is unknown, one procedure is to use an estimated value of $p_n(t; \lambda_0)$ under H_0 . Let $\hat{\lambda}$ be a consistent estimate of $p_n(t; \lambda_0)$ under H_0 ; for example, it could be $p_n(t; \hat{\lambda})$ or $p_\infty(t; \hat{\lambda})$ where $\hat{\lambda}$ is a consistent estimate of λ under H_0 . Now the large sample test, which we call the \hat{p} -test, is

$$\text{reject } H_0 \text{ if } \hat{p} \leq \alpha. \quad (4.19)$$

In small samples, this procedure must be used with caution because there is a possibility that the size of the test may exceed the nominal level α substantially. For example, for testing $\psi = \mathbf{0}$ against $\psi \geq \mathbf{0}$ based on a sample from $N(\psi, \lambda)$, the tail probability

$$\sum w_i(r, \hat{\lambda}) \text{pr}(\chi_i^2 / \chi_{n-r}^2 \geq t)$$

with $\hat{\lambda}$ being a consistent estimator of λ may turn out to be a poor estimate of

$$\sum w_i(r, \lambda_0) \text{pr}(\chi_i^2 / \chi_{n-r}^2 \geq t)$$

if the chi-bar-square weights, $w_i(r, \lambda)$, are sensitive to the form of the covariance matrix λ and $\hat{\lambda}$ is not close to λ_0 ; this was illustrated in Silvapulle (1996b).

4.3.3 Wald- and Score-type Tests with iid Observations

In the standard inference literature, Wald and score tests of $R\theta = 0$ against $R\theta \neq 0$ are known to be asymptotically equivalent to the likelihood ratio test. These tests may be generalized in different directions when there are inequality constraints. In this section, we will introduce Wald- and score-type tests that are asymptotically equivalent to the likelihood ratio test for testing against inequality constraints.

Wald-type tests

The usual Wald statistic for testing $R\theta = 0$ against $R\theta \neq 0$ is

$$n(R\hat{\theta})^T (R\hat{T}^{-1} R^T)^{-1} (R\hat{\theta}) \quad (4.20)$$

where \hat{T} is a consistent estimator of $\mathcal{I}(\theta)$. The following interpretations of this statistic are helpful; they are based on $\sqrt{n}(\hat{\theta} - \theta) \approx N\{0, \mathcal{I}(\theta)^{-1}\}$ and $\sqrt{n}(R\hat{\theta} - R\theta) \approx N(0, A)$ where $A = \{R\mathcal{I}(\theta)^{-1} R^T\}$:

1. It [i.e., (4.20)] is a measure of $\|R\hat{\theta}\|^2$. Since $R\hat{\theta}$ is an estimate of $R\theta$ obtained without imposing the constraints in the null hypothesis, it is reasonable to expect that $\|R\hat{\theta}\|^2$ would be small if $R\theta = 0$; and it would be large if $R\theta \neq 0$.

2. First, construct the LRT ($= n(R\hat{\theta})^T (R\mathcal{V} R^T)^{-1} (R\hat{\theta})$) for testing $R\theta = 0$ against $R\theta \neq 0$ based on a single observation of $\hat{\theta}$ under the assumption that $\hat{\theta}$ is exactly $N(\theta, V)$, where V is assumed known. Now, (4.20) is a large sample approximation to this statistic obtained by substituting $V = n^{-1}\hat{\mathcal{I}}_{-1}$.

3. It is a measure of $\{d(\hat{\theta}, H_0) - d(\hat{\theta}, H_2)\}$, where $d(\hat{\theta}, H_i)$ is a measure of the distance between $\hat{\theta}$ and the parameter space defined by the hypothesis H_i , ($i = 0, 1$); (4.20) is obtained by choosing

$$d(\hat{\theta}, H) = \min\{n(\hat{\theta} - \theta)^T \hat{\mathcal{I}}(\hat{\theta} - \theta) : \theta \in H\}.$$

For testing $R\theta = 0$ against $R\theta \geq 0$, a statistic that resembles (4.20) and is along the lines of the first interpretation (of the foregoing three) is

$$W = n(R\hat{\theta})^T (R\hat{T}^{-1} R^T)^{-1} (R\hat{\theta}) \quad (4.21)$$

where \hat{T} is an estimator of \mathcal{I} under $H_1 : R\theta \geq 0$; thus, \hat{T} can be $\mathcal{I}(\hat{\theta})$ or $\mathcal{I}(\hat{\theta})$. We shall call this a Wald-type statistic. Note that the definition in (4.21) does not follow from the asymptotic distribution of $R\hat{\theta}$ in the same way as (4.20) does from that of $R\hat{\theta}$ because $R\hat{\theta}$ is not asymptotically normal. It will be seen that the W in (4.21) is asymptotically equivalent to LRT .

A statistic along the lines of the third interpretation is

$$D = \inf_{R\theta=0} \{n(\hat{\theta} - \theta)^T \hat{\mathcal{I}}(\hat{\theta} - \theta)\} - \inf_{R\theta \geq 0} \{n(\hat{\theta} - \theta)^T \hat{\mathcal{I}}(\hat{\theta} - \theta)\}. \quad (4.22)$$

We shall call this a *Distance Statistic*.

Because the parameter θ enters the hypotheses only through $R\theta$, we can think of $R\theta$ as the parameter of interest, and define a statistic along the lines of the second interpretation with $R\theta$ playing the role of the parameter of interest. This leads to a statistic which is equal to D in (4.22) although its functional form is different:

$$D = \inf_{R\theta=0} \{n(R\hat{\theta} - R\theta)^T (R\hat{T}^{-1} R^T)^{-1} (R\hat{\theta} - R\theta)\} - \inf_{R\theta \geq 0} \{n(R\hat{\theta} - R\theta)^T (R\hat{T}^{-1} R^T)^{-1} (R\hat{\theta} - R\theta)\}. \quad (4.23)$$

It will be seen later that D is also asymptotically equivalent to the *LRT*.

A Global score statistic

Recall that the score statistic for testing $\theta = \theta_0$ against $\theta \neq \theta_0$ is

$$n^{-1} S(\theta_0)^T \mathcal{I}_{\theta_0}^{-1} S(\theta_0) \quad (4.24)$$

where $S(\theta) = (\partial/\partial\theta)\ell(\theta)$ is the score function. The motivation for this is that if θ_0 is the true value then $E\{\ell(\theta)\}$ has a global maximum and a stationary value at θ_0 , $n^{-1/2} S(\theta_0) \xrightarrow{d} N\{0, \mathcal{I}_{\theta_0}\}$, and $S(\theta_0)$ tends to be close to (respectively, far from) 0 when θ_0 is the true (respectively, not true) value. One possible way of extending this to the inequality constrained testing problem is to assess the extent to which $\{S(\bar{\theta}) - S(\tilde{\theta})\}$ is different from zero, where $\bar{\theta}$ and $\tilde{\theta}$ are estimates of θ under the null and alternative hypotheses, respectively. If H_0 is true, then $\bar{\theta}$ and $\tilde{\theta}$ are expected to be close and hence $\{S(\bar{\theta}) - S(\tilde{\theta})\}$ is expected to be close to zero; on the other hand, if H_0 is not true then $S(\bar{\theta})$ is expected to be close to zero but not $S(\tilde{\theta})$ and therefore, $\{S(\bar{\theta}) - S(\tilde{\theta})\}$ is expected to be away from zero. This suggests that a test can be based on the difference $\{S(\bar{\theta}) - S(\tilde{\theta})\}$. Instead of the score function $S(\theta)$ we may use the effective score, $R\mathcal{Z}(\theta_0)^{-1} S(\theta_0)$. This suggests the following:

$$S_G = [R\hat{T}^{-1}\{S(\bar{\theta}) - S(\tilde{\theta})\}]^T (R\hat{T}^{-1} R^T)^{-1} [R\hat{T}^{-1}\{S(\bar{\theta}) - S(\tilde{\theta})\}]. \quad (4.25)$$

The following form has also been suggested (see Robertson et al. (1988)):

$$\{S(\bar{\theta}) - S(\tilde{\theta})\}^T \mathcal{I}(\bar{\theta})^{-1} \{S(\bar{\theta}) - S(\tilde{\theta})\}.$$

The reason for referring to these as “Global” score statistics is that they use information about the shape of the score function over the whole parameter space. This is in contrast to the Rao score statistic and Neyman’s $C(\alpha)$ statistic, which are based only on the properties of the score function in a local neighborhood of the assumed true null value.

It is well known that the likelihood ratio, Wald, and score tests are asymptotically equivalent for testing $R\theta = 0$ against $R\theta \neq 0$. The next proposition says that a corresponding result also holds for the foregoing Wald-type statistic in (4.21), the distance statistic (4.22), and the score statistic in (4.25). In fact, this result holds even in some regression-type models such as the generalized linear models (see Silvapulle (1994)); this will be discussed in the next subsection.

Proposition 4.3.3 Suppose that Condition Q is satisfied. Then for testing $H_0 : R\theta = 0$ vs $H_1 : R\theta \geq 0$, we have the following:

$$LRT = W + o_p(1) = D + o_p(1) = S_G + o_p(1) \text{ under } H_0. \quad (4.26)$$

Consequently, the common asymptotic null distribution of *LRT*, *W*, *D*, and S_G is the chi-bar square in (4.9).

Proof: Follows from the local quadratic approximations of loglikelihood in Proposition 4.2.2 (see Gouriéroux and Monfort (1995) and Silvapulle (1994)). ■

It follows that, under some reasonable regularity conditions, the test statistics LRT , W , D , and S_G are asymptotically equivalent in the sense that they have the same asymptotic local power, more precisely Pitman asymptotic efficiency against local alternatives. Although the four statistics in (4.26) are asymptotically equivalent, very little is known about their relative advantages. From a computational point of view, D is easier to use than the others because it requires only the unconstrained estimator, $\hat{\theta}$, which is usually easier to compute than the constrained estimators. The statistics S_G and W do not have any significant computational advantage over LRT because all of them require the constrained estimator. In contrast to S_G , W , and D , LRT has the advantage that it is invariant under a reparameterization. If a test that is computationally simpler than the LRT is required, then D is worth considering. In view of these, at this stage, S_G and W do not appear to be serious competitors to LRT or D .

For testing $R\theta = 0$ against $R\theta \neq 0$, the classical score statistic, also known as Rao's score statistic and the Lagrange multiplier statistic, provides a simple way of testing when the full model is complicated; the implementation requires estimation of the model only under H_0 . This is important when the full model is significantly more complicated than the null model; this particular feature of the score statistic has been crucial for its popularity. Unfortunately, the global score statistic, S_G , does not have this simplicity. A natural generalization of the Rao's score statistic, which preserves the attractive feature (i.e., requires only the null model to be estimated) of the classical Rao's score statistic will be introduced in a later section.

4.3.4 Independent Observations with Covariates

Most of the foregoing results for independently and identically distributed observations extend in a natural way to other more general cases with non-*iid* observation, for example, linear and nonlinear regression models, generalized linear models, and some stochastic processes.

Let Y_1, \dots, Y_n be independent, $f_i(y_i; \theta)$ denote the density function of Y_i , and $\ell(\theta)$ denote the loglikelihood, $\sum \log f_i(Y_i; \theta)$. This includes the setting where the observations are (y_i, x_i) , $i = 1, \dots, n$, with x denoting a covariate; thus, regression-type models are included. Recall that in the previous subsections, we defined $\mathcal{I}(\theta)$ as the information per observation; it was the same for every observation because they were identically distributed. Now, since the observations are not identically distributed, we need to make some modifications. Let

$$\nabla = (\partial/\partial\theta) \text{ and } \nabla^2 = (\partial^2/\partial\theta\partial\theta^T).$$

Suppose that

$$n^{-1/2}\nabla\ell(\theta) \xrightarrow{d} N(\mathbf{0}, \mathcal{V}_\theta), \quad \text{and} \quad n^{-1}\nabla^2\ell(\theta) \xrightarrow{a.s.} -\mathcal{V}_\theta \quad (4.27)$$

for some positive definite matrix \mathcal{V}_θ . As was introduced at the beginning of this section, let $\hat{\theta}$ and $\bar{\theta}$ denote the maximizers of $\ell(\theta)$ over Θ , $\{R\theta \geq 0\}$ and $\{R\theta$

$0\}$, respectively. For testing

$$H_0 : R\theta = 0 \text{ against } H_1 : R\theta \geq 0, \quad (4.28)$$

the likelihood ratio statistic is

$$LRT = 2\{\ell(\tilde{\theta}) - \ell(\bar{\theta})\}. \quad (4.29)$$

We shall continue to assume that Assumption A1 on page 146 holds with appropriate modifications. Assume that (a) the quadratic approximation (4.3) in Proposition 4.2.2 (page 147) holds with $\mathcal{I}(\theta)$ replaced by $\mathcal{V}(\theta)$, (b) $\mathcal{V}(\theta)$ is continuous in θ and $(c)\hat{\theta}$ and $\bar{\theta}$ are $n^{1/2}$ -consistent. The proof of the quadratic approximation may involve lengthy technical details; for example, for the case of generalized linear models, see Fahrmeir and Kaufmann (1985). The \sqrt{n} -consistency of $\hat{\theta}$ and $\bar{\theta}$ would follow from the consistency of the global *mle* and the quadratic approximation, in much the same way as in the *iid* case. Thus, conceptually, the conditions required for the regression type setting is not that different from that for the *iid* setting but the technical details could be complicated.

Now, the proofs of Propositions 4.3.1 and 4.3.3 still hold with appropriate modifications; see Silvapulle (1994) for details.

Proposition 4.3.4 *The asymptotic null distribution of LRT for testing*

$$H_0 : R\theta = 0 \text{ against } H_1 : R\theta \geq 0$$

is equal to the distribution of

$$\inf_{R\theta=0} \{(Z - \theta)^T \mathcal{V}_{\theta_0}(Z - \theta)\} - \inf_{R\theta \geq 0} \{(Z - \theta)^T \mathcal{V}_{\theta_0}(Z - \theta)\}$$

where $Z \sim N(\mathbf{0}, \mathcal{V}_{\theta_0}^{-1})$. Consequently, as $n \rightarrow \infty$,

$$pr_{\theta_0}(LRT \geq c \mid H_0) \rightarrow \sum_{i=0}^r w_i(r, R\mathcal{V}_{\theta_0}^{-1}R^T) pr(\chi_i^2 \geq c). \quad \blacksquare$$

It appears that, if

1. $\sqrt{n}(\hat{\theta} - \theta_0)$ converges to $N(\mathbf{0}, \mathcal{V}_{\theta_0}^{-1})$, and
for testing $H_0 : R\theta = 0$ against $R\theta \neq 0$ is a chi-square,
2. the asymptotic null distribution of the likelihood ratio statistic, $2\{\ell(\hat{\theta}) - \ell(\bar{\theta})\}$

then it is almost certain that the conclusions of the foregoing proposition also hold. If the parameter space for θ is open and the asymptotic null distribution of $2\{\ell(\hat{\theta}) - \ell(\bar{\theta})\}$ for testing $H_0 : R\theta = 0$ against $H_2 : R\theta \neq 0$ is not a chi-square, then the results of Proposition 4.3.4 are unlikely to hold.

It is also possible to construct Wald-type, score-type, and distance-based tests as follows:

$$\begin{aligned} W &= n(\mathbf{R}\hat{\theta})^T (\mathbf{R}\hat{V})^{-1} \mathbf{R}^T)^{-1} (\mathbf{R}\hat{\theta}) \\ S_G &= n^{-1} \mathbf{U}^T (\mathbf{R}\hat{V})^{-1} \mathbf{R}^T)^{-1} \mathbf{U}, \quad \text{where } U = \mathbf{R}\hat{V}^{-1} \{\mathbf{S}(\hat{\theta}) - \mathbf{S}(\bar{\theta})\} \\ D &= \{d(\hat{\theta}, H_0) - d(\hat{\theta}, H_1)\}, \quad \text{where } d(\hat{\theta}, H) = \inf_{\theta \in H} n(\hat{\theta} - \theta)^T \hat{V}(\hat{\theta} - \theta). \end{aligned} \quad (4.30)$$

Now, it can be shown that (see Silvapulle (1994), Gourieroux and Monfort (1995))

$$LRT = W + o_p(1) = S_G + o_p(1) = D + o_p(1). \quad (4.31)$$

Therefore, these tests have the same asymptotic null distribution. Further, under some weak conditions, (4.31) also hold under Pitman-type local alternatives and hence the tests in (4.31) also have the same asymptotic local power.

Least squares in linear regression

The foregoing results extend in a natural way to some quasi-likelihood models as well. To illustrate the main ideas, let us first consider the linear regression model,

$$y_i = \mathbf{x}_i^T \boldsymbol{\theta} + \epsilon_i,$$

where the error terms are *iid* with common variance σ^2 . Let

$$R_n(\boldsymbol{\theta}) = (-1/2) \sum (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2 \quad \text{and} \quad \mathbf{W} = \lim(n^{-1} \mathbf{X}^T \mathbf{X}),$$

where \mathbf{X} is the design matrix $(\mathbf{x}_1, \dots, \mathbf{x}_n)^T$. Let $\hat{\boldsymbol{\theta}}$, $\bar{\boldsymbol{\theta}}$, and $\bar{\boldsymbol{\theta}}$ be the estimators corresponding to the objective function $R_n(\boldsymbol{\theta})$; thus $\hat{\boldsymbol{\theta}}$ is simply the unconstrained ordinary least squares estimator. Let $\mathbf{S}(\boldsymbol{\theta}) = \nabla R_n(\boldsymbol{\theta})$. Then we have that $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N(\mathbf{0}, \sigma^2 \mathbf{W}^{-1})$,

$$n^{-1/2} \mathbf{S}(\boldsymbol{\theta}_0) \xrightarrow{d} N(\mathbf{0}, \sigma^2 \mathbf{W}) \quad \text{and} \quad -n^{-1} \nabla^2 R_n(\boldsymbol{\theta}) \xrightarrow{p} \mathbf{W}.$$

Note that the asymptotic covariance of $\mathbf{S}(\boldsymbol{\theta}_0)$ and the limit of $-n^{-1} \nabla^2 R_n(\boldsymbol{\theta}_0)$ are proportional, but not equal; if R_n were the loglikelihood then they would have been equal. Now, for testing $H_0 : R\theta = 0$ against $H_2 : R\theta \neq 0$, an LR-type statistic and its asymptotic null distribution are given by

$$T_n = 2\{R_n(\hat{\boldsymbol{\theta}}) - R_n(\bar{\boldsymbol{\theta}})\}/\hat{\sigma}^2 \xrightarrow{d} \chi_r^2 \quad \text{under } H_0,$$

where $\hat{\sigma}$ is a consistent estimator of σ .

Note that

$$T_n = \{\|(\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\theta}})/\sigma\|^2 - \|(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\theta}})/\sigma\|^2\}(\sigma^2/\hat{\sigma}^2),$$

$(\sigma^2/\hat{\sigma}^2) \xrightarrow{p} 1$, and the distribution of the term in $\{\dots\}$ does not depend on σ . Thus, essentially σ^2 factors out and cancels with $\hat{\sigma}^{-2}$ in the limit. This illustrates that in the

linear model, we do not need to know the exact population distribution to construct tests of $H_0 : R\theta = 0$ against $H_1 : R\theta \neq 0$.

In what follows we extend Proposition 4.3.4 to incorporate this type of setting.

A quasi-likelihood method

Consider a general parametric model with unknown parameter $\boldsymbol{\theta}$. Let $R_n(\boldsymbol{\theta})$ denote an objective function which will be maximized for estimating $\boldsymbol{\theta}$. Let $\hat{\boldsymbol{\theta}}$, $\bar{\boldsymbol{\theta}}$ be the estimators corresponding to this objective function. Let $\mathbf{S}_n(\boldsymbol{\theta}) = \nabla R_n(\boldsymbol{\theta})$. Suppose that

$$n^{-1/2} \mathbf{S}(\boldsymbol{\theta}_0) \xrightarrow{d} N\{\mathbf{0}, \sigma^2 \mathbf{W}(\boldsymbol{\theta}_0)\} \quad \text{and} \quad n^{-1} \nabla^2 R_n(\boldsymbol{\theta}) \xrightarrow{p} -\mathbf{W}(\boldsymbol{\theta}). \quad (4.31)$$

Let $\mathbf{Z}_n(\boldsymbol{\theta}_0) = n^{-1/2} \mathbf{W}(\boldsymbol{\theta}_0)^{-1} \mathbf{S}(\boldsymbol{\theta}_0)$. Assume that R_n satisfies regularity conditions similar to those required for the *iid* setting. Therefore, we have $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N\{\mathbf{0}, \sigma^2 \mathbf{W}(\boldsymbol{\theta}_0)^{-1}\}$, $\mathbf{Z}_n(\boldsymbol{\theta}_0) \xrightarrow{d} N\{\mathbf{0}, \sigma^2 \mathbf{W}(\boldsymbol{\theta}_0)^{-1}\}$, and a quadratic approximation similar to (4.3) holds for $R_n(\boldsymbol{\theta})$ with the information matrix \mathcal{I} replaced by \mathbf{W} .

Now consider test of $H_0 : R\theta = 0$ vs $H_1 : R\theta \geq 0$. A suitable test statistic for this is

$$T_n = 2\{R_n(\hat{\boldsymbol{\theta}}) - R_n(\bar{\boldsymbol{\theta}})\}/\hat{\sigma}^2,$$

where $\hat{\sigma}^2$ is a consistent estimator of σ . Then we have

$$\begin{aligned} T_n &= 2\{R_n(\hat{\boldsymbol{\theta}}) - R_n(\bar{\boldsymbol{\theta}})\}/\sigma^2 + o_p(1) \\ &= D(\mathbf{Z}_n)/\sigma^2 + o_p(1) = D(\mathbf{Z}_n/\sigma) + o_p(1) \\ &\xrightarrow{d} D(U) \quad \text{under } H_0 \end{aligned} \quad (4.32)$$

where $D(\mathbf{z})$ is as in (4.8) with \mathcal{I}_{θ_0} replaced by $\mathbf{W}(\boldsymbol{\theta}_0)$, and $U \sim N\{\mathbf{0}, \mathbf{W}(\boldsymbol{\theta}_0)^{-1}\}$. Therefore,

$$T_n \xrightarrow{d} \bar{\chi}^2 \{(\mathbf{R}\mathbf{W}(\boldsymbol{\theta}_0)^{-1} \mathbf{R}^T, \mathbb{R}^{+r})\} \quad \text{under } H_0.$$

We can also define Wald, score, and distance type test statistics as in (4.30). All that we need to do is to introduce $\hat{\sigma}^{-2}$ on the RHS of the definitions of W , S_G , and D in (4.30) and replace \mathcal{V} by $\hat{\mathcal{V}}$, a consistent estimator of $\mathbf{W}(\boldsymbol{\theta}_0)$:

$$\begin{aligned} W &= \hat{\sigma}^{-2} n(R\hat{\boldsymbol{\theta}})^T (R\hat{\mathbf{W}}^{-1} R^T)^{-1}(R\hat{\boldsymbol{\theta}}) \\ S_G &= \hat{\sigma}^{-2} n^{-1} \mathbf{U}^T (R\hat{\mathbf{W}}^{-1} R^T)^{-1} \mathbf{U}, \quad \text{where } U = R\hat{\mathbf{W}}^{-1} \{\mathbf{S}(\hat{\boldsymbol{\theta}}) - \mathbf{S}(\bar{\boldsymbol{\theta}})\} \\ D &= \hat{\sigma}^{-2} \{d(\hat{\boldsymbol{\theta}}, H_0) - d(\hat{\boldsymbol{\theta}}, H_1)\}, \quad \text{where } d(\hat{\boldsymbol{\theta}}, H) = \inf_{\theta \in H} n(\hat{\boldsymbol{\theta}} - \theta)^T \hat{\mathbf{W}}(\hat{\boldsymbol{\theta}} - \theta). \end{aligned} \quad (4.33)$$

The rest of the details for the asymptotic equivalence of T_n , W , S_G , and D remain essentially the same; see Silvapulle (1994) and Gouriéroux and Monfort (1995) for details. Finally, let us note that, for a generalized linear model with canonical link, the foregoing results are adequate to apply quasi-likelihood ratio test in the presence of an overdispersion parameter, σ .

Robust test in the linear model

The foregoing results can be used for developing *robust tests* of $R\theta = 0$ vs $R\theta \geq 0$ in the linear regression model $\mathbf{Y} = \mathbf{X}\theta + \mathbf{E}$. Suppose that the objective function is

$$R_n(\theta) = -\sum \rho\{(y_i - \mathbf{x}_i^T \theta)/s\} s,$$

where ρ is a suitably chosen convex function and s is an estimate of scale, for example, a multiple of the median of the absolute values of the least square residuals. Now, it may be verified that the foregoing tests of $R\theta = 0$ vs $R\theta \geq 0$ are applicable. See Silvapulle (1992b) and Silvapulle (1992c) for details where proofs of the asymptotic equivalence of the tests are also given; these provide details for the statistics of the form T_n and D , respectively. It is shown that the robustness properties of the estimator carry over to the corresponding test statistics T_n and D . These results also hold in *group sequential analysis*; see Silvapulle and Sen (1993).

4.3.5 Examples

Example 4.3.1

An assay was carried out with the bacterium *E. coli* to evaluate the genotoxic effects of 9-AA and potassium chromate. Further background and the data are given in Example 1.2.7 of Chapter 1 (page 11). The objective is to test whether potassium chromate and 9-AA have a synergistic effect. The *simple independent action* (SIA) model has attracted interest for studies of the type in this example. To test SIA against synergism, let $\pi_{ij} = \text{pr}\{\text{positive response for a test unit in cell } (i, j)\}$, and

$$\log(1 - \pi_{ij}) = \mu + \alpha_i + \tau_j + \eta_{ij} \quad (4.34)$$

where $i = 1, 2$ and $j = 1, \dots, 5$. To ensure that the parameters in (4.34) are identified, let us impose the constraints, $\alpha_1 = \tau_1 = \eta_{11} = \eta_{1j} = 0$ for all (i, j) . Let $\theta = (\mu, \alpha_2, \tau_2, \tau_3, \tau_4, \tau_5, \eta_{22}, \eta_{23}, \eta_{24}, \eta_{25})^T$, $\lambda = (\mu, \alpha_2, \tau_2, \tau_3, \tau_4, \tau_5)^T$ and $\psi = (\eta_{22}, \eta_{23}, \eta_{24}, \eta_{25})^T$. It is easily seen that the number of distinct π_{ij} parameters, namely ten, is also the number of distinct parameters on the right-hand side of (4.34), and that (4.34) is a one-to-one transformation from π to θ , where $\pi = (\pi_{11}, \pi_{12}, \pi_{13}, \pi_{14}, \pi_{15}, \pi_{21}, \pi_{22}, \pi_{23}, \pi_{24}, \pi_{25})^T$. The SIA model is equivalent to $\psi = 0$. If $\psi \geq 0$ and $\psi \neq 0$, then there is synergism, if $\psi \leq 0$ and $\psi \neq 0$, then there is antagonism.

In this example, the question of interest is whether or not there is any synergism. Therefore, the statistical inference problem may be formulated as a test of

$$H_0 : \psi = 0 \text{ against } H_1 : \psi \geq 0. \quad (4.35)$$

This is within the general framework considered in this section; it is a special case of (4.28) with $R = [0 : I_4]$, where I_4 is the 4×4 identity matrix and 0 is the 4×6 matrix of zeroes. Now, the null and alternative hypotheses are $H_0 : R\theta = 0$ and $H_1 : R\theta \geq 0$, respectively.

For cell (i, j) , let n_{ij} and y_{ij} denote the number of units tested and the number of positive responses, respectively. Then the likelihood is the product of all the binomial probabilities, and therefore

$$\ell(\theta) = \sum_{i=1}^2 \sum_{j=1}^5 [y_{ij} \log \pi_{ij} + (n_{ij} - y_{ij}) \log(1 - \pi_{ij}) + c_{ij}]$$

where π_{ij} as a function of θ is given by (4.34) and the constant c_{ij} is the logarithm of the binomial coefficient, $n_{ij}! \{y_{ij}!(n_{ij} - y_{ij})!\}^{-1}$. By definition,

$$\begin{aligned} LRT &= 2\{\max_{\theta \in H_1} \ell(\theta) - \max_{\theta \in H_0} \ell(\theta)\} \\ &= 2\{\max_{\psi \geq 0} \ell(\theta) - \max_{\psi = 0} \ell(\theta)\}. \end{aligned} \quad (4.36)$$

In this model, the loglikelihood is concave, and by modifying the arguments in Silvapulle (1981), it can be verified that the *mle* exists at a finite point with probability one. The value of $\max\{\ell(\theta) : \psi \geq 0\}$ under H_1 can be computed using a nonlinear inequality constrained optimization program in software libraries, such as *IMSL*, *MATLAB*, or *NAG*. In this example, we did not need to use such a constrained optimization program, for reasons explained below.

Note that the equation (4.34) is a reparameterization of π in terms of θ . Therefore, to compute the global *mle*, $\hat{\theta}$, it suffices to solve (4.34) with $\pi = \hat{\pi}$ where $\hat{\pi}$, the vector of sample proportions of successes, is the global *mle* of π . The computed value of ψ turned to be nonnegative for each component. Therefore, the *mle* of θ under the constraint $\psi \geq 0$ is also the unconstrained estimator, $\hat{\theta}$. Therefore, the sample value of the *LRT* for testing $\psi = 0$ against $\psi \geq 0$ is equal to that for testing $\psi = 0$ against $\psi \neq 0$. Although the *LRT* for the two tests have the same numerical value for these data, they have different *p*-values, because their null distributions are different. The computed value of *LRT* is 60.5. A large sample approximation of the *p*-value for testing $\psi = 0$ against $\psi \geq 0$ is

$$\sup_{\lambda} \sum_{i=1}^4 w_i \{4, \mathcal{I}^{\psi\psi}(\lambda)\} \text{pr}(\chi_i^2 \geq 60.5). \quad (4.37)$$

Because this supremum is not easy to compute, it is worth applying a bounds test first. An upper bound for the foregoing large sample *p*-value is (see, (3.23))

$$0.5[\text{pr}(\chi_3^2 \geq 60.5) + \text{pr}(\chi_4^2 \geq 60.5)];$$

this is less than 0.001. Therefore, the large sample *p*-value corresponding to a sample value of 60.5 for the *LRT* is less than 0.001. Therefore, there is sufficient evidence to reject H_0 and accept H_1 ; in other words, reject the SIA model and accept that there is synergism between potassium chromate and 9-AA.

For illustrative purposes, let us apply the *LRT* for testing $\psi = 0$ against $\psi \neq 0$ as well. It was noted that the *LRT* for this is also 60.5. The large sample *p*-value for this unrestricted test is $\text{pr}(\chi_4^2 \geq 60.5)$, which is also less than 0.001. Therefore,

assuming that the nominal level is 0.05, we would reject $\psi = 0$ and accept $\psi \neq 0$. It is valid to apply this unrestricted test for testing $\psi = 0$ against $\psi \geq 0$ as well. Recall that the reason that we prefer to apply a one-sided/restricted test as opposed to a two-sided/unrestricted test for testing $\psi = 0$ against $\psi \geq 0$ is for power advantage; validity is not the issue because a 5% level test of $\psi = 0$ against $\psi \neq 0$ is also a 5% level test of $\psi = 0$ against $\psi \geq 0$. Since the p-value for $\psi = 0$ against $\psi \neq 0$ is small, we would reject $\psi = 0$ and accept $\psi \geq 0$. Therefore, in this example, the evidence against $\psi = 0$ is sufficiently strong and we need not have carried out a one-sided test to decide whether or not to reject H_0 . However, if it is known that $\psi \geq 0$ then a test of $H_0 : \psi = 0$ against $H_1 : \psi \geq 0$ would in general summarize the statistical evidence better than a test of $H_0 : \psi = 0$ against $H_2 : \psi \neq 0$.

The foregoing computations can be modified to incorporate an overdispersion parameter as well. ■

Example 4.3.2 Comparison of treatments with ordinal response data.

Example 1.2.8 of Chapter 1 (page 12) provides details of a trial to compare two treatments for ulcer; the data are given in Table 1.7. The objective of the study is to show that Treatment B is better than Treatment A. The hypotheses can be formulated in different ways. One possibility is to say that Treatment B is better than Treatment A if the local odds ratios are all greater than or equal to 1 with at least one of them being greater than 1, for interpretations and discussions about odds ratios, see Chapter 6. However, this may be an unnecessarily stringent requirement. A less stringent, and hence perhaps more desirable, formulation of the problem is

$$H_0 : \gamma_1 = \gamma_2 = \gamma_3 = 1 \quad \text{and} \quad H_1 : \gamma_1 \geq 1, \gamma_2 \geq 1, \gamma_3 \geq 1, \quad (4.38)$$

$$\text{where } \gamma_q = \left(\sum_{j=1}^q \pi_{1j} \right) \left(\sum_{j=q+1}^C \pi_{2j} \right) \left\{ \left(\sum_{j=1}^q \pi_{2j} \right) \left(\sum_{j=q+1}^C \pi_{1j} \right) \right\}^{-1} \quad (4.39)$$

and $\pi_{ij} = \text{pr}(\text{the response of a subject in row } i \text{ falls in column } j)$.

The γ_q 's are also certain odds ratios of interest as discussed in Example 1.2.8 of Chapter 1. The parameter $\boldsymbol{\pi} = (\pi_{11}, \pi_{12}, \pi_{13}, \pi_{21}, \pi_{22}, \pi_{23})^T$ defines the unrestricted model completely. The likelihood is simply the product of the two multinomial probabilities for Treatment A and Treatment B:

$$\text{Likelihood at } \boldsymbol{\pi} \propto \prod_{i=1}^{q+1} \prod_{j=1}^{C-q} \pi_{ij}^{y_{ij}}. \quad (4.40)$$

To apply the results of this section, it is important that the hypotheses involve only linear functions of the parameter that appears in the likelihood. However, in the foregoing formulation, this is not the case because the likelihood is expressed as a function of $\boldsymbol{\pi}$ (see 4.40) and the hypothesis involves nonlinear functions of $\boldsymbol{\pi}$. General methods of testing hypotheses involving nonlinear functions of parameters will be discussed later. Such methods can be applied for testing (4.38). However, in this particular example, there is a simple reparameterization such that the hypotheses involve only linear functions and hence the results of this section are applicable.

Let us parameterize the π_{ij} 's as

$$\log(\pi_{i1} + \dots + \pi_{iq}) = \lambda_q - \psi_{iq} \quad (4.41)$$

where $i = 1, 2$ and $q = 1, 2, 3$. Let $\boldsymbol{\theta} = (\lambda_1, \lambda_2, \lambda_3, \psi_{21}, \psi_{22}, \psi_{23})^T$, $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \lambda_3)^T$ and $\boldsymbol{\psi} = (\psi_{21}, \psi_{22}, \psi_{23})^T$. Then $\gamma_1 \geq 1, \gamma_2 \geq 1, \gamma_3 \geq 1$ is equivalent to $\boldsymbol{\psi} \geq 0$. Now,

Likelihood at $\boldsymbol{\theta} = L(\boldsymbol{\theta}) = \text{Const} \prod_{i=1, j=1}^{q+1} \prod_{j=1}^{C-q} \{\pi_{ij}(\boldsymbol{\theta})\}^{y_{ij}}$

where the value of $\pi_{ij}(\boldsymbol{\theta})$ in terms of $\boldsymbol{\theta}$ is given by (4.41). For the data in Table 1.7, $\hat{\gamma}_1 = (12 \times 27)/(5 \times 20) = 3.24$, $\hat{\gamma}_2 = (22 \times 19)/(13 \times 10) = 3.22$, 4 and $\hat{\gamma}_3 = (26 \times 11)/(21 \times 6) = 2.27$. Thus, each $\hat{\gamma}_i$ is greater than 1 and hence $\boldsymbol{\psi} \geq 0$. Again, as in the previous example, the unrestricted estimator of $\boldsymbol{\theta}$ satisfies the inequalities in H_1 . Consequently, the estimator of $\boldsymbol{\theta}$ under H_1 is also the unrestricted estimator of $\boldsymbol{\theta}$. Now, the LRT for testing $\boldsymbol{\psi} = 0$ against $\boldsymbol{\psi} \geq 0$ is 6.0. Under H_0 , we have that $\pi_{1j} = \pi_{2j}$ for $j = 1, \dots, 4$. Therefore, the mle of $(\pi_{11}, \pi_{12}, \pi_{13}, \pi_{24})$ under H_0 is the vector of sample proportions for column totals ($i = 1, 2$). Hence, the mle of $\boldsymbol{\pi}$ under H_0 is $\bar{\boldsymbol{\pi}} = (\bar{\pi}_1^T, \bar{\pi}_2^T)^T = (17/64, 18/64, 12/64, 17/64, 17/64, 18/64, 12/64, 17/64)^T$. Because the mle is invariant under transformations, the mle of $\boldsymbol{\lambda}$ under H_0 is obtained by solving (4.41) with $\boldsymbol{\pi} = \bar{\boldsymbol{\pi}}$ and $\boldsymbol{\psi} = 0$; this leads to $\bar{\boldsymbol{\lambda}} = \{\log(17/64), \log(35/64), \log(47/64)\}$. Now, an estimate of the large sample p-value for LRT is (see p-test)

$$\hat{p} = \sum_{i=1}^3 w_i \{3, \mathcal{T}^{\psi\psi}(\bar{\lambda})\} \text{pr}(\chi_i^2 \geq 6.0) = 0.054;$$

the formulas in (3.26) were used for computing $w_i \{3, \mathcal{T}^{\psi\psi}(\bar{\lambda})\}$. Using the same closed-form for the weights, the supremum of the tail probability over $\boldsymbol{\lambda}$ was computed as

$$\sup_{\boldsymbol{\lambda}} \sum_{i=1}^3 w_i \{3, \mathcal{T}^{\psi\psi}(\boldsymbol{\lambda})\} \text{pr}(\chi_i^2 \geq 6.0) = 0.083.$$

Therefore, a large sample approximation of the p-value is 0.083; as expected, this is larger than \hat{p} ($= 0.054$) and the difference between the two is not negligible. In this example, some of the observed cell frequencies (in Table 1.7) are not that large, and, therefore, it may be prudent to be conservative and use 0.083 ($= p_{sup}$) rather than $0.054 (= \hat{p})$ as a large sample p-value.

It would be of interest to compare these results with a test against the unrestricted alternative that Treatments A and B are different for which the null and alternative hypotheses take the form

$$H_0 : \boldsymbol{\psi} = 0 \text{ and } H_2 : \boldsymbol{\psi} \neq 0,$$

respectively. Because the unrestricted mle of $\boldsymbol{\psi}$ satisfies the constraint $\boldsymbol{\psi} \geq 0$, the LRT of $\boldsymbol{\psi} = 0$ against $\boldsymbol{\psi} \neq 0$ and that of $\boldsymbol{\psi} = 0$ against $\boldsymbol{\psi} \geq 0$ have the same numerical value, namely 6.0. However, the p-values are different because the null

distributions are different. For $H_0 : \psi = 0$ against $H_2 : \psi \neq 0$, the large sample p -value is $\text{pr}(\chi^2_3 \geqslant 6.0) = 0.112$. As expected (because the components of U are all positive), this p -value is larger than those obtained earlier for testing against $R\theta \geqslant 0$. As has been noted earlier, if the unrestricted estimate satisfies the inequalities in H_1 , then the LRT against $H_1 : \psi \geqslant 0$ and that against $H_2 : \psi \neq 0$ would have the same numerical value but the latter would have a larger p -value. This follows easily because the chi-bar-square distribution corresponding to the test against $\psi \geqslant 0$ has shorter tails than the χ^2 distribution corresponding to the test against $\psi \neq 0$. This is consistent with our intuition: if the unrestricted estimate, $\hat{\psi}$, is nonnegative, then the data are pointing us in the direction of $\psi \geqslant 0$ and the evidence against $\psi = 0$ from a test against $\psi \geqslant 0$ is likely to be stronger than that from a global test against $\psi \neq 0$. The corresponding phenomenon is well-known in the classical one-sided t -test on the mean of a normal distribution. ■

4.4 TESTS OF $h(\theta) = 0$

The results in Section 4.1 for linear constraints of θ can be extended to nonlinear constraints as well. For simplicity we shall first consider the case when the observations are *iid*. Throughout this section we shall assume that the natural parameter space Θ is open and that Condition Q is satisfied. Most of the results of this section can be extended to the setting when the observations are not *iid*. Some of these will be discussed later in this chapter.

Let $\mathbf{h}(\theta) = (h_1(\theta), \dots, h_r(\theta))^T$ be a continuously differentiable vector function. In the next subsection, we will define several tests of $\mathbf{h}(\theta) = 0$ vs $\mathbf{h}(\theta) \geqslant 0$. Then in subsection 4.4.2 we will consider a test against $\mathbf{h}_2(\theta) \geqslant 0$ where \mathbf{h}_2 is a subvector of \mathbf{h} .

4.4.1 Test of $\mathbf{h}(\theta) = 0$ Against $\mathbf{h}(\theta) \geqslant 0$

Let us consider the hypotheses

$$H_0 : \mathbf{h}(\theta) = 0 \quad \text{vs} \quad H_1 : \mathbf{h}(\theta) \geqslant 0. \quad (4.42)$$

Let $\mathbf{H}(\theta) = (\partial/\partial\theta^T)\mathbf{h}(\theta)$. Assume that

$$\mathbf{H}(\theta_0) \text{ has full row rank.} \quad (4.43)$$

As usual, let $\hat{\theta}$ denote the unconstrained mle, and $\bar{\theta}$ and $\tilde{\theta}$ denote the mle's under H_1 and H_0 , respectively. The definition of LRT does not require any new ideas: $LRT = 2\{\ell(\hat{\theta}) - \ell(\bar{\theta})\}$. The asymptotic null distribution of this turns out to be a chi-bar-square. Another convenient statistic for testing (4.42) is the distance statistic based on $\hat{\theta}$. Since $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{0, \mathcal{I}(\theta_0)^{-1}\}$, we define

$$D = d(\hat{\theta}, H_0) - d(\hat{\theta}, H_1), \quad \text{where } d(\hat{\theta}, H) = \inf_{\theta \in H} n(\hat{\theta} - \theta)^T \hat{\mathcal{I}}(\hat{\theta} - \theta)$$

where $\hat{\mathcal{I}}$ is a consistent estimator of $\mathcal{I}(\theta_0)$. It will be seen that the LRT and D are asymptotically equivalent and have the common asymptotic null distribution, $\chi^2\{(\mathbf{H}(\theta_0)\mathcal{I}(\theta_0)^{-1}\mathbf{H}(\theta_0)^T, \mathbb{R}^{+r})\}$. Even if $\mathbf{H}(\theta_0)$ has rank less than the dimension of \mathbf{h} , the asymptotic null distribution of LRT is still a chi-bar-square, but it has a different form.

Much of the existing literature concerning tests of $\mathbf{h}(\theta) = 0$ usually assumes that $\mathbf{H}(\theta_0)$ is of full row-rank. It is important to note that this may fail to be satisfied even in what appears to be well-behaved models.

An example where $\mathbf{H}(\theta_0)$ is not of full rank:

Let $\mathbf{A}(\theta)$ be the 2×2 matrix $(\theta_1, \theta_3 \mid \theta_3, \theta_2)$,

$$H_0 : \mathbf{A}(\theta) = 0 \quad \text{and} \quad H_1 : \mathbf{A}(\theta) \text{ is positive semi-definite.}$$

This may be restated in the equivalent form $H_0 : \mathbf{h}(\theta) = 0$ and $H_1 : \mathbf{h}(\theta) \geqslant 0$, where $\mathbf{h}(\theta) = (-\theta_1, \theta_1\theta_2 - \theta_3^2)^T$. Now, $(\partial/\partial\theta^T)\mathbf{h}(\theta) = [1, 0, 0|\theta_2, \theta_1, -2\theta_3]$; the last row becomes a zero vector at the null value, $\theta = 0$, and hence $(\partial/\partial\theta^T)\mathbf{h}(\theta)$ is not of full rank. Therefore, condition (4.43) is not satisfied.

To provide a more concrete example, we consider the regression model

$$y = g(\mathbf{x}; \theta) + \epsilon, \quad \text{where} \quad g(\mathbf{x}; \theta) = \theta_1 x_1^2 + 2\theta_3 x_1 x_2 + \theta_2 x_2^2 + \theta_4 x_1 + \theta_5 x_2 + \theta_6.$$

Let the hypotheses be

$$H_0 : g(\mathbf{x}; \theta) = \text{linear in } \mathbf{x} \quad \text{and} \quad H_1 : g(\mathbf{x}; \theta) \text{ is concave in } \mathbf{x}.$$

First note that the Hessian matrix of $g(\mathbf{x}; \theta)$ with respect to \mathbf{x} is $2[\theta_1, \theta_3|\theta_3, \theta_2]$. Now, concavity of $g(\mathbf{x}; \theta)$ is equivalent to the Hessian matrix, $(1/2)\mathbf{A}(\theta)$, of $g(\mathbf{x}; \theta)$ being negative semi-definite. Therefore, the hypotheses take the form $H_0 : \mathbf{A}(\theta) = 0$ and $H_1 : \mathbf{A}(\theta)$ is positive semi-definite.

Another example where this arises is when \mathbf{A} is the covariance matrix of a collection of random effects; see Example 1.2.10 on page 14. ■
We can also define tests that resemble the W and S_G introduced earlier for linear constraints of θ . All of these test statistics differ only by a $o_p(1)$ term under the null hypothesis. Thus, they have the same asymptotic null distribution. Further, they also have the same asymptotic local power against Pitman-type local alternatives; this would require some conditions, for example, contiguity. Other types of asymptotic efficiencies of these tests have not been studied in any detail. As noted earlier, not much is known about their relative advantages and disadvantages. For completeness, they are introduced here.

For simplicity, we shall adopt the following notation. A function, say $\bar{h}(\theta)$, evaluated at $\bar{\theta}$ and $\bar{\theta}$ will be denoted by \tilde{h} and \bar{h} , respectively; thus $\tilde{h} = h(\bar{\theta})$ and $\bar{h} = h(\bar{\theta})$. Let \tilde{H} and $\tilde{\mathcal{I}}$ denote consistent estimators of $\mathbf{H}(\theta)$ and $\mathcal{I}(\theta)$ respectively, under the null hypothesis; therefore, they could be obtained by substituting $\bar{\theta}$, $\bar{\theta}$ or $\bar{\theta}$ for θ .

(1) Wald-type test

This is very similar to the Wald-type test introduced earlier for linear constraints. Recall that the Wald-type statistic for testing $\mathbf{h}(\theta) = 0$ against $\mathbf{h}(\theta) \neq 0$ is

$n\hat{h}^T(\check{H}\check{T}^{-1}\check{H}^T)^{-1}\hat{h}$. This suggests the statistic

$$W = n\hat{h}^T(\check{H}\check{T}^{-1}\check{H}^T)^{-1}\hat{h}.$$

(2) Hausman-Wald type test

This simply evaluates the distance between the *mle*'s under H_0 and H_1 .

$$W_H = n(\tilde{\theta} - \bar{\theta})^T \check{I}(\tilde{\theta} - \bar{\theta}).$$

(3) Global score test:

Motivated by arguments given for the linear constraints setting, we define

$$S_G = n^{-1}U^T(\check{H}\check{T}^{-1}\check{H}^T)^{-1}U, \text{ where } U = \check{H}\check{T}^{-1}(\check{S} - \bar{S}).$$

Another, possible global score statistic is (see Robertson et al. (1988)),

$$n^{-1}(\check{S} - \bar{S})^T \check{I}^{-1}(\check{S} - \bar{S}).$$

(4) Kühn-Tücker multiplier test

First, it would be instructive to provide an insight into the so called *Lagrange multiplier test*. Let us consider the optimization problem

$$\text{Problem P1: } \max \ell(\theta) \quad \text{subject to } h(\theta) = b$$

where b is a given set of constraints. Let the Lagrangian function be defined as $\{\ell(\theta) + \lambda^T(h(\theta) - b)\}$, where $\lambda = (\lambda_1, \dots, \lambda_r)^T$. Let the optimal solution be $(\tilde{\theta}_b : \tilde{\lambda}_b)$; $\tilde{\lambda}_b$ is the Lagrange multiplier. Now, we have (for example, see Intriligator (1971), page 36)

$$\tilde{\lambda}_b = -(\partial/\partial b)\ell(\tilde{\theta}_b).$$

Thus, $\tilde{\lambda}_b$ is a measure of the sensitivity of the constrained maximum of $\ell(\theta)$ over $\{\theta : h(\theta) = b\}$ to the value b . If the unconstrained maximum of $\ell(\theta)$ is achieved over $\{\theta : h(\theta) = b\}$ then $\tilde{\lambda}_b = 0$.

For brevity of notation, let us assume that $b = 0$ and suppress the suffix b in λ_b . Aitchison and Silvey (1958) showed that, under $h(\theta) = 0$, $\tilde{\lambda}/\sqrt{n}$ is asymptotically $N\{0, V(\theta_0)\}$ where $V(\theta) = [\check{H}\check{T}^{-1}\check{H}^T]_{\theta=0}^{-1}$; it is asymptotically normal with nonzero mean under local alternatives. Therefore, it is reasonable to expect that the Lagrange multipliers have the potential to be building blocks for constructing a test of $h(\theta) = 0$ against $h(\theta) \neq 0$. The so called *Lagrange multiplier* statistic for testing $h(\theta) = 0$ vs $h(\theta) \neq 0$ is $n^{-1}\tilde{\lambda}^T V(\theta)^{-1}\tilde{\lambda}$. Since $\tilde{\lambda}$ satisfies $[(\partial/\partial\theta)\ell(\theta) + H(\theta)^T\tilde{\lambda}]_{\theta=\tilde{\theta}} = 0$ we have

$$n^{-1}\tilde{\lambda}^T V(\theta)^{-1}\tilde{\lambda} = n^{-1}[S(\theta)^T \mathcal{I}(\theta)^{-1} S(\theta)]_{\theta=\tilde{\theta}};$$

the expression on the right-hand side is Rao's score statistic.

The same idea can be extended to construct tests when there are inequality constraints. For example, consider the optimization problem

$$\text{Problem P2: } \max \ell(\theta) \quad \text{subject to } h(\theta) \geq b$$

Let the Lagrangian function be defined as above. The first-order Kühn-Tücker conditions (also known as Karush-Kühn-Tücker conditions) are

$$\nabla \ell(\tilde{\theta}) + \nabla h(\tilde{\theta})^T \tilde{\lambda} = 0, h(\tilde{\theta}) \geq 0, \tilde{\lambda} \geq 0, \lambda_j h_j(\tilde{\theta}) = 0 \text{ for } j = 1, \dots, r.$$

Let the solution be $(\tilde{\theta}_b, \tilde{\lambda}_b)$; $\tilde{\lambda}_b$ are known as Kühn-Tücker multipliers. Again $\tilde{\lambda}_b$ can also be interpreted as a measure of sensitivity of the constrained maximum of $\ell(\theta)$ to b , in much the same way as for the foregoing equality constrained optimization problem (for details about this interpretation and its relevance, see Intriligator (1971) page 60).

Let $\bar{\lambda}$ and $\tilde{\lambda}$ denote the Kühn-Tücker multipliers associated with the maximization of $\ell(\theta)$ subject to $h(\theta) = 0$ and $h(\theta) \geq 0$, respectively. Now, the KT-statistic is defined as

$$KT = n(\tilde{\lambda} - \bar{\lambda})^T (\check{H}\check{T}^{-1}\check{H})(\tilde{\lambda} - \bar{\lambda}).$$

Now, the next result says that the foregoing tests are asymptotically equivalent.

Proposition 4.4.1 *Under the null hypothesis, $H_0 : h(\theta) = 0$, we have*

$$LRT = W + o_p(1) = W_H + o_p(1) = S_G + o_p(1) = D + o_p(1) = KT + o_p(1).$$

Proof: The proofs of these, except that for D , are given in Gouriéroux and Monfort (1995, Chapter, 21); the equivalence of D to the *LRT* follows from the quadratic approximation (C) in Proposition 4.2.2. ■

An unattractive feature of W , W_H , S_G , and KT is that they depend on $H(\theta)$ explicitly; further, some of the technical arguments also assume that H is of full rank. By contrast, *LRT* and distance-based tests do not depend on the choice of the function h to define a given set of constraints.

In some cases, it may be difficult to compute the *mle* and/or the information matrix. In such cases, it may be of interest to construct a test based on an estimator that may be easier to obtain than the *mle*. Such a method is given in the next subsection.

4.4.2 Test of $h(\theta) = 0$ Against $h_2(\theta) \geq 0$

Let us consider the setting in Section 4.2. In particular, the parameter space Θ is open and Condition Q is satisfied. As in the previous subsection, let h be an $r \times 1$ vector function of θ ; it is assumed that h is continuously differentiable. Let $H(\theta)_{r \times p} = (\partial/\partial\theta^T)h(\theta)$ and assume that its rank is r for values of θ in the null parameter space.

Let h be partitioned as $(h_1^T, h_2^T)^T$. In this section, we shall introduce a test of $H(\theta) : h(\theta) = 0$ against $H_1 : h_2(\theta) \geq 0$

based on the basic idea that underlies the Wald test.

Let $\hat{\theta}$ denote an estimator of θ such that

$$n^{1/2}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{0, W(\theta_0)\}$$

for some $W(\theta)$. For example, $\hat{\theta}$ can be the unrestricted maximum likelihood estimator. It follows that

$$n^{1/2}\{h(\hat{\theta}) - h(\theta_0)\} \xrightarrow{d} N\{0, V(\theta_0)\}, \quad (4.45)$$

where $V(\theta) = H(\theta)W(\theta)H(\theta)^T$. Now, we can use some of the arguments used in subsection 4.3.3 for defining Wald-type and distance statistics. Let \hat{V} be a consistent estimator of $V(\theta)$ under the null hypothesis, for example, $\hat{V} = H(\hat{\theta})W(\hat{\theta})H(\hat{\theta})^T$.

A suitable statistic for testing H_0 against H_1 is

$$D = n[h(\hat{\theta})^T \hat{V}^{-1} h(\hat{\theta}) - \min_{\alpha} \{(h(\hat{\theta}) - \alpha)^T \hat{V}^{-1} (h(\hat{\theta}) - \alpha) : \alpha_2 \geq 0\}] \quad (4.46)$$

where α in \mathbb{R}^r is partitioned as $(\alpha_1 : \alpha_2)$ to conform with $\mathbf{h} = (\mathbf{h}_1 : \mathbf{h}_2)$; see Kodde and Palm (1986). Note that $D = d\{h(\hat{\theta}), H_0\} - d\{h(\hat{\theta}), H_1\}$, where $d\{h(\hat{\theta}), H\}$ is a measure of the distance between $h(\hat{\theta})$ and H . The test rejects H_0 if D is large.

Suppose that H_0 is satisfied. It follows from $\{V(\theta_0) - \hat{V}\} = o_p(1)$ that, under H_0 , the value of D would change by $o_p(1)$ if \hat{V} is replaced by $V(\theta_0)$ (see Lemma 4.10.2 in the Appendix to this chapter). Therefore, the limiting distribution of D in (4.46) is equal to the distribution of

$$T = \mathbf{X}^T V(\theta_0)^{-1} \mathbf{X} - \min_{\alpha} \{(\mathbf{X} - \alpha)^T V(\theta_0)^{-1} (\mathbf{X} - \alpha) : \alpha_2 \geq 0\}$$

where $\mathbf{X} \sim N\{0, V(\theta_0)\}$. Consequently, the asymptotic null distribution of D is given by

$$\text{pr}_{\theta_0}\{D \geq c \mid H_0\} \rightarrow \sum_{i=0}^q w_i \{q, V_{22}(\theta_0)\} \text{pr}(\chi^2_{r-q+i} \geq c), \quad (4.47)$$

where $V_{(22)}$ is the diagonal block of V corresponding to \mathbf{h}_2 . Note that the large sample null distribution of D depends on the assumed true value θ_0 of the parameter θ in the null parameter space. Approaches to dealing with such nuisance parameter problems were discussed in an earlier section.

If $\hat{\theta}$ is the mle, then it follows from the quadratic approximations of the loglikelihood, that LRT and the other tests introduced in this subsection are asymptotically equivalent and hence their asymptotic null distribution is given by (4.47) with $\mathcal{I}(\theta)$ substituted for $W(\theta)$ in the formula for $V(\theta)$.

4.5 AN OVERVIEW OF SCORE TESTS WITH NO INEQUALITY CONSTRAINTS

Although formal proofs are not given, the essential arguments are indicated. Details of the arguments may be found in Basawa (1985), Cox and Hinkley (1974), Sen and Singer (1993), and van der Vaart (2000). Hall and Mathiason (1990) provide a thorough account of many of these results at an intermediate level. These will be used in Section 4.6 to introduce local score and score-type tests against inequality constraints.

Let $f(y; \theta)$ denote the density of Y where $\theta \in \Theta \subset \mathbb{R}^p$, and Y_1, \dots, Y_n denote n iid observations on the random variable/vector Y . It is assumed that Θ is open in \mathbb{R}^p .

First, we shall consider the case of testing $\theta = \theta_0$ against $\theta \neq \theta_0$ and then consider the more general case when there are nuisance parameters. In this section we will be studying the distribution of statistics when H_0 is not true. Therefore, the true value may not be θ_0 .

4.5.1 Test of $H_0 : \theta = \theta_0$ Against $H_2 : \theta \neq \theta_0$

Let

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^n \log f(Y_i; \theta), & S(\theta) &= (\partial/\partial\theta)\ell(\theta) \\ \text{and } \mathcal{I}(\theta) &= E_{\theta}\{(\partial/\partial\theta)\log f(Y; \theta)(\partial/\partial\theta^T)\log f(Y; \theta)\} \end{aligned} \quad (4.48)$$

denote the loglikelihood, score function, and the information matrix, respectively. Assume that Condition Q in subsection 4.2 is satisfied. Let the null and alternative hypotheses be

$$H_0 : \theta = \theta_0 \quad \text{and} \quad H_2 : \theta \neq \theta_0, \quad (4.49)$$

respectively. Then, because $\log x \leq (x - 1)$ (or by Jensen's inequality), it follows that

$$E_{\theta} \log\{f(Y; \theta^*)/f(Y; \theta)\} \leq E_{\theta}\{f(Y; \theta^*)/f(Y; \theta)\} - 1 = 0,$$

and hence

$$\text{plim}_{\theta} n^{-1}\ell(\theta^*) = E_{\theta} \log f(Y; \theta^*) \leq E_{\theta} \log f(Y; \theta) = \text{plim}_{\theta} n^{-1}\ell(\theta)$$

where plim_{θ} is the probability limit when θ is the true value. Therefore, $\text{plim}_{\theta} n^{-1}\ell(\theta^*)$ achieves its global maximum at $\theta^* = \theta$. Further, $\theta^* = \theta$ is a stationary point of $\text{plim}_{\theta} n^{-1}\ell(\theta^*)$. This suggests that it may be possible to use the slope of loglikelihood at θ_0 to test $H_0 : \theta = \theta_0$. To this end, note that

$$n^{-1}\mathbf{S}(\theta_0) \xrightarrow{P} \mathbf{0} \text{ under } H_0, \text{ and } n^{-1}\mathbf{S}(\theta_0) \xrightarrow{P} \boldsymbol{\alpha}(\theta) \text{ under } H_2 \quad (4.50)$$

where $\boldsymbol{\alpha}(\theta) = E_{\theta}\{(\partial/\partial\theta)\log f(Y; \theta_0)\}$; assume that $\boldsymbol{\alpha}(\theta) \neq \mathbf{0}$ for $\theta \neq \theta_0$. We also have $n^{-1/2}\mathbf{S}(\theta_0) \xrightarrow{d} N\{0, \mathcal{I}(\theta_0)\}$ under H_0 . This suggests the score statistic

$$S_L = n^{-1}\mathbf{S}(\theta_0)^T \mathcal{I}(\theta_0)^{-1} \mathbf{S}(\theta_0) \quad (4.51)$$

for testing $H_0 : \theta = \theta_0$ vs $H_2 : \theta \neq \theta_0$. Because $S_L \xrightarrow{d} \chi_p^2$ under H_0 and $S_L \xrightarrow{P} \infty$ under H_2 , the score test rejects H_0 if S_L is large.

In this section, a brief summary of the relevant asymptotic results on likelihood ratio, Wald- and score-type tests is provided when there are no inequality constraints.

The Wald statistic for $H_0 : \theta = \theta_0$ against $H_2 : \theta \neq \theta_0$ is defined as

$$T_W = n(\hat{\theta} - \theta_0)^T \hat{\mathcal{I}}(\hat{\theta} - \theta_0), \quad (4.52)$$

where $\hat{\mathcal{I}}$ is a consistent estimator of $\mathcal{I}(\theta_0)$. In the usual definition of a Wald statistic, the normalizing matrix $\hat{\mathcal{I}}$ is chosen so that it would be consistent irrespective of whether the null hypothesis is true or not; however, for the validity of the test, it suffices $\hat{\mathcal{I}}$ to be consistent under the null hypothesis only. Clearly, $T_W \xrightarrow{d} \chi_p^2$ under $H_0 : \theta = \theta_0$ and $T_W \xrightarrow{p} \infty$ under $H_2 : \theta \neq \theta_0$, and therefore the Wald test rejects H_0 if T_W is large.

Under $H_2 : \theta \neq \theta_0$, the probability of rejecting H_0 tends to 1 as $n \rightarrow \infty$ for the foregoing score and Wald tests at a fixed level. Similarly, at any fixed level, the power of many other tests at $\theta \neq \theta_0$ also tend to 1. Therefore, to compare the performance of different tests, the limiting power of a test at a fixed level and at a fixed point in the alternative parameter space is not particularly helpful. The procedure adopted here to compare the asymptotic power of tests is the limiting power at *local alternatives* of the form $H_n : \theta = (\theta_0 + n^{-1/2}\delta)$ where δ is fixed and the level of the test is also fixed. The sequence of hypotheses $H_n : \theta = \theta_0 + n^{-1/2}\delta$ are also known as Pitman-type local alternatives; for a thorough discussion of these see Serfling (1980) or Hájek et al. (1999). Here, we shall indicate the main arguments and results.

Let $\theta_n = (\theta_0 + n^{-1/2}\delta)$ where δ is fixed. Let a sequence of local hypotheses be defined by $H_n : \theta = \theta_n$. Then, by a one-term Taylor expansion of $S(\theta_0)$ about θ_n , we have that

$$n^{-1/2} S(\theta_0) = n^{-1/2} S(\theta_n) + \mathcal{I}(\theta_n)\delta + o_p(1) \quad \text{under } H_n : \theta = \theta_n \quad (4.53)$$

Since $\mathcal{I}(\theta_n)\delta \xrightarrow{p} \mathcal{I}(\theta_0)\delta$ and $S(\theta_n)$ is a sum of *iid* random variables, we have that

$$n^{-1/2} S(\theta_0) \xrightarrow{d} N\{\mathcal{I}(\theta_0)\delta, \mathcal{I}(\theta_0)\} \quad \text{and} \quad S_L \xrightarrow{d} \chi_p^2(\Delta) \quad \text{under } H_n, \quad (4.54)$$

where $\Delta = \delta^T \mathcal{I}(\theta_0)\delta$ is the noncentrality parameter. Further, by modifying the arguments to accommodate that θ_n is the true value that may not be a fixed point, it can be shown that

$$n^{1/2}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{\delta, \mathcal{I}(\theta_0)^{-1}\} \quad \text{and} \quad T_W \xrightarrow{d} \chi_p^2(\Delta) \quad \text{under } H_n \quad (4.55)$$

where $\Delta = \delta^T \mathcal{I}(\theta_0)\delta$. Therefore, by (4.54) and (4.55), the score and Wald tests have the same asymptotic distribution under $H_n : \theta = \theta_n$ and hence the same limiting power against the sequences of local alternatives, H_n .

Now, assume that θ_0 is the *true value*. Then expansion of $n^{-1/2} S(\hat{\theta})$ about θ_0 leads to

$$n^{-1/2} \mathcal{I}(\theta_0)^{-1} S(\theta_0) = n^{1/2}(\hat{\theta} - \theta_0) + o_p(1). \quad (4.56)$$

Consequently, we can interchange $n^{-1/2} \mathcal{I}(\theta_0)^{-1} S(\theta_0)$ and $n^{1/2}(\hat{\theta} - \theta_0)$ in many cases without affecting the first-order asymptotic properties of the resulting test. Thus, we have that $T_W = S_L + o_p(1)$. An expansion of $\ell(\theta_0)$ about $\hat{\theta}$ leads to (see (4.4))

$$\ell(\theta_0) = \ell(\hat{\theta}) - (1/2)n(\theta_0 - \hat{\theta})^T \mathcal{I}(\hat{\theta})(\theta_0 - \hat{\theta}) + o_p(1). \quad (4.57)$$

It follows from (4.51), (4.52), (4.56), and (4.57) that

$$LRT = T_W + o_p(1) = S_L + o_p(1) \quad \text{under } H_0. \quad (4.58)$$

Consequently, LRT , T_W , and S_L , have the same asymptotic null distribution. It can also be shown that (for example, see Basawa (1985))

$$LRT = T_W + o_p(1) = S_L + o_p(1) \quad \text{under } H_n. \quad (4.59)$$

This follows easily if H_n is contiguous to H_0 ; otherwise quadratic approximations can be used. Consequently, we have that

$$LRT, T_W, S_L \xrightarrow{d} \chi^2(\Delta) \quad \text{under } H_n$$

where $\Delta = \delta^T \mathcal{I}_{\theta_0} \delta$, and hence the three tests have the same local power. Usually, once we have shown that the difference between two test statistics is $o_p(1)$ at θ_0 in the null parameter space, then the same would also hold under $H_n : \theta = \theta_0 + n^{-1/2}\delta$ as well. The proof of this may or may not be simple. For the *iid* setting, we can adopt an approach based on a quadratic approximation of $\ell(\theta)$ in $n^{-1/2}$ -neighborhoods. A different approach would be to use methods based on contiguity of hypotheses.

Essentially the same types of arguments are used when the null hypothesis is not simple. These are discussed in the next subsection.

4.5.2 Test of $H_0 : \psi = \psi_0$ Against $H_2 : \psi \neq \psi_0$

The discussions in the previous subsection are particularly useful because the simplicity of the null hypothesis with no nuisance parameters enabled us to bring out some important features without delicate technical details. In many practical situations, tests of hypotheses usually involve only some components of θ . Let θ be partitioned as $(\lambda : \psi)$; recall that when θ is partitioned as $(\lambda^T, \psi^T)^T$ it would be denoted by $(\lambda : \psi)$. Let the null and alternative hypotheses be

$$H_0 : \psi = \psi_0 \quad \text{and} \quad H_2 : \psi \neq \psi_0. \quad (4.60)$$

For this testing problem, LRT , Wald, and score tests are known to be asymptotically equivalent.

Let λ_0 be the true value of λ and $\theta_0 = (\lambda_0 : \psi_0)$ be the true value of θ when $H_0 : \psi = \psi_0$ holds. Different extensions of the score statistic in the previous section are available for this testing problem. Let the score function $S(\theta)$ in (4.48) be partitioned into $(S_\lambda : S_\psi)$ to conform with $(\lambda : \psi)$. Let $\{\mathcal{I}_{\lambda\lambda}, \mathcal{I}_{\lambda\psi}, \mathcal{I}_{\psi\lambda}, \mathcal{I}_{\psi\psi}\}$ be obtained by partitioning the information matrix \mathcal{I} to conform with $(\lambda : \psi)$; similarly, let $\{\mathcal{I}_{\lambda\lambda}^\psi, \mathcal{I}_{\lambda\psi}^\psi \mid \mathcal{I}_{\psi\lambda}^\psi, \mathcal{I}_{\psi\psi}^\psi\}$ be obtained by partitioning \mathcal{I}^{-1} . Let

$$\mathcal{I}_{\psi\psi,\lambda} = (\mathcal{I}_{\psi\psi} - \mathcal{I}_{\psi\lambda} \mathcal{I}_{\lambda\lambda}^{-1} \mathcal{I}_{\lambda\psi}); \quad \text{then } \mathcal{I}^{\psi\psi} = \mathcal{I}_{\psi\psi,\lambda}^{-1}.$$

Let $\bar{\lambda}$ denote the *mle* of λ under H_0 , defined by $S_\lambda(\bar{\theta}) = \mathbf{0}$ where $\bar{\theta} = (\bar{\lambda} : \psi_0)$; thus $\bar{\theta}$ is the *mle* of θ under H_0 . Let a sequence of local hypotheses be defined by

$$H_n : \psi = \psi_0 + n^{-1/2}\delta,$$

where $\delta \in \mathbb{R}^k$ is fixed. By employing Taylor series expansions similar to those leading to (4.53), we have

$$n^{-1/2} S_\psi(\bar{\theta}) \xrightarrow{d} N\{\mathcal{I}_{\psi\psi,\lambda}(\theta_0)\delta, \mathcal{I}_{\psi\psi,\lambda}(\theta_0)\} \quad \text{under } H_n. \quad (4.61)$$

Therefore, a score statistic for testing $H_0 : \psi = \psi_0$ vs $H_2 : \psi \neq \psi_0$ is

$$S_L = n^{-1}[S_\psi^T \mathcal{I}_{\psi\psi,\lambda}^{-1} S_\psi]_{\theta=\bar{\theta}} \quad (4.62)$$

where the suffix $\theta = \bar{\theta}$ for [.] indicates that [.] is evaluated at $\theta = \bar{\theta}$ (see Rao (1973), p 418). It follows from (4.61) that

$$S_L \xrightarrow{d} \chi_k^2(\Delta) \text{ under } H_n, \quad (4.63)$$

where $\Delta = \delta^T \mathcal{I}_{\psi\psi,\lambda}(\theta_0)\delta$. Because $\mathcal{I}(\theta_0)$ is positive definite, it follows that $\mathcal{I}(\theta_0)^{-1}$ and $\mathcal{I}_{\psi\psi,\lambda}(\theta_0)$ are positive definite, and hence $\Delta > 0$ for $\delta \neq 0$. Therefore, a test of $H_0 : \psi = \psi_0$ against $H_2 : \psi \neq \psi_0$ based on S_L in (4.62) rejects H_0 when it is large. This is the simple form of the score test of $\psi = \psi_0$ against $\psi \neq \psi_0$.

Because $n^{1/2}(\hat{\psi} - \psi_0) \xrightarrow{d} N\{0, \mathcal{I}_{\psi\psi,\lambda}^{-1}(\theta_0)\}$, the Wald statistic for testing $H_0 : \psi = \psi_0$ vs $H_2 : \psi \neq \psi_0$ is defined as

$$T_W = n(\hat{\psi} - \psi_0)^T \mathcal{I}_{\psi\psi,\lambda}(\hat{\theta})(\hat{\psi} - \psi_0). \quad (4.64)$$

Since $n^{1/2}(\hat{\psi} - \psi_0) \xrightarrow{d} N\{\delta, \mathcal{I}_{\psi\psi,\lambda}^{-1}(\theta_0)\}$ under H_n , it follows that

$$T_W \xrightarrow{d} \chi_k^2(\Delta) \text{ under } H_n.$$

Therefore, T_W in (4.64) and S_L in (4.62) have the same asymptotic power against H_n . In fact, by arguments essentially similar to those in the previous subsection, we have

$$LRT = S_L + o_p(1) = T_W + o_p(1) \text{ under } H_0 \text{ and under } H_n.$$

Thus, the three tests have the same asymptotic distribution under the null hypothesis, and also under the sequence of local alternatives, H_n ; consequently, they have the same local power as well.

4.5.3 Tests Based on Estimating Equations

As in the previous subsection, let θ be partitioned as $(\lambda : \psi)$, and let the null and alternative hypotheses be the same as in the previous subsection, namely

$$H_0 : \psi = \psi_0 \quad \text{and} \quad H_2 : \psi \neq \psi_0,$$

respectively. Instead of defining $S(\theta) = (\partial/\partial\theta)\ell(\theta)$, where $\ell(\theta)$ is the loglikelihood, a much larger class of tests can be developed by allowing $S(\theta)$ to be a suitable estimating function. For example, for large sample inference on β in the linear regression model, $y = x^T \beta + \epsilon$, there is no need to assume that the errors are normal; the estimator of β defined as the solution of the normal equations can be used for inference on β without assuming normal error distribution. Other important special cases of the estimating equation approach includes (i) the class of quasi-likelihood equations, which are particularly useful in generalized linear models (see McCullagh and Nelder (1998, Chapter 9), and (ii) the Generalized Estimating Equation (GEE) approach, which has found use in many areas of applications. The essential results are mentioned below; they are based on Basawa (1985); for a brief summary of these see Silvapulle and Silvapulle (1995) and for an in-depth and detailed coverage of quasi-likelihood see Heyde (1997).

Let $S(\theta) = \mathbf{0}$ be a $p \times 1$ estimating equation for θ ; $S(\theta)$ is called an *estimating function*. First, we shall state a set of conditions that are used in the literature to ensure that the choice of the estimating function is appropriate. However, later we shall relax this and state a simpler set of high-level sufficient conditions.

Assume that the following condition is satisfied

Condition A:

*There exist nonsingular matrices $G(\theta)$ and $V(\theta)$ such that, for $a > 0$, as $n \rightarrow \infty$,
and the convergences in (4.65) and (4.66) are uniform in θ_0 .*

$$n^{-1/2} S(\theta_0) \xrightarrow{d} N\{0, V(\theta_0)\}, \quad (4.65)$$

$$\sup_{\|\boldsymbol{u}\| \leq a} \|n^{-1/2}\{S(\theta_0 + n^{-1/2}\boldsymbol{u}) - S(\theta_0)\} + G(\theta_0)\boldsymbol{u}\| = o_p(1) \quad (4.66)$$

and the convergences in (4.65) and (4.66) are uniform in θ_0 . ■

As in the previous subsection, let $\theta_0 = (\lambda_0 : \psi_0)$ denote the true value of θ when $H_0 : \psi = \psi_0$ holds. Let $\hat{\theta}$ be defined by $S(\hat{\theta}) = \mathbf{0}$. Then

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{0, A(\theta_0)\}, \quad \text{where } A = G^{-1}V G^{-T} \text{ and } G^{-T} = (G^{-1})^T.$$

We can use $\hat{\theta}$ or $S(\hat{\theta})$ for testing hypotheses about θ . Let $(S_\lambda : S_\psi)$ be the partition of S to conform with $\theta = (\lambda : \psi)$. Similarly, let the partition of G and V be $\{G_{\lambda\lambda}, G_{\lambda\psi} \mid G_{\psi\lambda}, G_{\psi\psi}\}$ and $\{V_{\lambda\lambda}, V_{\lambda\psi} \mid V_{\psi\lambda}, V_{\psi\psi}\}$, respectively. Define

$$C(\theta) = [(V_{\psi\psi} - G_{\psi\lambda} G_{\lambda\lambda}^{-1} V_{\lambda\psi}) - (V_{\lambda\psi}^T - G_{\psi\lambda} G_{\lambda\lambda}^{-1} V_{\lambda\lambda})] G_{\lambda\lambda}^{-T} G_{\psi\lambda}^T \theta. \quad (4.67)$$

$$\text{and} \quad Z(\theta) = n^{-1/2}[S_\psi - G_{\psi\lambda} G_{\lambda\lambda}^{-1} S_\lambda]\theta. \quad (4.68)$$

The crucial result that enables us to construct local score-type tests is the following:

Proposition 4.5.1 *Let $\bar{\lambda}$ be the estimator of λ under H_0 defined by $S_\lambda(\bar{\lambda} : \psi_0) = \mathbf{0}$. Further, let $\bar{\theta} = (\bar{\lambda} : \psi_0)$ and let the sequence of local hypotheses be defined by $H_n : \psi = \psi_0 + n^{-1/2}\delta$ where δ is fixed. Then,*

- I. $Z(\bar{\theta}) = Z(\theta_0) + o_p(1)$ under H_0 .

2. $\mathbf{Z}(\boldsymbol{\theta}_0)$ and $\mathbf{Z}(\bar{\boldsymbol{\theta}}) \xrightarrow{d} N\{\mathbf{G}_{\psi\psi,\lambda}(\boldsymbol{\theta}_0)\boldsymbol{\delta}, \mathbf{C}(\boldsymbol{\theta}_0)\}$ under H_n .
3. For testing $\psi = \psi_0$ against $\psi \neq \psi_0$ a (local) score-type statistic is

$$\mathcal{S}_L = [\mathbf{Z}^T \mathbf{C}^{-1} \mathbf{Z}]_{\bar{\boldsymbol{\theta}}} \quad (4.69)$$

The asymptotic distribution of \mathcal{S}_L in (4.69), under H_n , is $\chi_k^2(\Delta)$, where $\Delta = \boldsymbol{\delta}^T [\mathbf{G}_{\psi\psi,\lambda}^T \mathbf{C}^{-1} \mathbf{G}_{\psi\psi,\lambda}]_{\boldsymbol{\theta}_0} \boldsymbol{\delta}$.

Proof: Let $\mathbf{B}(\boldsymbol{\theta}) = [-\mathbf{G}_{\psi\lambda} \mathbf{G}_{\lambda\lambda}^{-1}, \mathbf{I}_\theta]$. Then $\mathbf{Z}(\boldsymbol{\theta}) = n^{-1/2} \mathbf{B}(\boldsymbol{\theta}) \mathbf{S}(\boldsymbol{\theta})$ is a linear function of the statistic $\mathbf{S}(\boldsymbol{\theta})$. It follows from (4.65) that $\mathbf{Z}(\boldsymbol{\theta}_0) \xrightarrow{d} N\{0, \mathbf{C}(\boldsymbol{\theta}_0)\}$ under H_0 . One term Taylor expansion of $\mathbf{Z}(\bar{\boldsymbol{\theta}})$ about $\boldsymbol{\theta}_0$ leads to

$$\mathbf{Z}(\bar{\boldsymbol{\theta}}) - \mathbf{Z}(\boldsymbol{\theta}_0) \xrightarrow{p} \mathbf{0} \quad \text{under } H_0. \quad \blacksquare$$

It follows from (4.65) and (4.66) that $n^{-1/2} \mathcal{S}(\boldsymbol{\theta}_0) \xrightarrow{d} N\{\mathbf{G}(\boldsymbol{\theta}_0)\Delta, \mathbf{V}(\boldsymbol{\theta}_0)\}$ under $H_n : \psi = \psi_0 + n^{-1/2} \boldsymbol{\delta}$, where $\Delta = (0 : \boldsymbol{\delta})$ and therefore,

$$\mathbf{Z}(\bar{\boldsymbol{\theta}}) \xrightarrow{d} N\{\mathbf{G}_{\psi\psi,\lambda}(\boldsymbol{\theta}_0)\boldsymbol{\delta}, \mathbf{C}(\boldsymbol{\theta}_0)\} \text{ under } H_n : \psi = \psi_0 + n^{-1/2} \boldsymbol{\delta}. \quad (4.70)$$

Now, consider the special case when $\mathbf{S}(\boldsymbol{\theta})$ is the score function corresponding to loglikelihood. Let $\bar{\mathbf{S}} = \mathbf{S}(\bar{\boldsymbol{\theta}})$ and $\bar{\mathcal{I}} = \mathcal{I}(\bar{\boldsymbol{\theta}})$. Then $\mathbf{G}(\boldsymbol{\theta}) = \mathbf{V}(\boldsymbol{\theta}) = \mathcal{I}(\boldsymbol{\theta})$, and

$$\mathbf{Z}(\bar{\boldsymbol{\theta}}) = n^{-1/2} (\bar{\mathbf{S}}_\psi - \bar{\mathcal{I}}_{\psi,\lambda} \bar{\mathcal{I}}_{\lambda\lambda}^{-1} \bar{\mathbf{S}}_\lambda) = n^{-1/2} \bar{\mathbf{S}}_\psi.$$

In this case, the \mathcal{S}_L in (4.69) reduces to Rao's score statistic (4.62) (see Rao (1973), p 418),

$$n^{-1} \bar{\mathbf{S}}_\psi^T \bar{\mathcal{I}}_{\psi,\lambda}^{-1} \bar{\mathbf{S}}_\psi.$$

It follows from Condition A given at the beginning of this subsection that

$$n^{-1/2} \mathbf{G}^{-1}(\boldsymbol{\theta}_0) \mathbf{S}(\boldsymbol{\theta}_0) = n^{1/2} (\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + o_p(1).$$

Because of this asymptotic linear relationship between $\mathbf{S}(\boldsymbol{\theta}_0)$ and $(\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$ one would conjecture that a test based on $\mathbf{S}(\boldsymbol{\theta}_0)$ and one based on $\bar{\boldsymbol{\theta}}$ are likely to be equivalent; this in fact is the case, as will be seen. It follows from Condition A that

$$n^{1/2} (\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N\{0, \mathbf{A}(\boldsymbol{\theta}_0)\}, \quad \text{where } \mathbf{A} = [\mathbf{G}^{-1} \mathbf{V} \mathbf{G}^{-T}].$$

Therefore,

$$\sqrt{n}(\hat{\psi} - \psi_0) \xrightarrow{d} N\{0, \mathbf{A}_{\psi\psi}(\boldsymbol{\theta}_0)\},$$

and hence a Wald-type statistic for testing $H_0 : \psi = \psi_0$ against $H_1 : \psi \neq \psi_0$ is

$$T_W = n(\hat{\psi} - \psi_0)^T (\hat{\mathbf{A}}_{\psi\psi})^{-1} (\hat{\psi} - \psi_0). \quad (4.71)$$

where $\hat{\mathbf{A}}_{\psi\psi}$ is a consistent estimator of $\mathbf{A}_{\psi\psi}$, for example, $\mathbf{A}_{\psi\psi}(\hat{\boldsymbol{\theta}})$. Now we have the following:

Proposition 4.5.2 Under $H_n : \psi = \psi_0 + n^{-1/2} \boldsymbol{\delta}$, \mathcal{S}_L in (4.69) and T_W in (4.71) converge to $\chi_k^2(\Delta)$ where $\Delta = \boldsymbol{\delta}^T \mathbf{G}_{\psi\psi,\lambda}^T \mathbf{C}^{-1} \mathbf{G}_{\psi\psi,\lambda} \boldsymbol{\delta}$. ■

The definition of \mathcal{S}_L in (4.69) depends on $\bar{\boldsymbol{\theta}}$, which needs to be obtained by solving $\mathbf{S}_\lambda(\boldsymbol{\theta}) = \mathbf{0}$ under H_0 . Similarly, T_W depends on $\hat{\boldsymbol{\theta}}$, which is defined by $\mathbf{S}(\boldsymbol{\theta}) = \mathbf{0}$. In large samples, the $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$ in these definitions can be replaced by one-step estimators, which may be significantly easier to compute in complicated models. These one-step estimators are obtained by applying Newton-Raphson iteration once to a suitably chosen starting value. Usually they are asymptotically as efficient as the fully iterated solution and the difference between them is of order $o_p(n^{-1/2})$. The main steps are indicated below.

Let $t_{n,\lambda}$ be a preliminary estimator of λ such that $n^{1/2}(t_{n,\lambda} - \lambda_0) = O_p(1)$ and let $\bar{\lambda}^* = t_{n,\lambda} + n^{-1} G_{\lambda\lambda}^{-1}(\boldsymbol{\theta}^\dagger) \mathbf{S}_\lambda(\boldsymbol{\theta}^\dagger)$, where $\boldsymbol{\theta}^\dagger = (t_{n,\lambda} : \psi_0)$, and $\bar{\boldsymbol{\theta}}^* = (\bar{\lambda}^* : \psi_0)$. Then, under H_0 , $\mathbf{Z}(\bar{\boldsymbol{\theta}}^*) - \mathbf{Z}(\bar{\boldsymbol{\theta}}) = o_p(1)$. Consequently, the $\bar{\boldsymbol{\theta}}$ in (4.69) can be replaced by $\bar{\boldsymbol{\theta}}^*$ without affecting the large sample distribution of \mathcal{S}_L under H_0 . More details and references to the results of this subsection may be found in Basawa (1985).

4.6 LOCAL SCORE-TYPE TESTS OF $H_0 : \psi = 0$ AGAINST $H_1 : \psi \in \Psi$

$$H_1 : \psi \in \Psi$$

In this section, the main results mentioned in section 4.5 will be used to develop a local score-type test of $\psi = 0$ against $\psi \in \Psi$ where Ψ is a given subset of \mathbb{R}^k and $\boldsymbol{\theta} = (\lambda : \psi)$. Initially, we shall consider the case when $\Psi = C$, where C is a closed convex cone. In the last subsection, we shall consider the case when Ψ is not necessarily a cone. We shall consider the settings in subsections 4.5.1, 4.5.2, and 4.5.3 in turn. To introduce the ideas, we will first consider the special case of testing $\boldsymbol{\theta} = 0$ against $\boldsymbol{\theta} \in C$ where C is a cone, and provide sufficient motivation for the definition of a suitable statistic. This is then extended in subsection 4.6.2 to the case when nuisance parameters are present. Then in subsection 4.6.3 we extend these and define a statistic based on estimating equations. In all these, the score-type statistic will be denoted by \mathcal{S}_L , but the setting will become progressively more general. In the final subsection, it will be general enough to incorporate quasi-likelihood, partial likelihood, and estimating equations. This section is based on Silvapulle and Silvapulle (1995) with particular attention to the local score type tests introduced there.

4.6.1 Local Score Test of $H_0 : \boldsymbol{\theta} = 0$ Against $H_1 : \boldsymbol{\theta} \in C$

Let $f(y; \boldsymbol{\theta})$ be a density function and let Y_1, \dots, Y_n be independently and identically distributed with common density function $f(y; \boldsymbol{\theta})$ where $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p$. Assume that $\mathbf{0}$ is an interior point of Θ . Let \mathcal{I}_0 denote $\mathcal{I}(\mathbf{0})$. Assume that Condition Q in subsection

4.2 is satisfied. Let the null and alternative hypotheses be

$$H_0 : \theta = 0 \text{ and } H_1 : \theta \in C, \quad (4.72)$$

respectively, where C is a closed convex cone. Let $\mathbf{S}(\theta)$ denote the score function, $(\partial/\partial\theta)\ell(\theta)$, where $\ell(\theta)$ is the loglikelihood, and let

$$\mathbf{U} = n^{-1/2}\mathcal{I}_0^{-1}\mathbf{S}(0). \quad (4.73)$$

Let $\delta \in \mathbb{R}^p$ be a fixed point and $\theta_n = n^{-1/2}\delta$. It follows from (4.53) that

$$\mathbf{U} \xrightarrow{d} N(\delta, \mathcal{I}_0^{-1}) \text{ under } H_n : \theta = \theta_n. \quad (4.74)$$

This can be used as a guide to constructing a test of $\theta = 0$ against $H_1 : \theta \in C$. To this end, suppose that the true value of θ is $n^{-1/2}\delta$. Because $\theta \in C$ is equivalent to $\delta \in C$, it may be adequate to construct a test of $\delta = 0$ against $\delta \in C$ for testing $\theta = 0$ against $\theta \in C$. In view of (4.74), \mathbf{U} can be regarded as an estimator of δ . If the limiting distribution of \mathbf{U} in (4.74) were the exact distribution of \mathbf{U} for every n under H_n , then the likelihood ratio statistic for testing $\delta = 0$ against $\delta \in C$ based on a single observation of \mathbf{U} would be

$$S_L = \mathbf{U}^T \mathcal{I}_0 \mathbf{U} - \min\{(\mathbf{U} - \mathbf{a})^T \mathcal{I}_0 (\mathbf{U} - \mathbf{a}) : \mathbf{a} \in C\}. \quad (4.75)$$

The arguments leading to this definition of S_L are based on the local properties of \mathbf{U} and $\mathbf{S}(0)$. Therefore, these motivations are closely related to those leading to the classical Rao's score statistic. We shall refer to S_L as the *local score statistic*.

The foregoing local arguments were used only as a motivation and guide to constructing S_L . Whether or not it is a good statistic would depend on its properties and performance. Therefore, it would be of interest to compare the asymptotic performance of S_L with other statistics. Because $n^{1/2}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{0, \mathcal{I}(\theta_0)^{-1}\}$, where θ_0 is the true value, a distance statistic for testing $\theta = 0$ against $\theta \in C$ is

$$D = n[\hat{\theta}^T \hat{\mathcal{I}} \hat{\theta} - \min\{(\hat{\theta} - \mathbf{a})^T \hat{\mathcal{I}}(\hat{\theta} - \mathbf{a}) : \mathbf{a} \in C\}], \quad (4.76)$$

where $\hat{\mathcal{I}}$ is a consistent estimator of \mathcal{I}_0 . It follows from (4.56), (4.75), and (4.76) that $S_L = D + o_p(1)$ under H_0 . Therefore, S_L and D are asymptotically equivalent for testing $\theta = 0$ against $\theta \in C$. These are summarized in the following proposition.

Proposition 4.6.1 *Assume that Condition Q is satisfied. Let \mathbf{U} , S_L , D , and the hypotheses be as in (4.73), (4.75), (4.76), and (4.72), respectively. Then, under H_0 ,*

$$LRT = S_L + o_p(1) = D + o_p(1). \quad (4.77)$$

would be the case under fairly general conditions. In this sense all three tests are asymptotically equivalent. Since there are no nuisance parameters under H_0 , S_L may be computed without estimating the model. By contrast, the LRT requires the model to be estimated subject to the inequality constraints of H_1 , and D requires estimation of the full model without the constraints.

Since the Rao score statistic is based on the score $\mathbf{S}(0)$ as the basic building block and thinking of it as the slope of $\ell(\theta)$ at the null value, one might be tempted to define a score statistic for testing $\theta = 0$ vs $\theta \in C$ by replacing the role of \mathbf{U} in (4.75) by $\mathbf{S}(0)$. This leads to the statistic

$$R^* = \mathbf{S}_0^T \mathcal{I}_0^{-1} \mathbf{S}_0 - \min\{(\mathbf{S}_0 - \mathbf{a})^T \mathcal{I}_0^{-1} (\mathbf{S}_0 - \mathbf{a}) : \mathbf{a} \in C\}. \quad (4.78)$$

Since $n^{-1/2}\mathbf{S}(0) \xrightarrow{d} N(\mathcal{I}_0\delta, \mathcal{I}_0)$ under H_n , this essentially assumes that a test of $\theta = 0$ vs $\mathcal{I}_0\theta \in C$ would also be suitable for testing $\theta = 0$ vs $\theta \in C$. Since this is not a reasonable assumption, S_L appears to be better than R^* .

4.6.2 Local Score Test of $H_0 : \psi = 0$ Against $H_1 : \psi \in C$

It is very rarely that a null hypothesis would specify a value for every component of θ . Therefore, the testing problem considered in the previous subsection is very rare, although the discussions therein were meant to be instructive. In this subsection, we consider the more realistic setting in which the hypotheses impose constraints on some components only. The main ideas leading to the foregoing local score statistic, S_L , can also be applied when there are nuisance parameters.

Let $\theta = (\lambda : \psi)$, where ψ is $k \times 1$, and let the null and alternative hypotheses be

$$H_0 : \psi = 0 \quad \text{and} \quad H_1 : \psi \in C, \quad (4.79)$$

respectively, where C is a closed convex cone in \mathbb{R}^k . It is convenient to adopt the notation in subsection 4.5.2. Let $\mathbf{S}(\theta)$ denote the score function $(\partial/\partial\theta)\ell(\theta)$ and let it also be partitioned as $(\mathbf{S}_\lambda : \mathbf{S}_\psi)$. Let $\bar{\lambda}$ denote the *role* of λ under H_0 defined by $\mathbf{S}_\lambda(\bar{\lambda} : 0) = 0$ and let $\bar{\theta} = (\bar{\lambda} : 0)$ be the *role* of θ under H_0 . Define

$$\mathbf{U} = n^{-1/2}\mathcal{I}^\psi(\bar{\theta})\mathbf{S}_\psi(\bar{\theta}). \quad (4.79)$$

Let δ be a fixed point in \mathbb{R}^k . Then, since $\mathcal{I}_{\psi\psi,\lambda} = (\mathcal{I}_{\psi\psi} - \mathcal{I}_{\psi\lambda}\mathcal{I}_{\lambda\lambda}^{-1}\mathcal{I}_{\lambda\psi})$ and $(\mathcal{I}_{\psi\psi})^{-1} = \mathcal{I}_{\psi\psi,\lambda}$, we have that (see (4.61))

$$\mathbf{U} \xrightarrow{d} N\{\delta, \mathcal{I}_{\psi\psi,\lambda}^{-1}(\theta_0)\}, \text{ under } H_n : \psi = n^{-1/2}\delta, \quad (4.80)$$

where $\theta_0 = (\lambda_0 : 0)$, the true value of θ under H_0 . Now, by arguments similar to those leading to (4.75), we define the *local score statistic* as

$$S_L = \mathbf{U}^T \bar{\mathcal{I}}_{\psi\psi,\lambda} \mathbf{U} - \inf\{(\mathbf{U} - \mathbf{a})^T \bar{\mathcal{I}}_{\psi\psi,\lambda}(\mathbf{U} - \mathbf{a}) : \mathbf{a} \in C\} \quad (4.81)$$

Further, the common asymptotic null distribution of LRT, S_L , and D is $\chi^2(\mathcal{I}_0^{-1}, C)$.

A consequence of this result is that the local score statistic, the likelihood ratio statistic, and the distance statistic have the same asymptotic null distribution. Further, they also have the same local power if (4.77) holds under $H_n : \theta = n^{-1/2}\delta$, which

Because $n^{1/2}(\hat{\psi} - \psi_0) \xrightarrow{d} N\{0, \mathcal{I}_{\psi\psi,\lambda}^{-1}(\theta_0)\}$, a distance statistic for testing H_0 : $\psi = 0$ vs $H_1 : \psi \in C$ is

$$D = n[\hat{\psi}^T \bar{\mathcal{I}}_{\psi\psi,\lambda} \hat{\psi} - \min\{(\hat{\psi} - \mathbf{a})^T \bar{\mathcal{I}}_{\psi\psi,\lambda} (\hat{\psi} - \mathbf{a}) : \mathbf{a} \in C\}] \quad (4.82)$$

where $\bar{\mathcal{I}}_{\psi\psi,\lambda} = \mathcal{I}_{\psi\psi,\lambda}(\hat{\theta})$ which is a consistent estimator of $\mathcal{I}_{\psi\psi,\lambda}(\theta_0)$. Now, we have the following result, which says that the local score statistic S_L in (4.81) is asymptotically equivalent to the likelihood ratio and distance statistics; see Silvapulle and Silvapulle (1995) for a proof.

Proposition 4.6.2 *Let U, S_L, D , and the hypotheses be as in (4.79), (4.81), (4.82), and (4.78), respectively. Then, under H_0 ,*

$$LRT = S_L + o_p(1) = D + o_p(1). \quad (4.83)$$

Further, the common asymptotic null distribution of LRT, S_L , and D is $\bar{\chi}^2(\mathcal{I}_{\psi\psi,\lambda}^{-1}(\theta_0), C)$.

If the alternative hypothesis in (4.78) is $H_1 : \psi \neq 0$ then S_L in (4.81) reduces to the Rao score statistic; further, since S_L is asymptotically equivalent to the LRT and they both are based on the same building block, the S_L in (4.81) can be considered as a natural generalization of Rao's score statistic.

4.6.3 Local Score-Type Test Based on Estimating Equations

By employing arguments similar to those in the previous subsection, the main results in subsection 4.5.3 can be extended to introduce a test based on estimating equations. Let $\theta = (\lambda : \psi)$ where ψ is $k \times 1$ and $\mathbf{S}(\theta) = \mathbf{0}$ be an estimating equation for θ . Initially, we shall assume that Condition A on page 173 is satisfied. In the next subsection, we will relax this further to develop a large class of tests. Consequently, the procedure becomes more flexible; for example, quasi-likelihood, partial likelihood, and Generalized Estimating Equations (GEE) are incorporated.

Let the null and alternative hypotheses be

$$H_0 : \psi = \mathbf{0} \text{ and } H_1 : \psi \in C, \quad (4.84)$$

respectively, where C is a closed convex cone and $\theta_0 = (\lambda_0 : 0)$, the true value of θ when H_0 is true. Let $\bar{\theta} = (\bar{\lambda} : 0)$ denote the estimator obtained by solving $\mathbf{S}_\lambda(\lambda : 0) = \mathbf{0}$; i.e., solve $\mathbf{S}_\lambda(\theta) = \mathbf{0}$ under H_0 . Let

$$U = \bar{\mathbf{G}}_{\psi\psi,\lambda}^{-1} \bar{\mathbf{Z}}, \quad (4.85)$$

where $\bar{\mathbf{G}} = \mathbf{G}(\bar{\theta})$, $\bar{\mathbf{Z}} = \mathbf{Z}(\bar{\theta})$, and $\mathbf{G}(\cdot)$ and $\mathbf{Z}(\cdot)$ are as in subsection 4.5.3. Then, by Proposition 4.5.1,

$$U \xrightarrow{d} N\{\delta, \mathbf{A}_{\psi\psi}(\theta_0)\} \quad \text{under } H_n : \psi = n^{-1/2}\delta, \quad (4.86)$$

where $\mathbf{A}_{\psi\psi} = \mathbf{G}_{\psi\psi,\lambda}^{-1} \mathbf{C}(\mathbf{G}_{\psi\psi,\lambda}^{-1})^T$. Therefore, we define the local score-type test statistic for testing $H_0 : \psi = \mathbf{0}$ against $H_1 : \psi \in C$ as

$$S_L = \mathbf{U}^T \bar{\mathbf{A}}_{\psi\psi}^{-1} \mathbf{U} - \min\{(U - \mathbf{a})^T \bar{\mathbf{A}}_{\psi\psi}^{-1} (U - \mathbf{a}) : \mathbf{a} \in C\}, \quad (4.87)$$

where $\bar{\mathbf{A}}_{\psi\psi} = \mathbf{A}_{\psi\psi}(\bar{\theta})$. Further, the limiting distribution of S_L is $\bar{\chi}^2(\mathbf{A}_{\psi\psi}(\theta_0), C)$ under H_0 .

If $\mathbf{S}(\theta)$ is the likelihood score function $(\partial/\partial\theta)\ell(\theta)$, then $\mathbf{G}(\theta) = \mathbf{V}(\theta) = \mathcal{I}(\theta)$ and

$$\mathbf{U} = n^{1/2} \bar{\mathcal{I}}_{\psi\psi}(\bar{\mathbf{S}}_\psi - \bar{\mathcal{I}}_{\psi\lambda} \bar{\mathcal{I}}_{\lambda\lambda}^{-1} \bar{\mathbf{S}}_\lambda),$$

where $\bar{\mathbf{S}} = \mathbf{S}(\bar{\theta})$ and $\bar{\mathcal{I}} = \mathcal{I}(\bar{\theta})$; consequently, the S_L in (4.87) reduces to the local score statistic in the previous subsection (see 4.81). Further, if the alternative hypothesis is $\psi \neq 0$ then S_L reduces to

$$n^{-1} \{\bar{\mathcal{I}}^{\psi\psi}(\bar{\mathbf{S}}_\psi - \bar{\mathcal{I}}_{\psi\lambda} \bar{\mathcal{I}}_{\lambda\lambda}^{-1} \bar{\mathbf{S}}_\lambda)\}^T \bar{\mathcal{I}}_{\psi\psi,\lambda}^{-1} \{\bar{\mathcal{I}}^{\psi\psi}(\bar{\mathbf{S}}_\psi - \bar{\mathcal{I}}_{\psi\lambda} \bar{\mathcal{I}}_{\lambda\lambda}^{-1} \bar{\mathbf{S}}_\lambda)\},$$

which is Neyman's $C(\alpha)$ -statistic. In view of these observations, the local score-type statistic, S_L , in (4.87) may be seen as a generalization of the $C(\alpha)$ -statistic. Since $\mathbf{S}(\theta)$ is not necessarily the likelihood score function, the S_L in (4.87) and the LRT may not be equivalent.

4.6.4 A General Local Score-Type Test of $H_0 : \psi = 0$

This section assumes familiarity with the definition and properties of *approximating cone* introduced later in this chapter. However, the topic of this subsection fits in better in this section, because it is a direct extension of the results in the previous subsections. An approximating cone is purely a mathematical object. The introduction to this topic in Section 4.7 can be read independently; readers who are not familiar with it (i.e., approximating cone) may wish to read about it before reading the rest of this section.

So far we have assumed that the alternative hypothesis is of the form $\psi \in C$ where C is a closed convex cone. The results in the previous subsection extend in a natural way to provide a large class of tests of $\psi = 0$ against quite general alternatives based on quantities that are not necessarily the likelihood score function. Let $\mathbf{0} \in \Psi \subset \mathbb{R}^k$, and let the testing problem be

$$H_0 : \psi = \mathbf{0} \quad \text{vs} \quad H_1 : \psi \in \Psi.$$

Although this formulation allows $\mathbf{0}$ to be an interior point of Ψ , our main focus will be on the case when it is on the boundary of Ψ . Let \mathcal{A} denote the *approximating cone* of Ψ at the null value $\mathbf{0}$. If the null hypothesis is formulated as $H_0 : \psi = \psi_0$, then \mathcal{A} should be the approximating cone of Ψ at ψ_0 . Although we will continue to use the notation $\theta = (\lambda : \psi)$, here we allow λ to be an infinite dimensional nuisance parameter. Thus, *semiparametric* models are also incorporated in this subsection. Let $\theta_0 = (\lambda_0 : \mathbf{0})$ and let $\delta \in \mathcal{A}$ be arbitrary but fixed. Now, let us introduce the following condition:

Condition LS:

1. There exists a statistic U (i.e., a function of the data only) and a positive definite matrix $D(\theta_0)$ such that

$$U \xrightarrow{d} N\{\delta, D(\theta_0)\} \text{ under } H_n : \psi = n^{-1/2}\delta.$$

2. There exists a consistent estimator, \bar{D} of $D(\theta_0)$. ■

Assume that this condition is satisfied. Then, a local score-type statistic for testing $H_0 : \psi = 0$ against $H_1 : \psi \in \Psi$ is

$$S_L = U^T \bar{D}^{-1} U - \inf_{b \in \mathcal{A}} (U - b)^T \bar{D}^{-1} (U - b). \quad (4.88)$$

This is the general form of the local score-type statistic introduced in Silvapulle and Silvapulle (1995) and further developed in Silvapulle et al. (2002) for the semiparametric model. If \mathcal{A} is a closed convex cone then the limiting distribution of S_L is given by

$$S_L \xrightarrow{d} \chi^2\{D(\theta_0), \mathcal{A}\} \text{ under } H_0.$$

The asymptotic distribution of S_L under H_n is equal to the expression on the right-hand side (4.88) with $U \sim N\{\delta, D(\theta_0)\}$ and \bar{D} replaced by $D(\theta_0)$. If λ is finite dimensional and U is the efficient score then S_L is asymptotically equivalent to the likelihood ratio test [see Theorem 1 in Silvapulle and Silvapulle (1995)].

The test statistic in (4.88) provides a very flexible class of simple tests because its implementation requires only the null model to be estimated and it is based on essentially the same building blocks as those used for the traditional score tests. The applicability of the foregoing local score test in several semiparametric models (including linear regression, AR(p), ARCH, and nonlinear regression models) in which the error distribution is treated as an unknown infinite dimensional nuisance parameter, is discussed in Silvapulle et al. (2002) where *adaptive tests and efficient semiparametric tests* are developed. The results in Silvapulle et al. (2002) suggest that once an adaptive test is obtained for testing against an unconstrained hypothesis, it is likely that the same test can be extended to incorporate restricted alternatives as well. The reason for this is that a semiparametric test against an unrestricted alternative is likely to be based on a U satisfying Condition LS.

Let the testing problem be $H_0 : \psi = 0$ vs $H_1 : \psi \in C$, where C is a closed convex cone. Assume that there exists a vector c in C such that $c^T x > 0$ for every $x \in C \setminus \{0\}$. For example, if C is the positive orthant, then c could be the center direction $(1, \dots, 1)^T$; if C is a polyhedral generated by a_1, \dots, a_k then it suffices to find a c such that $a_i^T c > 0$ for $i = 1, \dots, k$. Now, the basic idea is to adopt the working assumption that $\psi = \eta c$, and test $\eta = 0$ against $\eta > 0$. It appears reasonable to expect that a test that can detect a departure in the direction of c would also be able to detect departures in other directions of C that are close to c as well. Let U be as in (4.85). Then, it follows from (4.86) that

$$c^T U \xrightarrow{d} N\{c^T \delta, c^T A_{\psi\psi}(\theta_0)c\} \text{ under } H_n : \psi = n^{-1/2}\delta. \quad (4.89)$$

Because c was chosen such that $c^T \delta \geq 0$ for $\delta \in C$, a simple test of $H_0 : \psi = 0$ against $H_1 : \psi \in C$ rejects H_0 for large values of T_c where

$$T_c = c^T U / \{c^T \bar{A}_{\psi\psi} c\}^{1/2}, \quad (4.90)$$

and $\bar{A}_{\psi\psi}$ is a consistent estimator of $A_{\psi\psi}(\theta_0)$. Clearly, this test is considerably simpler than the ones in the previous subsections. The cost of simplicity is that this test is unlikely to be as powerful as the local score test, S_L , on average, although T_c is likely to be more powerful than S_L in the direction of c .

A closely related procedure was suggested by King and Wu (1997); this test has some locally optimal properties. For this test, let $S(\theta)$ be the likelihood score function. It follows from the foregoing results that

$$n^{-1/2} c^T S_\psi(\bar{\theta}) \xrightarrow{d} N\{c^T \mathcal{I}_{\psi\psi,\lambda}(\theta_0)\delta, c^T \mathcal{I}_{\psi\psi,\lambda}(\theta_0)c\} \text{ under } H_n : \psi = n^{-1/2}\delta.$$

Now, suppose that the following is satisfied: $c^T \mathcal{I}_{\psi\psi,\lambda}(\theta_0)\delta > 0$ for $\delta \in C \setminus \{0\}$. This condition would be satisfied if the elements of $\mathcal{I}_{\psi\psi,\lambda}(\theta_0)$ are nonnegative and $c = (1, \dots, 1)^T$. Then a test of H_0 against H_1 rejects H_0 if

$$n^{-1/2} c^T S_\psi(\bar{\theta}) / \{c^T \bar{\mathcal{I}}_{\psi\psi,\lambda}(\bar{\theta})\}^{1/2}$$

is large compared with the critical value from the standard normal distribution. For an example to illustrate the application of this to ARCH/GARCH models in time series, see Lee and King (1993); see also King and Wu (1997).

4.6.6 A Data Example: Test Against ARCH Effect in an ARCH Model.

A common feature of many financial time series such as commodity prices, exchange rates, and stock returns is that periods of large changes tend to cluster together, and similarly periods of small changes tend to cluster together. Because high variability is associated with high risk, studying the variability of such time series is important. The AutoRegressive Conditional Heteroscedasticity (ARCH) model and its various generalizations/modifications are now widely used in modeling such volatility clustering patterns. For a comprehensive account of ARCH modeling see Gouriéroux (1997).

4.6.5 A Simple One-Dimensional Test of $\psi = 0$ Against $\psi \in C$

The local score-type statistic, S_L , in (4.87) is based on the distribution of the score vector, S , and its asymptotic null distribution would typically be a chi-bar-square. Although the computation of this may be considerably easier than that of the LRT, it is not as easy as the usual t -ratio for a single parameter. There is a need for a test that is simpler than S_L even if it is less powerful. To this end, we introduce the following simple test.

Let \mathcal{F}_t denote the information available up to time t , Y_t denote the dependent variable, and x_t be a column vector of explanatory variables. Typically, x_t includes lagged values of the dependent variable. The p th order ARCH model of Engle (1982) can be stated as

$$Y_t = \mathbf{x}_t^T \boldsymbol{\beta} + \epsilon_t, E(\epsilon_t | \mathcal{F}_{t-1}) = 0, \text{ and } h_t = \alpha_0 + \psi_1 \epsilon_{t-1}^2 + \dots + \psi_p \epsilon_{t-p}^2, \quad (4.91)$$

where $h_t = \text{var}(Y_t | \mathcal{F}_{t-1})$. Note that (4.91) provides a specification not only for the mean but also for the variance. If ψ_j in (4.91) is negative for some j ($j = 1, \dots, p$), then a large value for ϵ_{t-j} would result in a negative variance for $Y_t | \mathcal{F}_{t-1}$. Therefore, the admissible range for ψ_1, \dots, ψ_p is $\{\psi_1 \geq 0, \dots, \psi_p \geq 0\}$. In the context of this example, the following question is important: “is the variance, h_t , constant over time?” As a general rule, this testing problem is formulated as

$$H_0 : \psi_1 = \dots = \psi_p = 0 \quad \text{against} \quad H_1 : \psi_1 \geq 0, \dots, \psi_p \geq 0. \quad (4.92)$$

In this data example, we consider an ARCH model for an overall rate of return from stocks. Let $Y_t = \log(I_t/I_{t-1})$ where I_t is the *All Ordinaries Index (Australia)* on day t . This index is a weighted average of the prices of selected shares in Australia; it corresponds to the *Dow Jones Index* in the United States. The data that we used are for the period 5/January/1984 to 29/November/1985. This provides a total of 484 observations; see Silvapulle and Silvapulle (1995) for the data. For illustrative purposes, we chose $p = 3$; the method would be applicable for other values of p . The specific model that we consider is

$$\mathbf{x}_t^T \boldsymbol{\beta} = \beta_0 + \beta_1 Y_{t-1} + \dots + \beta_4 Y_{t-4} \text{ and } h_t = \alpha_0 + \psi_1 \epsilon_{t-1}^2 + \psi_2 \epsilon_{t-2}^2 + \psi_3 \epsilon_{t-3}^2. \quad (4.93)$$

Let $\boldsymbol{\theta} = (\boldsymbol{\lambda} : \boldsymbol{\psi})$ where $\boldsymbol{\lambda} = (\beta_0, \dots, \beta_4, \alpha_0)^T$ and $\boldsymbol{\psi} = (\psi_1, \psi_2, \psi_3)^T$. The null and alternative hypotheses are

$$H_0 : \boldsymbol{\psi} = \mathbf{0} \text{ and } H_1 : \boldsymbol{\psi} \geq 0.$$

Let $\mathbf{S}(\boldsymbol{\theta}) = (\partial/\partial\boldsymbol{\theta})L(\boldsymbol{\theta})$ where $L(\boldsymbol{\theta}) = -(1/2)\sum\{\log(2\pi h_t) + \epsilon_t^2/h_t\}$; thus $L(\boldsymbol{\theta})$ is the loglikelihood if we were to assume that $Y_t | \mathcal{F}_{t-1} \sim N(\mathbf{x}_t^T \boldsymbol{\beta}, h_t)$. Now, let us compute the local score-type test statistic S_L in (4.87). First, we need to estimate the nuisance parameter under the null hypothesis. Therefore, we solve the equation $\mathbf{S}_{n,\lambda}(\bar{\boldsymbol{\theta}}) = \mathbf{0}$ where $\bar{\boldsymbol{\theta}} = (\bar{\boldsymbol{\lambda}} : \mathbf{0})$. Under the null hypothesis, the error variance is constant, and therefore the null model is estimated by ordinary least squares. Thus we have that $\bar{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ and $\bar{a}_0 = n^{-1} \sum(Y_t - \mathbf{x}_t^T \bar{\boldsymbol{\beta}})^2$. Now, $\bar{\boldsymbol{\lambda}} = (\bar{\boldsymbol{\beta}} : \bar{a}_0)$.

Note that to apply the local score-type test statistic S_L , this is the only estimation required despite the fact that the ARCH model itself is rather involved. The next step is to compute \bar{G} , \bar{V} , \bar{Z} , U , and $\bar{A}_{\psi\psi}$. Closed form expressions for these quantities are obtained directly from the closed-form available for $\mathbf{S}(\boldsymbol{\theta})$; see Silvapulle and Silvapulle (1995) for these formulas.

The computed value of U is $(3.56, 1.30, 2.20)^T$. Because every component of U is positive, $\min\{(U - \mathbf{a})^T \bar{A}_{\psi\psi}^{-1} (U - \mathbf{a}) : \mathbf{a} \geq 0\} = 0$ and hence

$$S_L = U^T \bar{A}_{\psi\psi}^{-1} U - \min\{(U - \mathbf{a})^T \bar{A}_{\psi\psi}^{-1} (U - \mathbf{a}) : \mathbf{a} \geq 0\} = U^T \bar{A}_{\psi\psi}^{-1} U = 4.36.$$

Let

$$\zeta(\boldsymbol{\lambda}) = [\sum_{i=1}^3 w_i \{3, A_{\psi\psi}(\boldsymbol{\theta})\} \text{pr}(\chi_i^2 \geq 4.36)]_{\boldsymbol{\theta}=(\boldsymbol{\lambda}:0)}.$$

Using the closed forms given in the last chapter for the chi-bar square weights in this expression (see (3.25)), we can compute $\sup_{\boldsymbol{\lambda}} \zeta(\boldsymbol{\lambda})$. However, before we do so, let us apply a bounds test. It follows from (3.23) that

$$\zeta(\boldsymbol{\lambda}) \leq 0.5 \{\text{pr}(\chi_2^2 \geq 4.36) + \text{pr}(\chi_3^2 \geq 4.36)\} = 0.17.$$

This upper bound is too high to reject H_0 . Therefore, we computed $p_{sup} = \sup_{\boldsymbol{\lambda}} \zeta(\boldsymbol{\lambda})$ and $\hat{p} = \zeta(\bar{\boldsymbol{\lambda}})$. The computed values are 0.10 and 0.085, respectively. Thus, it appears that there may be some evidence in favor of ARCH effect, but it is weak.

For testing $H_0 : \boldsymbol{\psi} = \mathbf{0}$ against $H_1 : \boldsymbol{\psi} \geq 0$, it is also valid to apply an unrestricted test ignoring the inequality constraints in H_1 . To this end, let the testing problem be

$$H_0 : \boldsymbol{\psi} = \mathbf{0} \quad \text{vs} \quad H_2 : \boldsymbol{\psi} \neq \mathbf{0}.$$

Even though the ARCH model would not be well-defined with a negative value for any of the components of $\boldsymbol{\psi}$, it is valid to apply a score-type test of $H_0 : \boldsymbol{\psi} = \mathbf{0}$ against $H_2 : \boldsymbol{\psi} \neq \mathbf{0}$; in fact, this is the procedure that is found in much of the econometrics literature. The only information the score-type statistic $U^T \bar{A}_{\psi\psi}^{-1} U$ uses is that it is approximately $N(0, 1)$ under H_0 ; if H_0 is not true then it is expected to be large. The sample value of this is

$$U^T \bar{A}_{\psi\psi}^{-1} U = 4.36;$$

and its large sample p -value is $\text{pr}(\chi_3^2 \geq 4.36) = 0.22$. As expected, because the components of U turned out to be all positive, the statistics for the one-sided and two-sided tests have the same numerical value and the latter has a larger p -value.

4.7 APPROXIMATING CONES AND TANGENT CONES

Most of the results in the previous sections have generalizations to the case when the parameter space takes more general forms. Typically, such parameter spaces arise because the constraints are nonlinear in the parameter. It turns out that the main first-order results, namely the asymptotic distribution of an estimator, the asymptotic null distribution of a test statistic, and the asymptotic local power of a test are unchanged when the null and alternative parameter spaces are replaced by cones that approximate them at the true value of the parameter provided that the boundaries of the parameter spaces are sufficiently smooth. The objective of this section is to introduce the basics that are required to study the asymptotics when the parameter spaces are neither linear spaces nor cones; for example, the parameter space may have boundaries that are curved or the slope of the boundary may be discontinuous (i.e., kinked).

As an example, let $\Omega = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \theta_2 \geq g_1(\theta_1) \text{ and } \theta_2 \geq g_2(\theta_1)\}$ where $g_1(\theta_1) = \theta_1^2 + \theta_1$ and $g_2(\theta_1) = \theta_1^2 - \theta_1$. In Fig. 4.1, O is the origin; the two

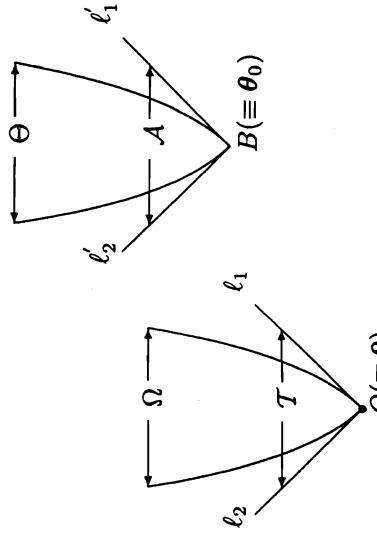


Fig. 4.2 The cone of tangents, T , and the approximating cone, \mathcal{A} , of Θ at B .

respectively, where $\theta_0 = (1, 1)^T$, $\Theta = \theta_0 + \Omega$ and Ω is as in Fig. 4.2. Let $\hat{\theta}_\Theta$, $\hat{\theta}_\Omega$ and $\hat{\theta}_T$ denote the mle over Θ , Ω , and T , respectively. Then

$$LRT = 2\{\ell(\hat{\theta}_\Theta) - \ell(\theta_0)\} = n\{\|\bar{Y} - \theta_0\|^2 - \|\bar{Y} - \Theta\|^2\}$$

where $\|\bar{Y} - \Theta\|^2 = \inf_{\theta \in \Theta} (\bar{Y} - \theta)^T (\bar{Y} - \theta)$ is the squared distance between \bar{Y} and Θ . Suppose that the true value of θ is θ_0 . Now, the results to be presented later in this section show that the asymptotic distribution of the constrained mle $\hat{\theta}_\Theta$ and that of the LRT of $\theta = \theta_0$ against $\theta \in \Theta$ remain the same even if the constrained parameter space Θ is replaced by \mathcal{A} , the approximating cone of Θ at θ_0 . More specifically,

$$\begin{aligned} LRT &= n\{\|\bar{Y} - \theta_0\|^2 - \|\bar{Y} - \Theta\|^2\} \\ &= n\{\|\bar{Y} - \theta_0\|^2 - \|\bar{Y} - \mathcal{A}\|^2\} + o_p(1), \end{aligned} \quad (4.94)$$

and $n^{1/2}(\hat{\theta}_\Theta - \hat{\theta}_\mathcal{A}) = o_p(1).$

Further, the constrained estimator $\hat{\theta}_\Theta$ and the projection of the unconstrained estimator $\hat{\theta}$ onto the approximating cone have the same asymptotic distributions in the following sense:

$$\sqrt{n}(\hat{\theta}_\Theta - \theta_0) \xrightarrow{d} \Pi(Z, T) \quad \text{and} \quad \sqrt{n}\{\Pi(\hat{\theta}, \mathcal{A}) - \theta_0\} \xrightarrow{d} \Pi(Z, T),$$

where $Z \sim N(0, I)$.

It is worth noting that the process of replacing Θ by \mathcal{A} is equivalent to replacing the nonlinear constraints defining Θ by their linear approximations at the null value.

In the following subsections, we will extend these ideas to more general settings. To do so, one important requirement is that the parameter space must satisfy certain regularity conditions. These conditions must be satisfied irrespective of the properties of the likelihood function and other stochastic conditions. These are discussed in the next subsection.

Before we introduce regularity conditions relating to parameter spaces, let us make a remark regarding our terminology. Recall that a cone is defined to have vertex at the

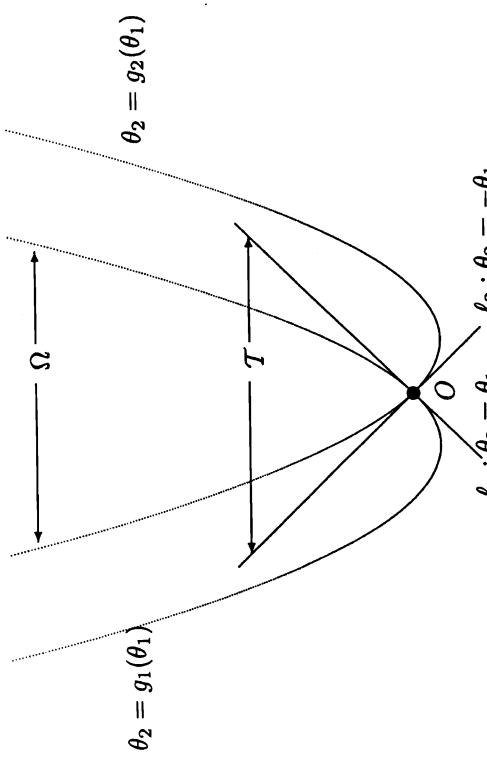


Fig. 4.1 The cone of tangents, T , of Ω at O .

parabolas, $\theta_2 = g_1(\theta_1)$ and $\theta_2 = g_2(\theta_1)$, are also shown. Further, Ω is the region bounded by these two parabolas, as shown.

Let ℓ_1 and ℓ_2 be the tangents at O to the two smooth curves that form the boundary of Ω ; thus the equations for ℓ_1 and ℓ_2 are $\theta_2 = \theta_1$ and $\theta_2 = -\theta_1$, respectively. Let $T = \{\theta : \theta_2 \geq \theta_1 \text{ and } \theta_2 \geq -\theta_1\}$. Locally at the origin, Ω is close to the cone T . In this example, T is known as the *linearizing cone* of Ω at O , because it is obtained by linearizing the boundaries defining the set Ω at O . The cone is also called the *contingent cone* and *cone of tangents* of Ω at O . The formal definitions of these terms will be introduced in the next subsection.

Let $\theta_0 (\neq 0)$ be a given fixed point in \mathbb{R}^p . Now, let us move $\{\Omega, \ell_1, \ell_2, T\}$ by θ_0 , and let the new objects be denoted by $\{\Theta, \ell'_1, \ell'_2, \mathcal{A}\}$ (see Fig. 4.2). Therefore, $\Theta = \theta_0 + \Omega$, $\mathcal{A} = \theta_0 + T$, ℓ'_1 and ℓ'_2 are the tangents to the boundaries of Θ at θ_0 , and \mathcal{A} is the cone formed by ℓ'_1 and ℓ'_2 . Then \mathcal{A} is close to Θ at θ_0 , in the same way as T is close to Ω at O . The cone \mathcal{A} with vertex at θ_0 is called the *approximating cone* of Θ at θ_0 . An important result relating to the approximating cone is that if Θ is the parameter space and θ_0 is the true value then Θ can be replaced by its approximating cone \mathcal{A} at θ_0 for obtaining the asymptotic null distribution of LRT and other first order asymptotic results when θ_0 is the true value.

Example 4.7.1 *Inference on normal mean when the parameter space is not a cone.*

Let Y_1, \dots, Y_n be iid as the bivariate normal distribution, $N(\theta, I)$, where $\theta \in \mathbb{R}^2$ and let $\ell(\theta)$ denote the kernel of the loglikelihood $-(n/2)\|\bar{Y} - \theta\|^2$. Let the null and alternative hypotheses be

$$H_0 : \theta = \theta_0 \quad \text{and} \quad H_1 : \theta \in \Theta,$$

origin unless the contrary is clear (see Appendix 1). According to our definition, the approximating cone does not necessarily have its vertex at the origin. Some authors refer to the \mathcal{T} in Fig. 4.2 as the approximating cone of Θ at θ_0 . We shall deviate from this terminology because it would be helpful to identify \mathcal{T} and \mathcal{A} separately using different terms; we shall refer to them as tangent cone of Θ at θ_0 and approximating cone of Θ at θ_0 , respectively.

4.7.1 Chernoff Regularity

Let $\Theta \subset \mathbb{R}^p$ and $\theta_0 \in \Theta$. A vector $w \in \mathbb{R}^p$ is said to be a *tangent* to Θ at θ_0 if there exists a sequence $\{\theta_n\}$ in Θ and a sequence of positive numbers $\{t_n\}$ converging to zero such that $t_n^{-1}(\theta_n - \theta_0) \rightarrow w$ as $n \rightarrow \infty$. This essentially means that $w \in \mathbb{R}^p$ is a tangent to Θ at θ_0 if either $w = 0$ or there exists a sequence of points in Θ converging to θ_0 such that the direction of the vector from θ_0 to θ_n converges to the direction of w . Therefore, a tangent at θ_0 is a vector rather than the set of points on a line passing through θ_0 . This definition of a tangent is more general than the usual one where it is seen as a line that touches a smooth curve at a point. The following examples illustrate this generalized notion of a tangent.

Example 4.7.2 Illustration of a generalized notion of tangents

1. For $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 = \theta_1^2\}$, the tangents at 0 are $k(1, 0)$ and $k(-1, 0)$ where $k \geq 0$ (see Fig. 4.3). We do not view the line $\theta_2 = 0$ as a tangent.
2. For $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 = \theta_1^2 \text{ and } \theta_1 \geq 0\}$ the only tangent at 0 is $k(1, 0)$ where $k \geq 0$; note that $(-1, 0)$ is not a tangent at 0.
3. For $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 \geq \theta_1^2\}$, the tangents at 0 are (θ_1, θ_2) where $\theta_2 \geq 0$ and $\theta_1 \in \mathbb{R}$ (see Fig. 4.4). Thus, the tangents to the curve $\theta_2 = \theta_1^2$, and those to the area bounded by the curve are different.

4. Let $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 = \sin(1/\theta_1) \text{ and } \theta_1 > 0\} \cup \{0\}$. The tangents to Θ at 0 are $k\theta$ where $\theta = (\theta_1, \theta_2), \theta_1 > 0$ and $k \geq 0$. Thus, even though $(d\theta_2/d\theta_1)$ does not exist at the origin, the set $\{\theta : \theta_2 = \sin(1/\theta_1)\}$ has an infinite number of tangents at the same point. ■

The *cone of tangents*, denoted $\mathcal{T}(\Theta; \theta_0)$, of Θ at θ_0 is defined as the set of all tangents to Θ at θ_0 . This cone is also called the *contingent cone*, *Boulingard tangent cone*, and the *ordinary tangent cone*. Note that $\mathcal{T}(\Theta; \theta_0)$ has vertex at the origin irrespective of whether or not θ_0 is the origin. A tangent v is said to be *derivable* if there exists a function $f : [0, \epsilon] \rightarrow \Theta$ for some $\epsilon > 0$ such that $f(0) = \theta_0$ and $(d/dt+)f(0) = v$. Intuitively, this means that there is a smooth curve in Θ with one end at θ_0 such that the slope of the curve at θ_0 is parallel to v . The set of all derivable tangents to Θ at θ_0 is called the *derived tangent cone* of Θ at θ_0 . It is also called the *cone of attainable directions*; see Bazaraa et al (1993).

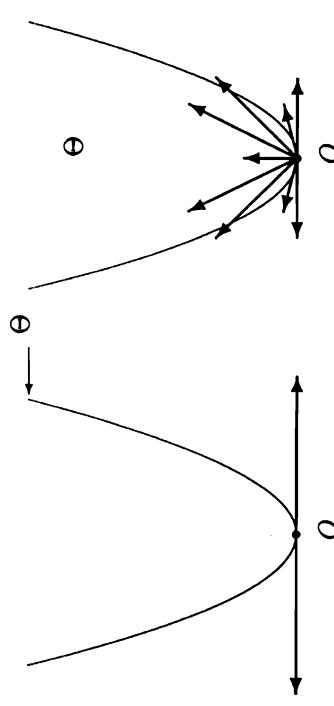


Fig. 4.3 The tangents to $\{\theta_2 = \theta_1^2\}$ at O .

Fig. 4.4 Tangents to $\{\theta_2 \geq \theta_1^2\}$ at O .

While it is not essential for us at this stage, let us also state equivalent definitions of the cones just introduced. The derived tangent cone is equal to

$$\{v : \forall t_n \downarrow 0, \exists \theta_n \in \Theta \text{ such that } \theta_n \rightarrow \theta_0 \text{ and } t_n^{-1}(\theta_n - \theta_0) \rightarrow v\},$$

and the cone of tangents is equal to

$$\{v : \exists t_n \downarrow 0, \exists \theta_n \in \Theta \text{ such that } \theta_n \rightarrow \theta_0 \text{ and } t_n^{-1}(\theta_n - \theta_0) \rightarrow v\}.$$

It is clear from these definitions that the derived tangent cone is contained in the cone of tangents, because the former must satisfy a condition for all $t_n \downarrow 0$ while the cone of tangents must satisfy the same condition for some $t_n \downarrow 0$. The following example of (Rockafellar and Wets, 1998, p 199) is helpful to illustrate the difference between the cone of tangents and the derived tangent cone.

Example 4.7.3 Illustration of a difference between cone of tangents and the derived tangent cone

Let $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 = \theta_1 \sin(1/\theta_1), \theta_1 > 0\} \cup \{0\}$. Note that there is no curve that lies in Θ and is smooth at 0. Therefore, the only member of the derived tangent cone is $\{0\}$. By contrast, the straight line $\theta_2 = k\theta_1$ where $|k| \leq 1$, intersects with Θ at a sequence of points $\{\theta_n\}$ and $\|\theta_n\|^{-1}\theta_n \rightarrow (1, k)$. Therefore, the cone of tangents is $\{\theta \in \mathbb{R}^2 : |\theta_2| \leq \theta_1, \theta_1 > 0\}$. ■

The foregoing definitions and the example suggest that if the cone of tangents and the derived tangent cone of Θ at θ_0 are not equal, then Θ is likely to be quite irregular near θ_0 . A set $\Theta \subset \mathbb{R}^p$ is said to be *Chernoff regular* at θ_0 if the cone of tangents and the derived tangent cone of Θ at θ_0 are equal. If Θ is Chernoff regular at θ_0 then we shall refer to the *cone of tangents* and the *derived tangent cone* as *tangent cone*. An intuitively simple way of verifying Chernoff regularity of Θ at θ_0 is the following (this is due to Rockafellar and Wets (1998), page 198). Consider the set $\mathcal{A}_t = t(\Theta - \theta_0)$ where $t > 0$. Now take the limit as t increases to ∞ (see Fig. 4.5). Note that \mathcal{A}_t provides a magnified view of Θ near θ_0 because it is like viewing Θ

through a magnifying glass focused at θ_0 . As t increases to ∞ , the magnification increases to ∞ and if the magnified view of Θ at θ_0 converges to a limit (i.e., if \mathcal{A}_t has a limit) then Θ is Chernoff regular at θ_0 , and the limit is the tangent cone of Θ at θ_0 .

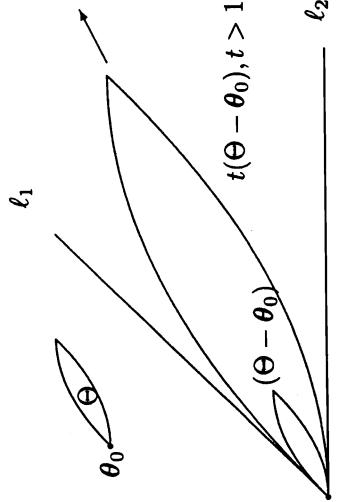


Fig. 4.5 The limit of $t(\Theta - \theta_0)$ as $t \rightarrow \infty$ is the cone between the lines ℓ_1 and ℓ_2 , and it is equal to the tangent cone, $T(\Theta; \theta_0)$, of Θ at θ_0 .

Reproduced from: *Variational Analysis*, Rockafellar and Wets, Copyright (1998), with permission from Springer, New York.

Let $\Theta \subset \mathbb{R}^p$ and $\theta_0 \in \Theta$. The set Θ is said to be *approximated by a cone \mathcal{A} at θ_0* if

- (a) $d(\theta, \mathcal{A}) = o(\|\theta - \theta_0\|)$ for $\theta \in \Theta$,
- and (b) $d(x, \Theta) = o(\|x - \theta_0\|)$ for $x \in \mathcal{A}$,

where d is the distance between a point and a set defined by

$$d(\theta, \mathcal{A}) = \inf_{x \in \mathcal{A}} \|\theta - x\|.$$

In this case, \mathcal{A} is called an *approximating cone of Θ at θ_0* and will be denoted by $\mathcal{A}(\Theta_0; \theta_0)$.

This definition was introduced by Chernoff (1954) and it is illustrated in Fig. 4.6 using two different shapes for Θ . In each of the two cases, let Θ be the region bounded by the two curves as shown and O be the point corresponding to θ_0 . In the diagram on the right-hand side, A is a point in Θ . Let $\theta = OA$ and B be the projection of A onto the approximating cone \mathcal{A} ; therefore, $\|AB\| = d(\theta, \mathcal{A})$. Condition (a) of the definition (4.95) says that $\|AB\|/\|OA\| \rightarrow 0$ as A converges to O while staying in Θ . In the diagram on the left-hand side, A is a point in \mathcal{A} and let $x = OA$. Let B be the projection of A on to Θ . Condition (b) of the definition (4.95) says that $\|AB\|/\|OA\| \rightarrow 0$ as A converges to O while staying in \mathcal{A} .

Andrews (1999, page 1358) noted that the foregoing definition of approximating cone can be stated using Hausdorff distance, see also Le Cam (1970, p 819). For two sets A and B , define the *Hausdorff distance between A and B* by

$$h(A, B) = \max\{\sup_{a \in A} d(a, B), \sup_{b \in B} d(b, A)\}.$$

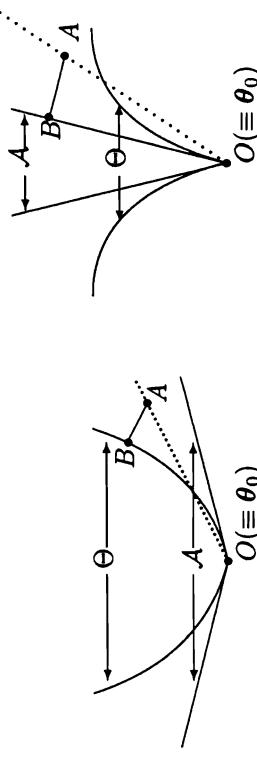


Fig. 4.6 Chernoff's definition of approximating cone \mathcal{A} of Θ at $\theta_0 : \|AB\| = o(\|OA\|)$

Now, the next result provides an equivalent form of definition (4.95).

Proposition 4.7.1 *Let the setting be as in (4.95) and $B_\epsilon = \{\theta : \|\theta - \theta_0\| < \epsilon\}$, where $\epsilon > 0$. Let \mathcal{A} be a cone with vertex at θ_0 . Then \mathcal{A} is an approximating cone of Θ at θ_0 as defined by (4.95) if and only if*

$$h\{\Theta \cap B_\epsilon, \mathcal{A} \cap B_\epsilon\} = o(\epsilon) \quad \text{as } \epsilon \searrow 0. \quad (4.96)$$

Proof: For a given $\epsilon > 0$ let $\Theta_\epsilon = \Theta \cap B_\epsilon$, $\mathcal{A}_\epsilon = \mathcal{A} \cap B_\epsilon$, $\theta_\epsilon = \arg \max_{\theta \in \Theta_\epsilon} d(\theta, \mathcal{A}_\epsilon)$, and $x_\epsilon = \arg \max_{x \in \mathcal{A}_\epsilon} d(x, \Theta_\epsilon)$. Then, the Hausdorff distance between Θ_ϵ and \mathcal{A}_ϵ takes the form

$$h(\Theta_\epsilon, \mathcal{A}_\epsilon) = \max\{d(\theta_\epsilon, \mathcal{A}_\epsilon), d(x_\epsilon, \Theta_\epsilon)\}.$$

Now suppose that \mathcal{A} is an approximating cone according to Chernoff's definition (4.95). Then

$$d(\theta_\epsilon, \mathcal{A}_\epsilon) = o(\|\theta_\epsilon - \theta_0\|) = o(\epsilon);$$

the first equality follows by (4.95) and the second follows since $\|\theta_\epsilon - \theta_0\| \leq \epsilon$. Similarly, $d(x_\epsilon, \Theta_\epsilon) = o(\epsilon)$. Therefore, $h(\Theta_\epsilon, \mathcal{A}_\epsilon) = o(\epsilon)$.

Now to prove the converse, suppose that (4.96) is satisfied. Let $\theta \in \Theta$ and $\delta = \|\theta - \theta_0\|$. Then,

$$d(\theta, \mathcal{A}) = d(\theta, \mathcal{A}_\delta) \leq d(\theta_\delta, \mathcal{A}_\delta) = o(\delta);$$

the first equality follows because $\|\Pi(\theta, \mathcal{A}) - \theta_0\| \leq \|\theta - \theta_0\|$ and hence $\Pi(\theta, \mathcal{A}) = \Pi(\theta_\delta, \mathcal{A})$, the middle inequality follows by the definition of θ_δ , and the third equality follows by (4.96). ■

For a discussion about Hausdorff distance in statistics, see van der Vaart and Wellner (1996, p 162). Andrews (1999, p 1358) extended the definition of Chernoff (1954) to define a sequence of sets being approximated by a cone; such a generalization may be required/helpful in more general settings. We say that a sequence of sets $\{\Theta_n\}$ is locally approximated at θ_0 by a cone \mathcal{A} if

- (a) $d(\theta_n, \mathcal{A}) = o(\|\theta_n - \theta_0\|)$ for $\theta_n \in \Theta_n$ and $\|\theta_n - \theta_0\| \rightarrow 0$
- and (b) $d(x_n, \Theta_n) = o(\|x_n - \theta_0\|)$ for $x_n \in \mathcal{A}$ and $\|x_n - \theta_0\| \rightarrow 0$.

Now the foregoing results for Θ also hold for Θ_n with appropriate modifications.

Throughout this book, parameter spaces will be assumed to be Chernoff regular unless the contrary is made clear; this is a very mild condition as was indicated earlier. The following important result, due to Geyer (1994, Theorem 2.1), establishes the link between the notion of approximating cone introduced by Chernoff and several tangent cones that have been known in the mathematics literature.

Proposition 4.7.2 (Geyer 1994) *The cone K is an approximating cone of Θ at θ_0 if and only if the derived tangent cone and the cone of tangents are equal and*

$$\text{Closure of } K = \theta_0 + \mathcal{T}(\Theta; \theta_0) = \theta_0 + \text{Derived tangent cone of } \Theta \text{ at } \theta_0.$$

■ This essentially says that Θ has an approximating cone at θ_0 if and only if Θ is Chernoff regular at θ_0 .

The following result provides a convenient way of obtaining the approximating cone of Θ at θ_0 when Θ is defined by a set of nonlinear equality and inequality constraints, in many cases; the proof is given in the appendix to this chapter.

Proposition 4.7.3 *Suppose that Ω is open and let*

$$\Theta = \{\theta \in \Omega : h_1(\theta) = \dots = h_\ell(\theta) = 0, h_{\ell+1}(\theta) \geq 0, \dots, h_k(\theta) \geq 0\},$$

where h_1, \dots, h_k are continuously differentiable. Let θ_0 be a point in Θ , and let $a_i = (\partial/\partial\theta)h_i(\theta_0)$ for $i = 1, \dots, k$, and $J(\theta_0) = \{i : h_i(\theta_0) = 0 \text{ and } \ell + 1 \leq i \leq k\}$. Assume that the following condition, known as the Mangasarian-Fromowitz constraint qualification (MF-CQ), is satisfied at θ_0 : There exists a nonzero $b \in \mathbb{R}^p$ such that $a_1^T b = \dots = a_\ell^T b = 0, a_1, \dots, a_\ell$ are linearly independent and $a_i^T b > 0$ for $i \in J(\theta_0)$. Then $\mathcal{A}(\Theta; \theta_0)$ is equal to

$$\{\theta \in \mathbb{R}^p : a_i^T(\theta - \theta_0) = 0 \text{ for } i = 1, \dots, \ell; a_i^T(\theta - \theta_0) \geq 0 \text{ for } i \in J(\theta_0)\}.$$

The foregoing result shows that if MF-CQ is satisfied then the approximating cone at θ_0 is obtained by substituting the first order approximation (i.e., linear approximation at θ_0), $\{h_i(\theta_0) + (\theta - \theta_0)^T(\partial/\partial\theta)h_i(\theta_0)\}$, of $h_i(\theta)$ at θ_0 if $h_i(\theta_0) = 0$; if $h_s(\theta_0) > 0$ (i.e., $h_s(\theta_0) \geq 0$ is inactive) then the constraint $h_s(\theta_0) \geq 0$ does not play a part in determining the approximating cone at θ_0 ($s = \ell + 1, \dots, k$).

The approximating cone obtained by applying the result in the previous proposition is a polyhedral. Thus, if the approximating cone is not a polyhedral then this method of obtaining the approximating cone would not be applicable in its current form. For example, in Fig. 4.1 if $\Theta = A \cup B$ where A is the region bounded by $\theta_2 = g_2(\theta_1)$ and $\theta_2 = \theta_1$, and B is the region bounded by $\theta_2 = g_1(\theta_1)$ and $\theta_2 = -\theta_1$, then the approximating cone of Θ at O is not convex. Therefore, for this Θ , the method in Proposition 4.7.3 is not directly applicable; in any case, this would be obvious because Θ is not of the form required by Proposition 4.7.3. It can be seen that A and B are

of the form required by Proposition 4.7.3 and $\mathcal{A}(A \cup B; \mathbf{0}) = \mathcal{A}(A; \mathbf{0}) \cup \mathcal{A}(B; \mathbf{0})$; now the method based on MF-CQ can be applied to A and B separately to obtain $\mathcal{A}(A \cup B; \mathbf{0})$. More generally, the foregoing method can be applied to a parameter space that is equal to the union of sets each of which satisfies the conditions of Proposition 4.7.3.

Example 4.7.4 Tangent Cones

(1) Let $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 \geq (1/2)\theta_1, \theta_2 \leq 2\theta_1\}$ (see Fig. 4.7). Then, we have the following:

$$\begin{aligned}\mathcal{T}(\Theta_0; \theta_0) &= \Theta, \text{ where } \theta_0 = (0, 0)^T. \\ \mathcal{T}(\Theta; \theta_a) &= \{x \in \mathbb{R}^2 : x_2 \geq (1/2)x_1\}, \text{ where } \theta_a = (2, 1)^T. \\ \mathcal{T}(\Theta; \theta_b) &= \mathbb{R}^2, \text{ where } \theta_b = (2, 2)^T. \\ \mathcal{T}(\Theta; \theta_c) &= \{x \in \mathbb{R}^2 : x_2 \leq 2x_1\}, \text{ where } \theta_c = (1, 2)^T.\end{aligned}$$

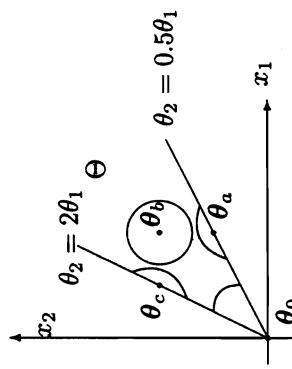


Fig. 4.7 The tangent cones of $\Theta = \{\theta : \theta_2 \leq 2\theta_1, \theta_2 \geq 0.5\theta_1\}$ at $\theta_0, \theta_a, \theta_b$ and θ_c .

Note that at θ_a , the inequality $\theta_2 \leq 2\theta_1$ holds strictly (i.e., it is not active) and therefore this inequality constraint can be ignored for determining the tangent cone at θ_a . At θ_b , both inequalities are satisfied strictly, and therefore both constraints can be ignored for determining the tangent cone. Since there are no other active constraints, the tangent cone is the full space \mathbb{R}^2 . At θ_c , the inequality $\theta_2 \geq (1/2)\theta_1$ holds strictly and hence it can be ignored for determining the tangent cone; the rest of the arguments are similar to those for θ_a . ■

(2) Let the matrices R_1 and R_2 be given and let $\Theta = \{\theta \in \mathbb{R}^p : R_1\theta \geq 0, R_2\theta \geq 0\}$. Then, by Proposition 4.7.3 we have the following:

$$\begin{aligned}R_1\theta_a > 0, R_2\theta_a = 0 &\Rightarrow \mathcal{T}(\Theta; \theta_a) = \{x \in \mathbb{R}^p : R_2x \geq 0\}. \\ R_1\theta_0 = 0, R_2\theta_0 = 0 &\Rightarrow \mathcal{T}(\Theta; \theta_0) = \{x \in \mathbb{R}^p : R_1x \geq 0, R_2x \geq 0\}. \\ R_1\theta_b > 0, R_2\theta_b > 0 &\Rightarrow \mathcal{T}(\Theta; \theta_b) = \mathbb{R}^p.\end{aligned}\quad (4.98)$$

(3) Let $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 \geq -\theta_1^2 + 2\theta_1, \theta_2 \geq -\theta_1^2 - 2\theta_1\}$ and $\theta_0 = \mathbf{0}$.

It is easily verified that the MF-CQ is satisfied at θ_0 and therefore the tangent cone is determined by linearizing the constraints at θ_0 as follows. Let the two constraints

be expressed as $g_1(\theta) \geq 0$ and $g_2(\theta) \geq 0$. Then the linearized versions of these constraints at θ_0 are

$g_1(\theta_0) + (\theta - \theta_0)^T (\partial/\partial\theta) g_1(\theta_0) \geq 0$ and $g_2(\theta_0) + (\theta - \theta_0)^T (\partial/\partial\theta) g_2(\theta_0) \geq 0$, respectively. These linearized versions are precisely $\theta_2 \geq 2\theta_1$ and $\theta_2 \geq -2\theta_1$, respectively. Since the two constraints defining Θ are active at θ_0 , we have $T(\Theta; \theta_0) = \{\theta : \theta_2 \geq 2\theta_1, \theta_2 \geq -2\theta_1\}$. In this particular example, we can sketch the shape of Θ easily and the tangent cone can be written down quickly; the above analytical approach is meant to be instructive.

(4) Let $\Theta = \{\theta \in \mathbb{R}^2 : \theta_2 \geq \theta_1^2, \theta_2 \leq \theta_1\}$. Let $g_1(\theta) = \theta_2 - \theta_1^2$ and $g_2(\theta) = \theta_1 - \theta_2$. Then $\nabla g_1(\theta) = (-2\theta_1, 1)^T$, $\nabla g_2(\theta) = (1, -1)^T$. Now, applying the same technique as in the previous example, we have $T(\Theta; \theta_0) = \{\theta \in \mathbb{R}^2 : \theta_1 \geq \theta_2 \geq 0\}$, and Θ is Chernoff regular at 0 (see Fig. 4.8). ■

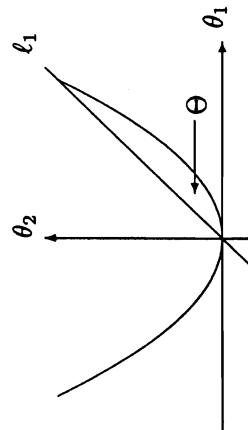


Fig. 4.8 The tangent cone of $\Theta = \{\theta_2 \geq \theta_1^2, \theta_2 \leq \theta_1\}$ at the origin is the cone between ℓ_1 and the θ_1 axis.

The next result is useful in asymptotic derivations because it enables us to substitute the approximating cone for the parameter space at the true value even if it (i.e., the true parameter value) is on the boundary which may or may not be smooth.

Proposition 4.7.4 Let $\Theta \subset \mathbb{R}^p$, V be a positive definite matrix of order $p \times p$ and $\theta_0 \in \Theta$ and A be an approximating cone of Θ at θ_0 . Then

$$\|\mathbf{y} - \Theta\|_V^2 - \|\mathbf{y} - A(\Theta; \theta_0)\|_V^2 = o(\|\mathbf{y} - \theta_0\|_V^2) \text{ as } \mathbf{y} \rightarrow \theta_0.$$

This completes the proof of the lemma.

Now, the proof of the first part of the corollary follows from the foregoing Lemma and Proposition 4.7.4 by substituting $g(\mathbf{y}) = \|\mathbf{y} - \Theta\|_V^2 - \|\mathbf{y} - A(\Theta; \theta_0)\|_V^2$, $\beta = 2$ and $\gamma = 1/2$.

To prove the second part, let $Z_n = n^{1/2}(T_n - \theta_0)$ and note that

$$\begin{aligned} n\|\mathbf{T}_n - \Theta\|_V^2 &= n\|\mathbf{T}_n - A(\Theta; \theta_0)\|_V^2 + o_p(n^{-1}) \\ &= \|Z_n - T(\Theta; \theta_0)\|_V^2 + o_p(1) \xrightarrow{d} \|\mathbf{Z} - T(\Theta; \theta_0)\|_V^2 \end{aligned}$$

A useful corollary to the foregoing proposition is the following.

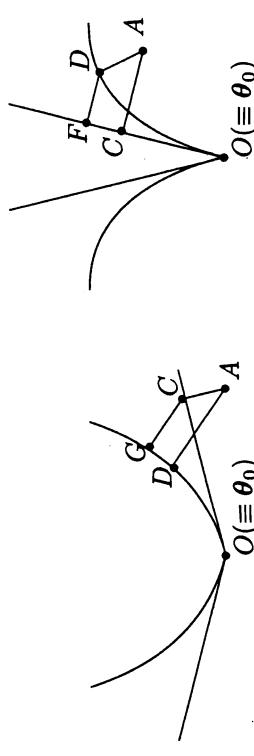


Fig. 4.9 $\|\mathbf{y} - \Theta\|^2 - \|\mathbf{y} - A(\Theta; \theta_0)\|^2 = o(\|\mathbf{y} - \theta_0\|^2)$ [i.e., $|AD^2 - AC^2| = o(OA^2)$]

Corollary 4.7.5 Let $\Theta \subset \mathbb{R}^p$, $\theta_0 \in \Theta$, V be a $p \times p$ positive definite matrix and Θ be Chernoff regular at θ_0 . Then we have the following:

1. If $n^{1/2}(T_n - \theta_0) = O_p(1)$ then

$$\|\mathbf{T}_n - \Theta\|_V^2 - \|\mathbf{T}_n - A(\Theta; \theta_0)\|_V^2 = o_p(n^{-1}).$$

2. If $n^{1/2}(T_n - \theta_0) \xrightarrow{d} \mathcal{Z}$ then $n\|\mathbf{T}_n - \Theta\|_V^2 \xrightarrow{d} \|\mathbf{Z} - T(\Theta; \theta_0)\|_V^2$

Proof: Let us first state a lemma.

Lemma 4.7.6 Let $g : \mathbb{R}^p \rightarrow \mathbb{R}$ be such that $g(\mathbf{y}) = o(\|\mathbf{y} - \theta_0\|^\beta)$ as $\|\mathbf{y} - \theta_0\| \rightarrow 0$, where $\beta > 0$. Let T_n be a sequence of random p -vectors such that $T_n - \theta_0 = O_p(n^{-\gamma})$ where $\gamma > 0$. Then $g(T_n) = o_p(n^{-\beta\gamma})$.

Proof: Given $\epsilon > 0$, there exists a δ such that $|g(\mathbf{y})|/\|\mathbf{y} - \theta_0\|^\beta < \epsilon$ for $\|\mathbf{y} - \theta_0\| < \delta$. Therefore,

$$\text{pr}\{|g(T_n)|/\|T_n - \theta_0\|^\beta < \epsilon\} \geq \text{pr}(\|T_n - \theta_0\| < \delta) \rightarrow 1 \text{ as } n \rightarrow \infty.$$

Thus, $g(T_n)/\|T_n - \theta_0\|^\beta = o_p(1)$. Now,

$$g(T_n) = \|T_n - \theta_0\|^\beta o_p(1) = |n^{-\gamma} O_p(1)|^\beta o_p(1) = o_p(n^{-\beta\gamma}).$$

This completes the proof of the lemma.

Now, the proof of the first part of the corollary follows from the foregoing Lemma and Proposition 4.7.4 by substituting $g(\mathbf{y}) = \|\mathbf{y} - \Theta\|_V^2 - \|\mathbf{y} - A(\Theta; \theta_0)\|_V^2$, $\beta = 2$ and $\gamma = 1/2$.

To prove the second part, let $Z_n = n^{1/2}(T_n - \theta_0)$ and note that

The corollary essentially says that $\mathcal{A}(\Theta; \theta_0)$ can be substituted for Θ in most of the expressions; in particular, if T_n is a \sqrt{n} -consistent estimator, θ_0 is the true value and $\theta_0 \in \Theta$ then the squared distance between T_n and Θ is equal to that between T_n and the approximating cone of Θ at θ_0 except for a $o_p(n^{-1})$ term. The corollary can also be extended to the case when Θ_n is a sequence of sets such that $\theta_0 \in \Theta_n$ and $\mathcal{A}(\Theta_n; \theta_0)$ converges to a cone, or Θ_n converges to a cone.

4.8 GENERAL TESTING PROBLEMS

Let $f(y; \theta)$ be a density function where $\theta \in \Theta \subset \mathbb{R}^p$ and $\Theta_0 \subset \Theta_1 \subset \Theta$. Let Y be a random variable with density function $f(y; \theta)$. In this section, test of

$$H_0 : \theta \in \Theta_0 \text{ against } H_1 : \theta \in \Theta_1 \quad (4.99)$$

is studied when Θ_0, Θ_1 and Θ may have boundary points and the boundaries of these sets may or may not be smooth; for example, the boundary of a cone at its vertex is not smooth. The special case when Θ_0 is a linear space and Θ_1 is a polyhedral was studied in Section 4.3, where the special structure on the parameter spaces enabled us to obtain explicit expressions for the asymptotic null distributions of the test statistics.

Now, we consider more general cases and obtain general results. It will be assumed that Θ_0 and Θ_1 have approximating cones at every point in Θ_0 ; in the terminology introduced in Section 4.7, this is same as saying that Θ_0 and Θ_1 are Chernoff regular at every $\theta \in \Theta_0$. For simplicity, it is assumed that the observations on Y are independently and identically distributed; however, the main results of this section hold under appropriate regularity conditions, for much more general stochastic processes although such generalizations are not considered.

Let Y_1, \dots, Y_n be independently and identically distributed as Y . As in the earlier sections, let

$$\ell(\theta) = \sum f(Y_i; \theta), \quad S(\theta) = (\partial/\partial\theta)\ell(\theta), \quad (4.100)$$

$$\text{and } \mathcal{I}(\theta) = E\{(\partial/\partial\theta)\log f(Y; \theta)(\partial/\partial\theta^T)\log f(Y; \theta)\} \quad (4.101)$$

denote the loglikelihood, the score function, and the information matrix, respectively. Some important fundamental results on the *LRT* for testing problems of the form (4.99) were obtained by Chernoff (1954). The technical details in Chernoff (1954) assumed that Θ was open and that inference was based on likelihood. Chernoff's (1954) results have natural extensions to more general cases although modifications would be required to the regularity conditions and technical details; for example, see Self and Liang (1987), Andrews (1998), Vu and Zhou (1997), Shapiro (1985, 1989, 2000a), Geyer (1994), Silvapulle (1992a, 1992b, 1992c, 1994), and El Barmi (1996) among others. The areas covered include the cases when (i) Θ is not necessarily open (ii) observations are not necessarily independently and identically distributed, and (iii) inference is based on quasi-likelihood, partial likelihood, or general objective functions such as those used in, for example, M -estimation, minimum distance methods, and empirical likelihood.

In the next subsection, we consider the simple case of likelihood ratio test when Θ is open. The simplicity of this case makes it easier to convey the essentials. On the other hand, this case is also of independent interest because many practical problems fall into this category. Then, in subsection 4.8.2, we consider likelihood inference when Θ is not necessarily open and the true value of θ can be a boundary point of Θ , Θ_0 , and Θ_1 simultaneously. The distribution and implementation of likelihood ratio test are virtually the same whether or not θ_0 is a boundary point of Θ provided some conditions are satisfied; essentially all that we require is that it should be possible to approximate the objective function by a quadratic similar to that in (4.2) and (4.3). In this sense, the results in subsection 4.8.1 are not very different from those in subsection 4.8.2. In subsection 4.8.3, we consider inference based on general objective functions rather than just likelihood. Some results obtained in subsections 4.8.1 and 4.8.2 do not necessarily remain the same in this general setting. In subsection 4.8.5, some simple examples in two dimensions are considered to illustrate the theory. In subsection 4.8.5, a general form of the Type B testing problem, test of $\mathbf{h}(\theta) \geq 0$ against $\mathbf{h}(\theta) \not\geq 0$, is studied; this particular type of problems arises in economics, for example, the interest may be to test whether or not a regression function deviates from concavity. As usual, we adopt the notation $\hat{\theta}, \bar{\theta}$, and $\tilde{\theta}$ for the estimators over Θ, Θ_1 , and Θ_0 respectively.

4.8.1 Likelihood Approach When Θ is Open

Assume that Θ is open and that $H_0 : \theta \in \Theta_0$ holds. Let θ_0 denote the true value of θ in Θ_0 . The set of regularity conditions required for this setting are virtually the same as those usually required for establishing that $(\hat{\theta} - \theta_0)$ is asymptotically normal and the asymptotic null distribution of the *LRT* of $R\theta = 0$ against $R\theta \neq 0$ is chi-squared. There are different sets of conditions available for these results; some involve third-order partial derivatives of the log density, some involve only the second-order partial derivatives, and some replace the differentiability by other special types of differentiability (quadratic mean, stochastic, etc.). For simplicity, we shall work with conditions on the third-order derivatives of the log density, but we recognize that our technical conditions can be improved. Let

$$\begin{aligned} \ell(\theta) &= \sum f(Y_i; \theta), & S(\theta) &= (\partial/\partial\theta)\ell(\theta), \\ \text{and } \mathcal{I}(\theta) &= E\{(\partial/\partial\theta)\log f(Y; \theta)(\partial/\partial\theta^T)\log f(Y; \theta)\} \end{aligned} \quad (4.102)$$

$$W_D = \min_{\theta \in \Theta_0} n(\hat{\theta} - \theta)^T \hat{\mathcal{I}}(\hat{\theta} - \theta) - \min_{\theta \in \Theta_1} n(\hat{\theta} - \theta)^T \hat{\mathcal{I}}(\hat{\theta} - \theta), \quad (4.103)$$

Define,

where $\hat{\mathcal{I}} = \mathcal{I}(\hat{\theta})$; however, in general $\hat{\mathcal{I}}$ can be any consistent estimator of $\mathcal{I}(\theta_0)$. Note that W_D is a measure of $\{dist(\hat{\theta}, H_0) - dist(\hat{\theta}, H_1)\}$, where $dist(\hat{\theta}, H)$ is a measure of distance between $\hat{\theta}$ and the parameter space defined by H . Hence it is a suitable statistic for testing $H_0 : \theta \in \Theta_0$ against $H_1 : \theta \in \Theta_1$. As in the case of test against the unrestricted alternative, *LRT* and W_D are asymptotically equivalent. A simple method of obtaining their common asymptotic distribution is given in the

next result; essentially, it says that the value of LRT changes by a term of order $o_p(1)$ if $\mathcal{A}(\Theta_i; \theta_0)$ is substituted for Θ_i ($i = 0, 1$).

Proposition 4.8.1 *Assume that Θ is open, Condition Q is satisfied, and $H_0 : \theta \in \Theta_0$ holds. Then, with the notation in (4.102), we have the following:*

1. $LRT = W_D + o_p(1)$, where W_D is given in (4.103).
2. The common asymptotic null distribution of LRT and W_D is equal to the distribution of

$$\|\mathbf{X} - \mathcal{A}(\Theta_0; \theta_0)\|^2 - \|\mathbf{X} - \mathcal{A}(\Theta_1; \theta_0)\|^2, \quad (4.104)$$

which is also equal to the distribution of

$$\|\mathbf{Z} - \mathcal{T}(\Theta_0; \theta_0)\|^2 - \|\mathbf{Z} - \mathcal{T}(\Theta_1; \theta_0)\|^2 \quad (4.105)$$

where $\mathbf{X} \sim N(\theta_0, \mathcal{I}(\theta_0)^{-1})$, $\mathbf{Z} \sim N(\mathbf{0}, \mathcal{I}(\theta_0)^{-1})$ and $\|\cdot\|$ is as in (4.102).

Proof: Let $\mathcal{A}_i = \mathcal{A}(\Theta_i; \theta_0)$ and $\mathcal{T}_i = \mathcal{T}(\Theta_i; \theta_0)$ for $i = 0, 1$. Now, by using the quadratic approximation (4.4), we have that

$$\begin{aligned} LRT &= 2[\sup_{\theta \in \Theta_1} \ell(\theta) - \sup_{\theta \in \Theta_0} \ell(\theta)] \\ &= n\|\hat{\theta} - \Theta_0\|^2 - n\|\hat{\theta} - \Theta_1\|^2 + o_p(1) \\ &= W_D + o_p(1); \end{aligned} \quad (4.106)$$

the last step, which says that $\|\cdot\|_{\mathcal{T}(\theta_0)}$ can be replaced by $\|\cdot\|_{\mathcal{Z}}$, follows from Lemma 4.10.2 in the Appendix.

Now, let us prove the second part. By Corollary 4.7.5, Θ_i in (4.106) can be replaced by its approximating cone ($i = 0, 1$). Let $\mathbf{Z}_n = \sqrt{n}(\hat{\theta} - \theta_0)$. Now, by applying essentially the same arguments as those used in the proof of the first part, we have

$$\begin{aligned} LRT &= n\{\|\hat{\theta} - \mathcal{A}_0\|^2 - \|\hat{\theta} - \mathcal{A}_1\|^2\} + o_p(1), \\ &= \|\mathbf{Z}_n - \mathcal{T}_0\|^2 - \|\mathbf{Z}_n - \mathcal{T}_1\|^2 + o_p(1), \\ &\xrightarrow{d} \|\mathbf{Z} - \mathcal{T}_0\|^2 - \|\mathbf{Z} - \mathcal{T}_1\|^2; \end{aligned} \quad (4.107)$$

the last step follows since $\|\mathbf{z} - \mathcal{T}_i\|^2$ is continuous in \mathbf{z} and $\mathbf{Z}_n \xrightarrow{d} \mathbf{Z}$. ■

The foregoing results show that, if the null hypothesis is true and θ_0 denotes the true value, then the asymptotic distribution of the LRT of $\theta \in \Theta_0$ against $\theta \in \Theta_1$ is equal to the distribution of

1. The LRT of $\theta \in \mathcal{A}(\Theta_0; \theta_0)$ against $\theta \in \mathcal{A}(\Theta_1; \theta_0)$ based on a single observation of \mathbf{X} , where $\mathbf{X} \sim N\{\theta, \mathcal{I}(\theta)^{-1}\}$ and the true value of θ is θ_0 .
2. The LRT of $\alpha \in \mathcal{T}(\Theta_0; \theta_0)$ against $\alpha \in \mathcal{T}(\Theta_1; \theta_0)$ based on a single observation of \mathbf{Z} where $\mathbf{Z} \sim N\{\alpha, \mathcal{I}(\alpha)^{-1}\}$ and the true value of α is 0 .

Thus, as in classical inference and in Section 4.3, the asymptotic distribution of the test statistics can be expressed in terms of a single observation from the normal distribution. However, this apparently simpler problem may be still quite complicated, as was noted in Section 4.3, mainly due to the presence of the nuisance parameter in the null distribution; the nuisance parameter appears because the asymptotic null distribution of the LRT depends on $\mathcal{I}(\theta_0)$ where θ_0 is the assumed true value in the null parameter space. This is in sharp contrast to the case where the asymptotic null distribution of LRT for testing $R\theta = 0$ against $R\theta \neq 0$ is a chi-squared and hence does not depend on θ_0 in the null parameter space.

If $\mathcal{T}(\Theta_0; \theta_0)$ is a linear space and $\mathcal{T}(\Theta_1; \theta_0)$ is a closed convex cone, then it follows from the foregoing proposition that the asymptotic null distribution of the LRT is a chi-bar-square. For example, if the null and alternative hypotheses are $R\theta = 0$ and $R\theta \geq 0$, respectively, then the asymptotic null distribution of the LRT is chi-bar-square; this was also obtained in Proposition 4.3.1. In fact, Proposition 4.3.1 on the asymptotic null distribution of LRT is a special case of the foregoing proposition. However, in general, the asymptotic null distribution of the LRT is not a chi-bar-square.

The results in the foregoing proposition provide elegant asymptotic representations of the LRT and they lead to its asymptotic distribution, which has a simpler structure because the original parameter spaces have been replaced by their approximating cones and the limiting distribution is a function of Z where $Z \sim N(\mathbf{0}, \mathcal{I}_{\theta_0}^{-1})$. Although Proposition 4.8.1 was established for the case when Θ is open, similar results are available even when Θ is not open and θ_0 is a boundary point of Θ .

4.8.2 Likelihood Approach When Θ Is not Open

Now assume that Θ is not necessarily open and that θ_0 may not be an interior point of Θ . Then $(\partial/\partial\theta)\ell(\hat{\theta})$ may not be zero and $n^{1/2}(\hat{\theta} - \theta_0)$ may not be asymptotically normal. Consequently, the arguments in the previous subsection that depend on the asymptotic normality of $n^{1/2}(\hat{\theta} - \theta_0)$ are not valid; also, arguments based on Taylor series may not be valid. However, most of the results still hold. For example, the asymptotic null distribution of LRT remains unchanged even if θ_0 is a boundary point Θ , but the technical details leading to it are different from those in the previous subsection. *Assume that the null hypothesis holds.* In most practical situations, $\ell(\theta)$ would have directional derivatives at θ_0 , and hence would admit a quadratic approximation in a neighborhood of θ_0 . Further, the maximizers of $\ell(\theta)$ over Θ_0 and over Θ_1 are likely to be $n^{1/2}$ -consistent. This is because consistency of the estimator is not affected by the true parameter being on the boundary, and \sqrt{n} -consistency follows once we have a quadratic approximation of $\ell(\theta)$ similar to (4.2) by arguments indicated at the beginning of this chapter. If the objective function $\ell(\theta)$ does not have a quadratic approximation similar to (4.2) or the estimators are not \sqrt{n} -consistent, then the general approach of this section will need modifications to obtain the asymptotic distribution of LRT .

The Condition Q2 defined below is a slight modification of Condition Q; the modification is made in view of the fact that θ_0 may be a boundary point of Θ .

Condition Q2. There exist $S_*(\theta)$ and $\mathcal{I}_*(\theta)$ such that (4.3) holds with $S(\theta)$ and $\mathcal{I}(\theta)$ replaced by $S_*(\theta)$ and $\mathcal{I}_*(\theta)$ respectively. Further,

$$1. n^{-1/2} S_*(\theta_0) \xrightarrow{d} N\{\mathbf{0}, \mathcal{I}_*(\theta_0)\}.$$

$$2. \text{The maximizers of } \ell(\theta) \text{ over } \Theta_0 \text{ and over } \Theta_1 \text{ are } \sqrt{n}\text{-consistent for } \theta_0. \blacksquare$$

There is room for fine tuning this condition; for example, \sqrt{n} -consistency would follow from (4.3) and consistency. However, it is convenient for us to state them explicitly.

If θ_0 is an interior point of Θ , then S_* and \mathcal{I}_* would be the S and \mathcal{I} in the previous subsection. If Θ is a rectangle in the Euclidean space, θ_0 is a boundary point of Θ , and $\ell(\theta)$ has partial derivatives at θ_0 from the appropriate directions then Condition Q2 could be established using Taylor series expansions (for example, see Self and Liang (1987) and Andrews (1999)).

If Θ does not have the shape of a rectangle and θ_0 is a boundary point of Θ , then it is likely that Condition Q2 can be established using directional derivatives (see Shapiro (2000b)). Directional derivatives of $\ell(\theta)$ at θ_0 would be required only in directions that take us into the approximating cone of Θ at θ_0 . The proof of the next result is essentially the same as that for the previous one, hence omitted.

Proposition 4.8.2 Assume that Condition Q2 is satisfied and that $H_0 : \theta \in \Theta_0$ holds. Let θ_0 denote the true value of θ in Θ_0 , $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathcal{I}_*(\theta_0) \mathbf{x}$, $\mathbf{Z} \sim N\{\mathbf{0}, \mathcal{I}_*(\theta_0)^{-1}\}$, $\mathbf{X} = \theta_0 + \mathbf{Z}$, and $\mathbf{U}_n = n^{-1/2} \mathcal{I}_*(\theta_0)^{-1} \mathbf{S}_*(\theta_0)$. Then,

$$1. LRT = \|\mathbf{U}_n - \mathcal{T}(\theta_0; \theta_0)\|^2 - \|\mathbf{U}_n - \mathcal{T}(\Theta_1; \theta_0)\|^2 + o_p(1).$$

2. The asymptotic null distribution of LRT is equal to the distribution of

$$\|\mathbf{X} - \mathcal{A}(\Theta_0; \theta_0)\|^2 - \|\mathbf{X} - \mathcal{A}(\Theta_1; \theta_0)\|^2, \quad (4.108)$$

which is also equal to the distribution of

$$\|\mathbf{Z} - \mathcal{T}(\Theta_0; \theta_0)\|^2 - \|\mathbf{Z} - \mathcal{T}(\Theta_1; \theta_0)\|^2. \quad (4.109)$$

Now, a likelihood ratio type statistic for testing $\theta \in \Theta_0$ against $\theta \in \Theta_1$ is

$$T_n = 2\left\{\sup_{\theta \in \Theta_1} R_n(\theta) - \sup_{\theta \in \Theta_0} R_n(\theta)\right\}. \quad (4.110)$$

The next proposition shows that the asymptotic null distribution of T_n is obtained by an argument similar to those in the previous subsection.

Proposition 4.8.3 Assume that Condition Q3 is satisfied and that $H_0 : \theta \in \Theta_0$ holds. Let θ_0 denote the true value of θ in Θ_0 , $\|\mathbf{x}\|^2 = \mathbf{x}^T G_*(\theta_0) \mathbf{x}$, and $\mathbf{U}_n = n^{-1/2} G_*(\theta_0)^{-1} \mathbf{S}_*(\theta_0)$. Then

$$1. T_n = \|\mathbf{U}_n - \mathcal{T}(\Theta_0; \theta_0)\|^2 - \|\mathbf{U}_n - \mathcal{T}(\Theta_1; \theta_0)\|^2 + o_p(1).$$

2. The asymptotic null distribution of T_n is equal to that of

$$\|\mathbf{U} - \mathcal{T}(\Theta_0; \theta_0)\|^2 - \|\mathbf{U} - \mathcal{T}(\Theta_1; \theta_0)\|^2 \quad (4.111)$$

where $U \sim N\{0, G_*(\theta_0)^{-1}\mathcal{I}_*(\theta_0)G_*(\theta_0)^{-1}\}$.

Proof: It follows from Condition Q3(1) that, for $\theta - \theta_0 = O(n^{-1/2})$,

$$R_n(\theta) = R_n(\theta_0) - (n/2)\|n^{1/2}U_n - (\theta - \theta_0)\|^2 + (1/2)n^{-1}\|U(\theta_0)\|^2 + o_p(1).$$

Now,

$$\begin{aligned} T_n &= 2[\sup\{R_n(\theta) : \theta \in \Theta_1\} - \sup\{R_n(\theta) : \theta \in \Theta_0\}] \\ &= n\|n^{-1/2}U_n - (\theta_0 - \theta_0)\|^2 - n\|n^{-1/2}U_n - (\theta_1 - \theta_0)\|^2 + o_p(1) \\ &= \|U_n - \mathcal{T}(\theta_0; \theta_0)\|^2 - \|U_n - \mathcal{T}(\theta_1; \theta_0)\|^2 \\ &\xrightarrow{d} \|U - \mathcal{T}(\theta_0; \theta_0)\|^2 - \|U - \mathcal{T}(\theta_1; \theta_0)\|^2. \end{aligned}$$

Note that the quadratic form $\|U - \alpha\|^2$ in (4.111), which is equal to $(U - \alpha)^T G_*(\theta_0)(U - \alpha)$, is with respect to the matrix $G_*(\theta_0)^{-1}$ but the covariance matrix of U is $G_*(\theta_0)^{-1}\mathcal{I}_*(\theta_0)G_*(\theta_0)^{-1}$. Therefore, if $\mathcal{I}_*(\theta_0) \neq G_*(\theta_0)$, then the quadratic form $\|U - \alpha\|^2$ may not be with respect to the covariance matrix of U and hence the results in Chapter 2 for normal distributions are not applicable for obtaining the asymptotic distribution of T_n . In particular, (i) for testing $\theta = 0$ against $\theta \neq 0$, the asymptotic null distribution of T_n may not be a chi-square, and (ii) for testing $R\theta = 0$ against $R\theta \geq 0$, the asymptotic null distribution of T_n in (4.111) may not be a chi-bar-square.

Suppose that Θ is open and that $H_0 : \theta \in \Theta_0$ holds. As usual, let $\hat{\theta}$ denote the unrestricted estimator of θ . Then by standard arguments, $n^{1/2}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{0, V(\theta_0)\}$, where $V(\theta_0) = G_*(\theta_0)^{-1}\mathcal{I}_*(\theta_0)G_*(\theta_0)^{-1}$. Therefore, a W_D -type statistic for testing $\theta \in \Theta_0$ against $\theta \in \Theta_1$ is

$$W_D = \min_{\theta \in \Theta_0} n(\hat{\theta} - \theta)^T \hat{V}^{-1}(\hat{\theta} - \theta) - \min_{\theta \in \Theta_1} n(\hat{\theta} - \theta)^T \hat{V}^{-1}(\hat{\theta} - \theta), \quad (4.112)$$

where \hat{V} is a consistent estimator of $V(\theta_0)$, for example, $V(\hat{\theta})$. The asymptotic distribution of the foregoing W_D is equal to that of

$$\min_{\theta \in \mathcal{T}(\theta_0; \theta_0)} (\mathbf{Z} - \theta)^T \mathbf{V}(\theta_0)^{-1}(\mathbf{Z} - \theta) - \min_{\theta \in \mathcal{T}(\theta_1; \theta_0)} (\mathbf{Z} - \theta)^T \mathbf{V}(\theta_0)^{-1}(\mathbf{Z} - \theta)$$

where $\mathbf{Z} \sim N\{0, V(\theta_0)\}$.

Note that the quadratic form $(\mathbf{Z} - \theta)^T \mathbf{V}(\theta_0)^{-1}(\mathbf{Z} - \theta)$, which appears in the asymptotic null distribution of W_D is with respect to the covariance matrix $V(\theta_0)$ of \mathbf{Z} . Consequently, it follows that for testing $R\theta = \mathbf{0}$ against $R_1\theta \geq 0$, the asymptotic null distribution of W_D is a chi-bar-square. Therefore, the asymptotic null distribution of W_D in (4.112) and that of the likelihood ratio type statistic T_n in (4.110) are not necessarily the same, and hence are unlikely to be asymptotically equivalent in this type of setting.

For some partial-likelihood and quasi-likelihood, $G_*(\theta)$ and $\mathcal{I}_*(\theta)$ are equal except for a scalar multiplicative constant. Since the multiplicative constant factors out from matrix products, it (i.e., multiplicative constant) can be replaced by a consistent estimator of it and the required results can be obtained (see quasi-likelihood method on page 159 and Silvapulle (1994)).

4.8.4 Examples

Example 4.8.1 Asymptotic distribution of the LRT of $g(\theta) = 0$ vs $g(\theta) \geq 0$.

Let $\mathbf{Y} = (Y_1, Y_2)^T \sim N(\boldsymbol{\theta}, \mathbf{V})$ where $\boldsymbol{\theta} \in \Theta = \mathbb{R}^2$ and \mathbf{V} is known and does not depend on $\boldsymbol{\theta}$. Let the null and alternative hypotheses be

$$H_0 : g(\boldsymbol{\theta}) = 0 \text{ and } H_1 : g(\boldsymbol{\theta}) \geq 0,$$

respectively, where $\mathbf{g} = (g_1, g_2)^T$, $g_1(\boldsymbol{\theta}) = \theta_2 - \theta_1^2 - \theta_1$ and $g_2(\boldsymbol{\theta}) = \theta_2 - \theta_1^2 + \theta_1$; see Fig. 4.1. Let us obtain the asymptotic null distribution of the LRT given by Proposition 4.8.1. Let $\bar{\mathbf{Y}}$ denote the mean of n independently and identically distributed observations on \mathbf{Y} ; then $\bar{\mathbf{Y}} \sim N(\boldsymbol{\theta}, n^{-1}\mathbf{V})$. Now, $\Theta_0 = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \mathbf{g}(\boldsymbol{\theta}) = \mathbf{0}\}$ and $\Theta_1 = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \mathbf{g}(\boldsymbol{\theta}) \geq 0\}$, Θ_1 is shown as Ω in Fig. 4.1 on page 184. Clearly, $\Theta_0 = \{\mathbf{0}\}$. Assume that H_0 is true; then the true value $\boldsymbol{\theta}_0$ is $\mathbf{0}$, $\mathcal{T}(\boldsymbol{\theta}_0; \boldsymbol{\theta}_0) = \{\mathbf{0}\}$, and $\mathcal{T}(\boldsymbol{\theta}_1; \boldsymbol{\theta}_0) = \{\boldsymbol{\theta} : \theta_2 \geq \theta_1, \theta_2 \geq -\theta_1\} = \{\boldsymbol{\theta} : R\boldsymbol{\theta} \geq 0\}$ where R is the 2×2 matrix $[-1, 1][1, 1]$; $\mathcal{T}(\boldsymbol{\theta}_1; \mathbf{0})$ is shown as \mathcal{T} in Fig. 4.1. Let $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{V})$. Now, since $\ell(\boldsymbol{\theta}) = (-1/2)n(\bar{\mathbf{Y}} - \boldsymbol{\theta})^T \mathbf{V}^{-1}(\bar{\mathbf{Y}} - \boldsymbol{\theta}) + const$ and $\mathcal{I}(\boldsymbol{\theta}) = \mathbf{V}^{-1}$, we have the following:

$$\begin{aligned} LRT &= n\|\bar{\mathbf{Y}}\|_{\mathbf{V}}^2 - n\|\bar{\mathbf{Y}} - \boldsymbol{\theta}_1\|_{\mathbf{V}}^2 \\ &= n\|\bar{\mathbf{Y}}\|_{\mathbf{V}}^2 - n\|\bar{\mathbf{Y}} - \mathcal{A}(\boldsymbol{\theta}_1; \mathbf{0})\|_{\mathbf{V}}^2 + o_p(1) \\ &\xrightarrow{d} \|\mathbf{Z}\|_{\mathbf{V}}^2 - \|\mathbf{Z} - \mathcal{T}(\boldsymbol{\theta}_1; \mathbf{0})\|_{\mathbf{V}}^2 \sim \chi^2\{\mathbf{V}, \mathcal{T}(\boldsymbol{\theta}_1; \mathbf{0})\}. \end{aligned}$$

Therefore, by (3.8)

$$\begin{aligned} \text{pr}(LRT \geq c | H_0) &\rightarrow 0.5\text{pr}(\chi_1^2 \geq c) + (0.5 - q)\text{pr}(\chi_2^2 \geq c) \quad c > 0 \\ &\text{where } (0.5 - q) \text{ is as in (3.8).} \end{aligned}$$

Example 4.8.2 Asymptotic null distribution of the LRT of $g(\theta) \geq 0$ vs $g(\theta) \geq 0$.

Let $\mathbf{Y}, \bar{\mathbf{Y}}, \mathbf{V}$, and g be as in the previous example. Let us obtain the asymptotic null distribution of the LRT for testing

$$H_0 : g(\boldsymbol{\theta}) \geq 0 \text{ against } H_1 : g(\boldsymbol{\theta}) \not\geq 0.$$

The relevant asymptotic result is given in Proposition 4.8.1. By definition, $\Theta_0 = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \mathbf{g}(\boldsymbol{\theta}) \geq 0\}$ and $\Theta_1 = \mathbb{R}^2$. Assume that the null hypothesis holds and let $\boldsymbol{\theta}_0$ denote the true value of $\boldsymbol{\theta}$. Clearly, $\mathcal{T}(\boldsymbol{\theta}_1; \boldsymbol{\theta}_0) = \mathbb{R}^2$. However, $\mathcal{T}(\boldsymbol{\theta}_0; \boldsymbol{\theta}_0)$ depends on exactly which of the inequalities defining Θ_0 are active at $\boldsymbol{\theta}_0$ (recall that $g_i(\boldsymbol{\theta}) \geq 0$ is said to be active at $\boldsymbol{\theta}_0$ if $g_i(\boldsymbol{\theta}_0) = 0$). We shall consider four different cases corresponding to different tangent cones of Θ_0 at $\boldsymbol{\theta}_0$.

Case 1: $g_1(\boldsymbol{\theta}_0) > 0$ and $g_2(\boldsymbol{\theta}_0) > 0$.

The condition “ $g_1(\boldsymbol{\theta}_0) > 0$ and $g_2(\boldsymbol{\theta}_0) > 0$ ” is equivalent to saying that $\boldsymbol{\theta}_0$ is an interior point of Θ_0 . Therefore, $\mathcal{T}(\boldsymbol{\theta}_0; \boldsymbol{\theta}_0) = \mathbb{R}^2$ and hence $\mathcal{T}(\boldsymbol{\theta}_0; \boldsymbol{\theta}_0) = \mathcal{T}(\boldsymbol{\theta}_1; \boldsymbol{\theta}_0)$.

Now, it follows from (4.105) that $LRT \xrightarrow{p} 0$. Therefore, the probability of rejecting H_0 when the true value is an interior point of Θ_0 converges to zero. Intuitively, this is to be expected, because it follows from $\hat{\theta} \xrightarrow{p} \theta_0$ that $LRT = 0$ and hence the LRT will not reject the null hypothesis with probability going to one.

Case 2: $g_1(\theta_0) = 0, g_2(\theta_0) > 0$.

The tangent cone at θ_0 is obtained by replacing the curved boundaries at θ_0 by their linear approximations. Equivalently, by Proposition 4.7.3, we have that $\mathcal{A}(\Theta_0; \theta_0) = \{\theta : (\theta - \theta_0)^T (\partial/\partial\theta) g_1(\theta_0) \geq 0\}$ and $\mathcal{T}(\Theta_0; \theta_0) = \mathcal{A}(\Theta_0; \theta_0) - \theta_0$. Therefore, $\mathcal{T}(\Theta_0; \theta_0) = \{\mathbf{x} : \mathbf{x}^T (\partial/\partial\theta) g_1(\theta_0) \geq 0\} = \{\mathbf{x} : x_2 \geq x_1(1 + 2\theta_{01})\}$. Now, with $Z \sim N(0, V)$ we have that

$$\text{pr}_{\theta_0}(LRT \geq c \mid H_0) \rightarrow \text{pr}_{\theta_0}(\|Z - \mathcal{T}(\Theta_0; \theta_0)\|_V^2 \geq c) = (1/2)\text{pr}(\chi_1^2 \geq c).$$

Case 3: $g_1(\theta_0) = g_2(\theta_0) = 0$.

As was shown in the previous example, $\theta_0 = 0$ and $\mathcal{T}(\Theta_0; \theta_0) = \{d \in \mathbb{R}^2 : d_2 \geq d_1, d_2 \geq -d_1\} = \{d : Rd \geq 0\}$, where R is the 2×2 nonsingular matrix $[-1, 1 | 1, 1]$. Therefore, by (4.105), $LRT \xrightarrow{d} \|Z - \mathcal{T}(\Theta_0; \theta_0)\|_V^2$ where $Z \sim N(0, V)$. Now by (3.39) it follows that

$$\text{pr}_{\theta_0}(LRT \geq c) \rightarrow 0.5\text{pr}(\chi_1^2 \geq c) + q\text{pr}(\chi_2^2 \geq c), \quad c > 0$$

where

$$q = (2\pi)^{-1} \cos^{-1}[(\mathbf{R}_1 V \mathbf{R}_2^T)/\{(\mathbf{R}_1 V \mathbf{R}_2^T)(\mathbf{R}_2 V \mathbf{R}_1^T)\}^{1/2}].$$

Case 4: $g_1(\theta_0) > 0$ and $g_2(\theta_0) = 0$.

This is similar to Case 2 and hence $\text{pr}(LRT \geq c \mid H_0) \rightarrow (1/2)\text{pr}(\chi_1^2 \geq c)$.

It follows from the different cases considered above that $\text{pr}(LRT \geq c \mid H_0)$ is a maximum when $g_1(\theta_0) = g_2(\theta_0) = 0$. Therefore, Case 3 is the least favorable null configuration, and hence $\theta_0 = 0$ is the least favorable null value, and

$$\lim \text{pr}(LRT \geq c \mid H_0) \leq 0.5\text{pr}(\chi_1^2 \geq c) + q\text{pr}(\chi_2^2 \geq c).$$

Since q is known, we can apply the LRT easily in this case. It will be seen later that if the null hypothesis involves several inequality constraints, for example, say $g_1(\theta) \geq 0, \dots, g_k(\theta) \geq 0, (k \geq 3)$, then a least favorable null value may or may not be a point in $\{\theta : g_1(\theta) = \dots = g_k(\theta) = 0\}$. The foregoing example illustrates how the asymptotic null distribution of the LRT depends on the shape of the parameter space at the null value. ■

4.8.5 Test of $\mathbf{h}(\theta) \geq 0$ Against $\mathbf{h}(\theta) \not\geq 0$

Let $\mathbf{h}(\theta)$ be a continuously differentiable $k \times 1$ vector function. Let the null and alternative hypotheses be

$$H_0 : \mathbf{h}(\theta) \geq 0 \text{ and } H_1 : \mathbf{h}(\theta) \not\geq 0, \quad (4.113)$$

In what follows, we shall assume that θ_0 is a boundary point of Θ_0 . Let us consider a setting that is simpler, yet general enough for many practical situations. Assume that $\mathbf{h}(\theta)$ is continuously differentiable and that the Mangasarian-Fromowitz condition

be less than k . Therefore, even if $\sqrt{n}(\hat{\theta} - \theta_0)$ converges to $N(\mathbf{0}, V)$ where V is positive definite, the asymptotic covariance of $\mathbf{h}(\hat{\theta})$ may not be positive definite. The foregoing null and alternative hypotheses may also be stated as

$$H_0 : \theta \in \Theta_0 \quad \text{and} \quad H_0 : \theta \in \Theta_1,$$

respectively, where $\Theta_0 = \{\theta : \mathbf{h}(\theta) \geq 0\}$ and $\Theta_1 = \mathbb{R}^p$. Therefore, $LRT = \sup_{\theta \in \Theta_0} \ell(\theta)$. As usual, assume that Θ_0 is Chernoff regular. Further, assume that Condition Q2 in Section 4.8.2 is satisfied.

To study the asymptotic null distribution of the LRT , assume that the null hypothesis holds and let θ_0 denote the true value. Wolak (1987, 1989, 1991) studied this testing problem and obtained important results; he also highlighted some of the difficulties that arise in this type of problems. Wolak (1988) pointed out that there is no guarantee that $\mathbf{h}(\theta) = 0$ would correspond to the least favorable null value. It follows from Section 4.8.2 that $LRT \xrightarrow{d} \|Z - \mathcal{T}(\Theta_0; \theta_0)\|^2$ where $Z \sim N\{\mathbf{0}, \mathcal{I}(\theta_0)^{-1}\}$ and $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathcal{I}(\theta_0) \mathbf{x}$.

For simplicity, let us consider the case when $\mathcal{T}(\Theta_0; \theta_0)$ is a closed convex cone. Now, from Theorem 3.4.2 and the fact that

$$\|Z - \mathcal{T}(\Theta_0; \theta_0)\|^2 = \|Z\|^2 - \|Z - \mathcal{T}^o(\Theta_0; \theta_0)\|^2,$$

it follows that

$$\text{pr}_{\theta_0}(LRT \geq c \mid H_0) \rightarrow \sum_{i=0}^p w_i \text{pr}(\chi_i^2 \geq c) \quad (4.114)$$

where $w_i = w_i\{p, \mathcal{I}(\theta_0)^{-1}, \mathcal{T}^o(\Theta_0; \theta_0)\}$ is the usual chi-bar square weight and T^o denotes the polar cone of \mathcal{T} with respect to $\mathcal{I}(\theta_0)^{-1}$. In general, the $\bar{\chi}^2$ -distribution in (4.114) depends on θ_0 through the weights, $w_i\{p, \mathcal{I}(\theta_0)^{-1}, \mathcal{T}^o(\Theta_0; \theta_0)\}$. Therefore, to implement a large sample test, we need to find

$$\sup_{\theta \in \Theta_0} \sum_{i=0}^p w_i \{p, \mathcal{I}(\theta_0)^{-1}, \mathcal{T}^o(\Theta_0; \theta_0)\} \text{pr}(\chi_i^2 \geq c), \quad c > 0. \quad (4.115)$$

Computation of this supremum can be a difficult task; however, we can always obtain bounds for it.

If θ_0 is an interior point of Θ_0 then $\mathcal{T}(\Theta_0; \theta_0) = \mathbb{R}^p$ and hence $\text{pr}(LRT \geq c \mid \theta_0) \rightarrow 0$ for $c > 0$. Therefore, the supremum in (4.115) cannot be achieved at an interior point of Θ_0 , and hence for the purposes of computing (4.115), we may restrict the attention to values of θ_0 on the boundary of Θ_0 .

It can be shown that at the least favorable null value, at least two of the inequality constraints in $\mathbf{h}(\theta_0) \geq 0$ must be active; this was shown by Wolak (1991). Thus, it follows that if there are only two inequality constraints, say $h_1(\theta) \geq 0$ and $h_2(\theta) \geq 0$, then they both must be active at the least favorable null value. ■

In what follows, we shall assume that θ_0 is a boundary point of Θ_0 . Let us consider a setting that is simpler, yet general enough for many practical situations. Assume that $\mathbf{h}(\theta)$ is continuously differentiable and that the Mangasarian-Fromowitz condition

(see Proposition 4.7.3) is satisfied. Because θ_0 is a boundary point of Θ_0 , it follows that $h_s(\theta_0) = 0$ for at least one s in $\{1, \dots, k\}$. Let a_i denote $(\partial/\partial\theta)h_i(\theta_0)$ and let $\{j_1, \dots, j_m\}$ denote $\{i : h_i(\theta_0) = 0\}$, the set of all indices of the active constraints at θ_0 ; note that m , the number of active constraints at θ_0 , is itself a function of θ_0 . Then

$$\mathcal{T}(\Theta_0; \theta_0) = \{\alpha : a_{j_1}^T \alpha \geq 0, \dots, a_{j_m}^T \alpha \geq 0\}.$$

Because $\mathcal{T}(\Theta_0; \theta_0)$ is a polyhedral, it is a closed convex cone, and hence it follows that $\|\mathbf{Z} - \mathcal{T}(\Theta_0; \theta_0)\|^2$, the limiting distribution of LRT , is a chi-bar-square (see Theorem 3.8.2).

Now, let $\mathbf{H}(\theta_0) = [a_{j_1}, \dots, a_{j_m}]^T$. Let us consider the case when $\text{rank}\{\mathbf{H}(\theta_0)\} = m(\theta_0)$. Since $LRT \xrightarrow{d} \|\mathbf{Z} - \mathcal{T}(\Theta_0; \theta_0)\|^2$, we have that

$$\Pr(LRT \geq c | \theta = \theta_0) \rightarrow \sum_{i=0}^m w_{m-i} \{m, V(\theta_0)\} \Pr(\chi_i^2 \geq c), \quad (4.116)$$

where $\mathbf{V}(\theta_0) = \mathbf{H}(\theta_0) \mathcal{I}(\theta_0)^{-1} \mathbf{H}(\theta_0)^T$. With t denoting the sample value of LRT , the large sample p -value, p_{sup} , is

$$\sup_{\theta_0 \in \Theta_0} \sum_{i=0}^m w_{m-i} \{m, V(\theta_0)\} \Pr(\chi_i^2 \geq t).$$

In general, the value of θ_0 at which this supremum is achieved may also depend on t . Consequently, there may not be a single least favorable null value for all possible values of t . It is not easy to maximize the tail probability of a chi-bar square distribution numerically using Newton-type iteration, because analytical expressions are not available for the derivatives of the tail probability of (4.115). Recent developments using simulation approaches may be helpful (see Shapiro, 2000); the utility of this approach remains to be investigated. In general, there is no easy way to compute the supremum in (4.115).

If computation of the supremum in (4.115) does not appear feasible, we can apply a bounds test using,

$$\begin{aligned} \sup_{\theta_0 \in \Theta_0} \sum_{i=0}^p w_i \{p, \mathcal{I}(\theta_0)^{-1}, \mathcal{T}^o(\theta_0; \theta_0)\} \Pr(\chi_i^2 \geq c) \\ \leq 0.5 \{\Pr(\chi_{p-1}^2 \geq t) + \Pr(\chi_p^2 \geq t)\}. \end{aligned}$$

If the maximum number of inequalities in the null hypothesis that can be satisfied simultaneously as equalities is, say q , and $q \leq p$, then the foregoing upper bound could be sharpened to $0.5 \{\Pr(\chi_{q-1}^2 \geq t) + \Pr(\chi_q^2 \geq t)\}$. Thus, while we can write down the asymptotic distribution of statistics for testing $\mathbf{h}(\theta) \geq 0$ vs $\mathbf{h}(\theta) \not\geq 0$, there is scope for developing practically feasible computational procedures to implement them.

The foregoing result extends in a natural way for testing

$$H_0 : \mathbf{h}_1(\theta) = 0, \mathbf{h}_2(\theta) \geq 0 \quad \text{vs} \quad H_1 : \text{No restrictions on } \theta$$

where $\mathbf{h} = (\mathbf{h}_1 : \mathbf{h}_2)$ is a partition of \mathbf{h} into two subvectors; only minor changes are required for this extension (see Kodde and Palm (1986) and Wolak (1989a)).

Example 1. Let $\theta \in \mathbb{R}^2$. Consider a statistical model satisfying Condition Q2 in Section 4.8.2. Let the null and alternative hypotheses be

$$H_0 : \theta \geq 0 \text{ and } H_1 : \theta \not\geq 0,$$

respectively, where $\theta \in \mathbb{R}^2$. Then $LRT = \sup\{\ell(\theta) : \theta \geq 0\}$. Suppose that the null hypothesis holds and let θ_0 denote the true null value. Then, the asymptotic null distribution of LRT is equal to that of $\|\mathbf{Z} - \mathcal{T}(\Theta_0; \theta_0)\|^2$ where Θ_0 is the positive orthant, $\mathbf{Z} \sim N\{0, \mathcal{I}(\theta_0)^{-1}\}$ and $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathcal{I}(\theta_0) \mathbf{x}$. Clearly, $\mathcal{T}(\Theta_0; \theta_0)$, and hence the foregoing null asymptotic distribution of LRT , depend on θ_0 . We shall consider four different cases.

Case 1: $\theta_0 = (0, 0)^T$. It is easily seen that $\mathcal{T}(\Theta_0; \theta_0)$ is the positive orthant and therefore,

$$\Pr(LRT \geq c | \theta = 0) \rightarrow (1/2) \Pr(\chi_1^2 \geq c) + q \Pr(\chi_2^2 \geq c),$$

where $q = (2\pi)^{-1} \cos^{-1}\{\mathcal{I}^{12}/(\mathcal{I}^{11}\mathcal{I}^{22})^{1/2}\}$ and $\mathcal{I}(0)^{-1} = (\mathcal{I}^{ij})_{2 \times 2}$.

Case 2: $\theta_0 = (0, a)^T$, where $a > 0$. Because $\mathcal{T}(\Theta_0; \theta_0)$ is the half-space $\{\alpha \in \mathbb{R}^2 : \alpha_1 \geq 0\}$,

$$\Pr\{LRT \geq c | \theta = (0, a)\} \rightarrow (1/2) \Pr(\chi_1^2 \geq c).$$

Case 3: $\theta_0 = (a, 0)^T$, where $a > 0$. The asymptotic null distribution is the same as that in Case 2.

Case 4: $\theta_0 = (a, b)^T$, where $a > 0$ and $b > 0$. Now, $\mathcal{T}(\Theta_0; \theta_0) = \mathbb{R}^2$, and therefore, $\|\mathbf{Z} - \mathcal{T}(\Theta_0; \theta_0)\| = 0$ and $\Pr\{LRT = 0 | \theta = (a, b)\} \rightarrow 1$.

It follows from the foregoing four cases that $\theta = 0$ is the least favorable null value. Therefore, the large sample p -value corresponding to $LRT = t$ is $(1/2) \Pr(\chi_1^2 \geq t) + q \Pr(\chi_2^2 \geq t)$. ■

Example 2. Let the testing problem be

$$H_0 : \mathbf{h}(\theta) \geq 0 \text{ vs } H_1 : \mathbf{h}(\theta) \not\geq 0$$

where $\mathbf{h} = (\mathbf{h}_1 : \mathbf{h}_2)$. Suppose that $\mathbf{h}_1(\theta) = \mathbf{h}_2(\theta) = 0$ has a solution, say at θ_0 and $\text{rank}[\mathbf{H}(\theta_0)] = 2$. As in the previous example, it would be convenient to consider four different cases.

Case 1: $\mathbf{h}_1(\theta_0) = 0$ and $\mathbf{h}_2(\theta_0) = 0$

$$\begin{aligned} \Pr\{LRT \geq c | \mathbf{h}_1(\theta_0) = 0, \mathbf{h}_2(\theta_0) = 0\} \\ \rightarrow (1/2) \Pr(\chi_1^2 \geq c) + q(\theta_0) \Pr(\chi_2^2 \geq c), \end{aligned}$$

where

$$q(\theta) = (2\pi)^{-1} \cos^{-1}\{V_{12}/(V_{11}V_{22})^{1/2}\}, \text{ and } \mathbf{V} = \mathbf{V}(\theta_0).$$

Case 2: $h_1(\theta_0) > 0$ and $h_2(\theta_0) = 0$

$$\text{pr}\{LRT \geq c \mid h_1(\theta_0) > 0, h_2(\theta_0) = 0\} \rightarrow (1/2)\text{pr}(\chi_1^2 \geq c).$$

Case 3: $h_1(\theta_0) = 0$ and $h_2(\theta_0) > 0$

$$\text{pr}\{LRT \geq c \mid h_1(\theta_0) = 0, h_2(\theta_0) > 0\} \rightarrow (1/2)\text{pr}(\chi_1^2 \geq c);$$

Case 4: $h_1(\theta_0) > 0$ and $h_2(\theta_0) > 0$

$$\text{pr}\{LRT = 0 \mid h_1(\theta_0) > 0, h_2(\theta_0) > 0\} \rightarrow 1.$$

Therefore, the least favorable null value is a point at which $h_1(\theta) = h_2(\theta) = 0$. If $h_1(\theta) = h_2(\theta) = 0$ has more than one solution, then the solution at which $q(\theta)$ is the maximum is the least favorable null value. ■

Example 3 (Wolak (1991)). Let $\theta = (\theta_1, \theta_2)$ and

$$H_0: h_1(\theta) \geq 0 \text{ and } h_2(\theta) \geq 0,$$

where $h_1(\theta) = \theta_2 - \theta_1^2$ and $h_2(\theta) = \theta_1 - \theta_2$; the parameter space is shown in Fig. 4.8. As in the previous example, there are four cases to be considered. It follows that the least favorable null value is a θ_0 satisfying $h_1(\theta_0) = 0$ and $h_2(\theta_0) = 0$. Let t denote the sample value of LRT . Because there are two values of θ_0 , namely $(0, 0)$ and $(1, 1)$, satisfying these equations, and

$$\text{pr}(LRT \geq t \mid h(\theta_0) = 0) \rightarrow (1/2)\text{pr}(\chi_1^2 \geq t) + q(\theta_0)\text{pr}(\chi_2^2 \geq t),$$

the least favorable null value is the θ_0 in $\{(0, 0), (1, 1)\}$ for which

$$(1/2)\text{pr}(\chi_1^2 \geq t) + q(\theta_0)\text{pr}(\chi_2^2 \geq t)$$

is a maximum. If $\mathcal{I}(\theta_0)$ is known for $\theta_0 = (0, 0)$ and $(1, 1)$, then we can compute the foregoing tail probability for each of the two values of θ_0 and choose the largest as the p -value. ■

In the foregoing examples with only two inequality constraints in the null hypothesis, the least favorable null value turned out to be a point at which all the inequalities in the null hypothesis are active. We also observed in Chapter 2 that for testing $H_0: \theta \geq 0$ against $H_1: \theta \not\geq 0$, based on observations from $N(\theta, V)$ where V does not depend on θ , the least favorable null value, namely 0 , is the point at which every inequality constraint in the null hypothesis is active. However, in general, if the null hypothesis is $H_0: h_1(\theta) \geq 0, \dots, h_k(\theta) \geq 0$, where $k > 2$, then the least favorable null value may not be a point at which all the inequalities are active. The following example from Wolak (1991) illustrates this.

Example 4 (Wolak (1991)). Let $\mathbf{X} = (X_1, X_2)^T \sim N(\mathbf{0}, \mathbf{V})$ where

$$\mathbf{V} = \begin{pmatrix} \sigma^2 & \rho\sigma\tau \\ \rho\sigma\tau & \tau^2 \end{pmatrix}.$$

Let $\theta = (\sigma^2, \tau^2, \rho)^T$,

$$\Theta = \{\theta : \theta_1 > 0, \theta_2 > 0, -1 < \theta_3 < 1\},$$

$$\theta^A = (1, 1, 0.95)^T \text{ and } \theta^B = (1, 1, 0)^T. \text{ Let the null hypothesis be } H_0 : \theta \leq \theta^A.$$

To express this testing problem in the notation of (4.113), let $\mathbf{h}(\theta) = (1 - \theta_1, 1 - \theta_2, 0.95 - \theta_3)$; then $\mathbf{H}(\theta_0) = -\mathbf{I}$. The inverse of the information matrix and the corresponding correlation matrix are given below (see Lehmann (1983), p 441).

$$\mathcal{I}(\theta)^{-1} = \begin{bmatrix} 2\sigma^4 & 2\rho^2\sigma^2\tau^2 & \rho(1-\rho^2)\sigma^2 \\ 2\rho^2\sigma^2\tau^2 & 2\tau^4 & \rho(1-\rho^2)\tau^2 \\ \rho(1-\rho^2)\sigma^2 & \rho(1-\rho^2)\tau^2 & (1-\rho^2)^2 \end{bmatrix}$$

$$\text{corr}\{\mathcal{I}(\theta)^{-1}\} = \begin{bmatrix} 1 & \rho^2 & \rho/\sqrt{2} \\ \rho^2 & 1 & \rho/\sqrt{2} \\ \rho/\sqrt{2} & \rho/\sqrt{2} & 1 \end{bmatrix}.$$

Note that Θ is open, $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\{\mathbf{0}, \mathcal{I}(\theta_0)^{-1}\}$, and Condition Q2 in subsection 4.8.2 is satisfied. Now,

$$\lim_{n \rightarrow \infty} \text{pr}(LRT \geq c \mid \theta = \theta^A) = w_1\{3, \mathcal{I}(\theta_0)\}\text{pr}(\chi_3^2 \geq c) + w_2\{3, \mathcal{I}(\theta_0)\}\text{pr}(\chi_2^2 \geq c) + w_3\{3, \mathcal{I}(\theta_0)\}\text{pr}(\chi_1^2 \geq c), \quad (4.117)$$

where w_1, w_2 , and w_3 are determined by the formulas in (3.26). By direct substitution into these formulas, we have $w_1 = 0.0153, w_2 = 0.1682, w_3 = 0.4847$,

$$\lim_{n \rightarrow \infty} \text{pr}\{LRT \geq c \mid \theta = \theta^B\} = 0.5\{\text{pr}(\chi_1^2 \geq c) + 0.5\text{pr}(\chi_2^2 \geq c)\}, \quad (4.118)$$

$$\text{and} \quad \lim_{n \rightarrow \infty} \text{pr}(LRT \geq c \mid \theta = \theta^A) < \text{pr}\{LRT \geq c \mid \theta = \theta^B\}.$$

Therefore, the least favorable null value is not θ^A , the only point at which all the inequalities in the null hypothesis are active.

Let θ^* be a point at which exactly two of those inequalities are active. Then,

$$\lim \text{pr}\{LRT \geq c \mid \theta = \theta^*\} = 0.5\text{pr}(\chi_1^2 \geq c) + q\text{pr}(\chi_2^2 \geq c),$$

for some q in the range, $0 \leq q \leq 0.5$. Let θ^+ be a point at which exactly one of the inequalities in H_0 is active. Then,

$$\lim \text{pr}\{LRT \geq c \mid \theta = \theta^+\} = 0.5\text{pr}(\chi_1^2 \geq c).$$

Since this limit is smaller than the RHS of (4.118), it follows that $\theta = \theta^B$ is the least favorable null value. ■

Example 4.8.3 (Stram and Lee, 1994) Testing for the presence of random effects.

Consider the growth curve model in Example 1.2.10 (p. 14). The data consist of measurements of the distance (in mm) from the center pituitary to the pterygomaxillary fissure for 27 children (11 girls and 16 boys). Every subject has four measurements, taken at ages 8, 10, 12 and 14 years. The object is to model the growth in the distance as a function of age and sex. To allow for individual variability in the growth function, the model being considered is

$$\mathbf{y}_i = (1, \text{Age}, \text{Sex}, \text{Age} \cdot \text{Sex})\boldsymbol{\alpha} + (1, \text{Age})\boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i$$

where \mathbf{y}_i is 4×1 for the four measurements of the i th subject, $\text{Age} = (8, 10, 12, 14)^T$, $\text{Sex} = (a, a, a, a)^T$ where a is 0 for boys and 1 for girls. The parameter $\boldsymbol{\alpha}(4 \times 1)$ captures the fixed effects and $\boldsymbol{\beta}_i(2 \times 1)$ captures the random effect of subject i . Assume that $\boldsymbol{\epsilon}_i \sim N(0, \sigma^2 I)$, $\boldsymbol{\beta}_i \sim N(0, \Psi)$ for some 2×2 positive semi-definite matrix Ψ and that $\boldsymbol{\epsilon}_i$ and $\boldsymbol{\beta}_i$ are independent. The following models are of interest:

1. Model $M1 : \Psi = \mathbf{0}$. There are no random effects.
2. Model $M2 : \Psi = (\psi_{11}, 0 \mid 0, 0)$. The individual variability is captured by a random intercept; therefore, the regression planes are parallel for different subjects.
3. Model $M3 : \Psi = (\psi_{11}, \psi_{12} \mid \psi_{21}, \psi_{22})$. Individual variability results in not only a random intercept but also a different slope parameter for age.

To simplify notation, let us write $\Psi = (\psi_1, \psi_2 \mid \psi_2, \psi_3)$. In the process of modeling these data, it is possible that we may be interested to know whether or not the fit for Model $M3$ is significantly better than that for the simpler Model $M2$. Therefore, we formulate the testing problem as

$$H_0 : \psi_1 \geq 0, \psi_2 = \psi_3 = 0 \quad \text{vs} \quad H_1 : \psi_1 \geq 0, \psi_3 \geq 0, \psi_1\psi_3 - \psi_2^2 \geq 0.$$

The regression model can be written as

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\alpha} + \mathbf{Z}_i \boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i.$$

With the assumptions as introduced earlier for the distributions of $\boldsymbol{\epsilon}_i$ and $\boldsymbol{\beta}_i$, we have

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^n (-n_i/2) \log(2\pi) - (1/2) \log\{\det(\boldsymbol{\Sigma}_i^{-1})\} - (1/2) \text{trace}(\boldsymbol{\Sigma}_i^{-1} \mathbf{S}_i)$$

where $\boldsymbol{\Sigma}_i = \sigma^2 I + \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i^T$ and $\mathbf{S}_i = (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\alpha})^T (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\alpha})$. Now, under the null,

$$LRT \xrightarrow{d} \|U - \mathcal{T}(\boldsymbol{\Theta}_0; \boldsymbol{\theta}_0)\|^2 - \|U - \mathcal{T}(\boldsymbol{\Theta}_1; \boldsymbol{\theta}_0)\|^2$$

where $\mathbf{U} \sim N(\mathbf{0}, \mathcal{I}_{\boldsymbol{\Theta}_0}^{-1})$,

$$\boldsymbol{\Theta}_0 = \{\boldsymbol{\psi} \in \mathbb{R}^3 : \psi_2 = \psi_3 = 0, \psi_1 \geq 0\},$$

$$\text{and } \boldsymbol{\Theta}_1 = \{\boldsymbol{\psi} \in \mathbb{R}^3 : \psi_1 \geq 0, \psi_3 \geq 0, \psi_1\psi_3 - \psi_2^2 \geq 0\}.$$

There are two possible values of $\boldsymbol{\theta}_0$ in $\boldsymbol{\Theta}_0$ that we need to consider separately: $(0, 0, 0)^T$ and $(a, 0, 0)^T$ where $a > 0$.

Case 1: $\boldsymbol{\theta}_0 = (0, 0, 0)^T$.

It may be verified that

$$\mathcal{T}(\boldsymbol{\Theta}_0; \boldsymbol{\theta}_0) = \{\boldsymbol{\psi} \in \mathbb{R}^3 : \psi_2 = \psi_3 = 0, \psi_1 \geq 0\}.$$

To determine the tangent cone of $\boldsymbol{\Theta}_1$, note that

$$h_3(\boldsymbol{\theta}) = \psi_1\psi_3 - \psi_2^2, \quad \text{and} \quad (\partial/\partial\boldsymbol{\theta})h_3(\boldsymbol{\theta}) = (\psi_3, -2\psi_2, \psi_1)^T.$$

At $\boldsymbol{\theta}_0 = (0, 0, 0)^T$ we have $h_3(\boldsymbol{\theta}_0) = 0$ and $(\partial/\partial\boldsymbol{\theta})h_3(\boldsymbol{\theta}_0) = \mathbf{0}$. Therefore, MF-CQ is not satisfied at $\boldsymbol{\theta}_0$, and hence we cannot use Proposition 4.7.3 to write down the tangent cone; thus we should not simply replace the nonlinear constraints by their linear approximations at $\boldsymbol{\theta}_0$.

Since $\boldsymbol{\Theta}_1$ is a closed convex cone with vertex at the origin it follows that $\mathcal{T}(\boldsymbol{\Theta}_1; \boldsymbol{\theta}_0) = \boldsymbol{\Theta}_1$. Since neither $\mathcal{T}(\boldsymbol{\Theta}_1; \boldsymbol{\theta}_0)$ nor $\mathcal{T}(\boldsymbol{\Theta}_0; \boldsymbol{\theta}_0)$ is a linear space, the asymptotic null distribution of LRT may not be a chi-bar square. Exact distribution of LRT for this case is nonstandard. A bound can be obtained by replacing $\mathcal{T}(\boldsymbol{\Theta}_1; \boldsymbol{\theta}_0)$ by the larger set $\{\psi_1 \geq 0, \psi_3 \geq 0\}$. In this case, the distribution is a mixture of chi-square distributions.

Case 2: $\boldsymbol{\theta}_0 = (a, 0, 0)^T$ where $a > 0$.

It may be verified that

$$\begin{aligned} \mathcal{T}(\boldsymbol{\Theta}_0; \boldsymbol{\theta}_0) &= \{\boldsymbol{\psi} \in \mathbb{R}^3 : \psi_2 = \psi_3 = 0\}, \\ \text{and} \quad \mathcal{T}(\boldsymbol{\Theta}_1; \boldsymbol{\theta}_0) &= \{\boldsymbol{\psi} \in \mathbb{R}^3 : \psi_3 \geq 0\}, \end{aligned}$$

the Mangasarian-Fromowitz constraint qualification is satisfied at $\boldsymbol{\theta}_0$ and therefore we can use Proposition 4.7.3 to write down the tangent cones. Since $\mathcal{T}(\boldsymbol{\Theta}_0; \boldsymbol{\theta}_0)$ is a linear space and $\mathcal{T}(\boldsymbol{\Theta}_1; \boldsymbol{\theta}_0)$ is a closed convex cone it follows that the asymptotic null distribution of LRT is a chi-bar square. ■

4.9 PROPERTIES OF THE MLE WHEN THE TRUE VALUE IS ON THE BOUNDARY

We continue to consider the *iid* setting for simplicity although most of the results in this section would hold in more general settings as will be indicated later. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ denote *iid* observations from a population with density function $f(\mathbf{x}; \boldsymbol{\theta})$. Let $\ell(\boldsymbol{\theta})$ denote the loglikelihood, where $\boldsymbol{\theta} \in \boldsymbol{\Theta} \subset \mathbb{R}^p$. Let $\Omega \subset \boldsymbol{\Theta}$ where Ω is not necessarily open. In this section we study large sample properties of local and global mle's of $\boldsymbol{\theta}$. Let us denote an *mle*, global or local, by $\hat{\boldsymbol{\theta}}$. The true value, denoted by $\boldsymbol{\theta}_0$, will be assumed to be a boundary point of Ω unless the contrary is made clear. *Assume that Condition Q is satisfied.* For simplicity, we shall restrict to maximum likelihood. However, most of the results of this section hold for M -estimators as well; for details of this case see Geyer (1994) and Shapiro (2000a).

Distribution of the *mle* over a restricted parameter space that is not open has not attracted as much attention as the corresponding hypothesis testing problems. If θ_0 were an interior point then $\sqrt{n}\hat{\mathcal{I}}^{-1/2}(\hat{\theta} - \theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{I})$. Since this limiting distribution does not depend on any unknown parameters (i.e., $\sqrt{n}\hat{\mathcal{I}}^{-1/2}(\hat{\theta} - \theta_0)$ is a pivotal quantity), we can use this result to construct a confidence region for θ . Unfortunately, the asymptotic distribution of $\sqrt{n}\hat{\mathcal{I}}^{-1/2}(\hat{\theta} - \theta_0)$ is a discontinuous function of θ_0 . Further, it is unknown whether or not there is a pivotal quantity based on $\hat{\theta}$ that would lend itself for constructing a confidence region. Consequently, statistical inference based on the constrained estimator $\tilde{\theta}$ has not attracted much attention. Nevertheless, it is of interest to study properties of the constrained *mle* because it has an important role to play in the theory of statistical inference. In this section, we shall consider some examples to illustrate the main results. Then we shall state the main results and provide references to sources where detailed statements of the regularity conditions and proofs may be found.

At this stage it would be helpful to reconsider Example 3.3.1.

Example 4.9.1 Asymptotic distribution of mle when $\Omega = \mathbb{R}^{+2}$ and $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{I})$.

Let $\mathbf{X} \sim N(\boldsymbol{\theta}, \mathbf{I})$ and $\Omega = \mathbb{R}^{+2}$. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be iid as \mathbf{X} . Then,

$$\ell(\boldsymbol{\theta}) = (-n/2)\|\bar{\mathbf{X}} - \boldsymbol{\theta}\|^2$$

and $\tilde{\theta}$ is the point in \mathbb{R}^{+2} that is closest to $\bar{\mathbf{X}}$. Thus,

$$\tilde{\theta} = \Pi(\bar{\mathbf{X}}, \mathbb{R}^{+2}) \text{ and } \tilde{\theta} \xrightarrow{P} \theta_0.$$

If θ_0 is an interior point of \mathbb{R}^{+2} then $\Pr(\bar{\mathbf{X}} = \tilde{\theta}) \rightarrow 1$ and hence $\sqrt{n}(\tilde{\theta} - \theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{I})$. Let $\mathcal{T} = \mathcal{T}(\mathbb{R}^{+2}, \theta_0)$, $Z_n = \sqrt{n}(\bar{\mathbf{X}} - \theta_0)$ and $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$. Suppose that θ_0 lies on the boundary of \mathbb{R}^{+2} . Then it is easily seen that

$$(\tilde{\theta} - \theta_0) = \Pi(\bar{\mathbf{X}} - \theta_0, \mathcal{T}) \text{ with probability approaching 1.}$$

Therefore,

$$\sqrt{n}(\tilde{\theta} - \theta_0) = \Pi(\mathbf{Z}_n, \mathcal{T}) \text{ with prob approaching 1.}$$

Since $\Pi(z, \mathcal{T})$ is continuous in z and $\mathbf{Z}_n \xrightarrow{d} \mathbf{Z}$ it follows that $\Pi(\mathbf{Z}_n, \mathcal{T}) \xrightarrow{d} \Pi(\mathbf{Z}, \mathcal{T})$. Therefore, we have

$$\sqrt{n}(\tilde{\theta} - \theta_0) \xrightarrow{d} \Pi(\mathbf{Z}, \mathcal{T}).$$

This may be verified directly (in fact, it would be instructive to do so) by considering the following three cases separately: (i) θ_0 is on the positive θ_1 -axis (this is the case shown in Fig. 4.10), (ii) θ_0 is on the positive θ_2 -axis and (iii) θ_0 is at the origin. It is perhaps easier to visualize this result as

$$(\tilde{\theta} - \theta_0) \approx \Pi(\mathbf{Z}/\sqrt{n}; \mathcal{T}),$$

for large n as shown in Fig. 4.10. ■

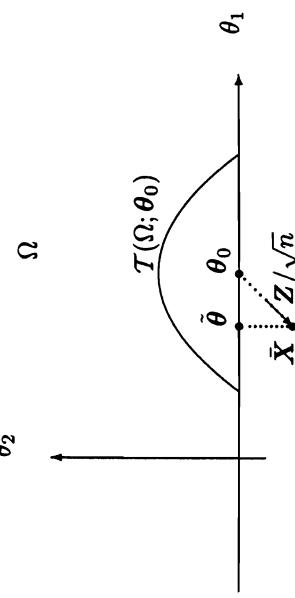


Fig. 4.10 $(\tilde{\theta} - \theta_0) \approx \Pi\{\mathbf{Z}/\sqrt{n}; \mathcal{T}(\Omega; \theta_0)\}$, where $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$.

When θ_0 is an interior point, a consistent *mle* $\hat{\theta}$, whether it is a global or local *mle*, satisfies $n^{-1/2}\hat{\mathcal{I}}^{-1/2}\mathcal{S}(\theta_0) = \sqrt{n}(\hat{\theta} - \theta_0) + o_p(1)$. This is a neat representation of $\hat{\theta}$ as a sum of independently and identically distributed variables ($\mathcal{S}(\theta_0)$ is a sum of iid variables). This is a direct result of applying a one-term Taylor expansion on the first-order condition $\nabla\ell(\hat{\theta}) = \mathbf{0}$. If θ_0 is not an interior point, then the first-order conditions, such as the Kuhn-Tucker conditions, involve inequalities. Consequently, $\sqrt{n}(\tilde{\theta} - \theta_0)$ does not have a simple representation as for $\sqrt{n}(\hat{\theta} - \theta_0)$. In fact, the properties of local *mle*'s are very subtle. Let us consider a simple example to illustrate this.

Example 4.9.2 Nonconvex parameter space

Let \mathbf{X} have the bivariate normal distribution $N(\mathbf{0}, \mathbf{I})$ and let X_1, \dots, X_n be iid as \mathbf{X} . Then $\ell(\boldsymbol{\theta}) = (-n/2)\|\bar{\mathbf{X}} - \boldsymbol{\theta}\|^2$. Let Ω be the union of the θ_1 -axis and the θ_2 -axis; thus $\Omega = \{(\theta_1, \theta_2) : \theta_1 = 0 \text{ or } \theta_2 = 0\}$. In this case, there are two local *mle*'s: $\tilde{\theta}^a = (\bar{X}_1, 0)$ and $\tilde{\theta}^b = (0, \bar{X}_2)$. It is clear that their large sample distributions are different. Further, when θ_0 is on the θ_1 -axis, we have the following with probability $\rightarrow 1$:

$$\tilde{\theta}^a - \theta_0 = \text{projection of } (\bar{\mathbf{X}} - \theta_0) \text{ onto the } \theta_1\text{-axis,}$$

and $\sqrt{n}(\tilde{\theta}^a - \theta_0)$ converges to $\Pi(\mathbf{Z}, \mathcal{P})$ where $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$ and $\mathcal{P} = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \theta_2 = 0\}$. In fact, we have

$$\sqrt{n}(\tilde{\theta}^a - \theta_0) \xrightarrow{d} \begin{cases} \Pi(\mathbf{Z}, \mathcal{P}) & \text{if } \theta_0 \in \theta_1\text{-axis} \\ \begin{bmatrix} \sqrt{n}\bar{X}_1 \\ -\sqrt{n}\theta_{02} \end{bmatrix} & \text{if } \theta_0 \notin \theta_1\text{-axis.} \end{cases}$$

A similar comment applies to $\tilde{\theta}^b$ by symmetry.

By contrast, with $\boldsymbol{\theta}$ denoting the global *mle*, it is easily seen that

$$(\tilde{\theta} - \theta_0) = \Pi(\bar{\mathbf{X}} - \theta_0, \mathcal{T}(\Omega; \theta_0)) \text{ with prob approaching 1.}$$

Therefore,

$$\sqrt{n}(\tilde{\theta} - \theta_0) \xrightarrow{d} \Pi(Z, \mathcal{T}(\Omega; \theta_0)), \text{ where } Z \sim N(0, I).$$

Thus, local and global mle's have different asymptotic distributions. If the parameter space is \mathbb{R}^{+2} then there is only one maximum, hence local and global mle's are the same. This example illustrates that if the parameter space is not convex, or at least nearly convex in a way to be explained later, the global and local mle's may have different properties. A consequence of such a difference between local mle's is that if an iterative algorithm with Newton-Raphson type steps were to be used to find the global maximum, it could be trapped around a local maximum. ■

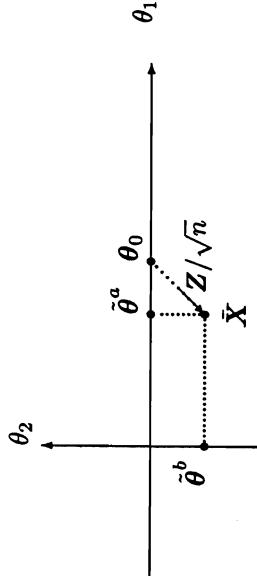


Fig. 4.11 Two local mles when Ω is the union of θ_1 -axis and the θ_2 -axis, and θ_0 is not on the θ_2 -axis; $\sqrt{n}(\tilde{\theta}^a - \theta_0)$ converges in distribution but $\|\sqrt{n}(\tilde{\theta}^b - \theta_0)\| \rightarrow \infty$.

The standard method that is adopted in the literature to study large sample properties of consistent mle's is very similar to that was adopted so far to study the properties of LRT. To illustrate the main idea, let us first consider the case when the parameter space Ω is a closed convex cone with its vertex at θ_0 . Let $\mathcal{T} = \Omega - \theta_0$; thus \mathcal{T} is the tangent cone of Ω at θ_0 . Suppose that Condition Q is satisfied. Let us write the quadratic approximation of $\ell(\theta)$ as (see (4.3))

$$\ell(\theta) = K_n - 2^{-1}(Z_n - \mathbf{u})^T I_{\theta_0}(Z_n - \mathbf{u}) + \delta_n(\mathbf{u}),$$

where $\mathbf{u} = \sqrt{n}(\theta - \theta_0)$ and K_n does not depend on θ . Let

$$\theta^\dagger = \operatorname{argmin}_{\theta \in \Omega} (Z_n - \sqrt{n}(\theta - \theta_0))^T I_{\theta_0}(Z_n - \sqrt{n}(\theta - \theta_0)).$$

A consequence of $\delta_n(\mathbf{u}) = n^{-1/2}\|\mathbf{u}\|^3 O_p(1)$ is that

$$n^{1/2}(\tilde{\theta} - \theta^\dagger) = o_p(1)$$

(see Self and Liang (1987, Lemma 2)). Therefore, $n^{1/2}(\tilde{\theta} - \theta_0)$ and $n^{1/2}(\theta^\dagger - \theta_0)$ have the same asymptotic distribution. Now, recall the following notation for projection:

$$\Pi_W(z, \mathcal{T}) = \operatorname{argmin}_{u \in \mathcal{T}} (z - u)^T V^{-1}(z - u).$$

Since \mathcal{T} is convex and projection onto a convex set is distance reducing, it follows that $\Pi_W(z, \mathcal{T})$ is continuous in (z, W) .

$$\|z^t(W_1 - W_2)z\| \leq \|W_1 - W_2\| \|z\|^2$$

that

$$\Pi_W(z, \mathcal{T}) \text{ is continuous in } (z, W).$$

Since $Z_n \xrightarrow{d} Z \sim N(\mathbf{0}, I_{\theta_0}^{-1})$ and $\sqrt{n}(\theta^\dagger - \theta_0) = \Pi(Z_n, \mathcal{T})$, it follows that

$$\sqrt{n}(\theta^\dagger - \theta_0) \xrightarrow{d} \Pi(Z, \mathcal{T}).$$

Therefore, the asymptotic distribution of the global mle is given by

$$\sqrt{n}(\tilde{\theta} - \theta_0) \xrightarrow{d} \Pi(Z, \mathcal{T}).$$

These arguments can also be modified when Ω is not necessarily a cone, but provided $\mathcal{T}(\Omega, \theta_0)$ is convex; see Andrews (1999) and Le Cam (1970, p 820) for details. This result has also been extended to the case when $\mathcal{T}(\Omega, \theta_0)$ is not convex, but Ω is Chernoff regular at θ_0 . The next result states the most general form known for the iid setting.

Proposition 4.9.1 Suppose that Ω is Chernoff regular at θ_0 , and let $\tilde{\theta}$ denote the global mle. Suppose that Condition Q is also satisfied. Then

$$\sqrt{n}(\tilde{\theta} - \theta_0) \xrightarrow{d} \Pi_{\mathcal{T}(\theta_0)^{-1}}\{Z, \mathcal{T}(\Omega; \theta_0)\}, \text{ where } Z \sim N(0, \mathcal{I}(\theta_0)^{-1}).$$

Proof: See Self and Liang (1987, Theorem 2) and Geyer (1994, Theorem 4.4); the proof in Geyer (1994) is applicable for M -estimators as well. ■

To state the foregoing result differently, let $Y \sim N(\theta, \mathcal{I}(\theta_0)^{-1})$ where $\theta \in \Omega$. Let G denote the distribution of the mle of θ based on a single observation of Y when $\theta = \theta_0$. Then the asymptotic distribution of $\sqrt{n}(\tilde{\theta} - \theta_0)$ is G . This is illustrated in Fig. 4.12. It shows that the first order asymptotic behavior of $(\tilde{\theta} - \theta_0)$ is the same as that of $\Pi\{Z/\sqrt{n}, \mathcal{T}(\Omega; \theta_0)\}$, where $Z \sim N(0, \mathcal{I}(\theta_0)^{-1})$.

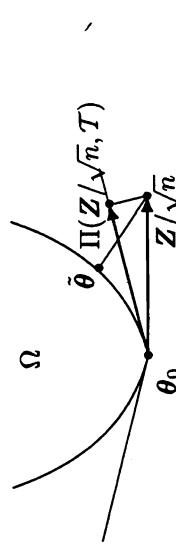


Fig. 4.12 $(\tilde{\theta} - \theta_0) \approx \Pi\{Z/\sqrt{n}, \mathcal{T}\}$, where $Z \sim N(0, \mathcal{I}(\theta_0))$ and $\mathcal{T} = \mathcal{T}(\Omega; \theta_0)$.

To study properties of local mle's, let us introduce the terms *near convexity of a set* and *pro-regular*.

Definition: We say that the set Ω is *nearly convex* at θ_0 if there exists a neighborhood V of θ_0 and a function $k(\theta, \theta^*)$ tending to zero as $\theta \rightarrow \theta_0$, $\theta^* \rightarrow \theta_0$, such that

$$\text{dist}\{\theta^* - \theta, T(\Omega; \theta)\} \leq k(\theta, \theta^*)\|\theta^* - \theta\|, \quad \forall \theta, \theta^* \in \Omega \cap V.$$

We say that the set Ω is *pro-regular* at θ_0 if there exists a neighborhood V of θ_0 and a constant K such that

$$\text{dist}\{\theta^* - \theta, T(\Omega; \theta)\} \leq K\|\theta^* - \theta\|^2, \quad \forall \theta, \theta^* \in \Omega \cap V.$$

If Ω is convex at θ_0 then it is nearly convex and pro-regular at θ_0 . The concept of *near convexity* was introduced by Shapiro and Al-Khayyal (1993). *Pro-regularity* appears in Shapiro (1994) under the name “O(2)-Convexity”; for related details on this topic see Rockafellar and Wets (1998).

Another concept that is related to this is *monotonicity of normals*. This is closely related to what is known as *Clark regularity*, convexity of tangent cones and near convexity. See Shapiro (2000a) for details about these concepts. Geyer (1994) provides an excellent discussion of *Clark regularity*; for more detailed discussions see Rockafellar and Wets (1998).

If Ω satisfies the Mangasarian and Fromowitz constraint qualification at θ_0 , then Ω is nearly convex at θ_0 . This is particularly relevant because many of the constrained inference problems in statistics are likely to be specified by equality and inequality constraints on functions of θ , and the Mangasarian and Fromowitz condition is likely to be satisfied in many cases. Further, if the constraining functions $h_i(\theta)$ are Lipschitz continuous in a neighborhood of θ_0 then Ω is pro-regular at θ_0 .

Proposition 4.9.2 (Shapiro (2000a)). Suppose that Ω is nearly convex at θ_0 and Condition Q is satisfied. Let $\tilde{\theta}^a$ and $\tilde{\theta}^b$ be two local mle's. Then,

$$\sqrt{n}(\tilde{\theta}^a - \tilde{\theta}^b) = o_p(1). \blacksquare$$

Shapiro (2000a) showed, by constructing a counter example, that even if the tangent cone of Ω at θ_0 is convex, two local mle's may fail to be asymptotically equivalent [i.e., $\sqrt{n}(\tilde{\theta}^a - \tilde{\theta}^b)$ may not be $o_p(1)$]. In fact, Shapiro (2000a) showed that even if Ω is *Clark regular*, $\sqrt{n}(\tilde{\theta}^a - \tilde{\theta}^b)$ may not be $o_p(1)$. Clark regularity is weaker than near convexity, but it is a considerably stronger than Chernoff regularity; for example, it ensures that the tangent cone is convex.

Suppose that Ω is pro-regular at θ_0 and $\tilde{\theta}^a$ and $\tilde{\theta}^b$ are strongly consistent. Then, under some regularity conditions on the smoothness of and uniformity of convergence of $n^{-1}\ell(\theta)$ we have $\tilde{\theta}^a = \tilde{\theta}^b$ with probability approaching one (see Shapiro (2000a) for details).

Throughout this section we restricted our discussions to the simple *iid* setting. However, these results hold under much more general conditions. The developments in Geyer (1994) and Shapiro (2000a) include M -estimation; Andrews (1999) studied properties of the global mle under conditions that are more general than the usual

iid setting. For example, his framework includes many time series models used in econometric modeling. He considered several important econometric models and illustrated the relevance of the theoretical results. In these papers quadratic approximations are used even if the objective function is not differentiable; see also Shapiro (1989).

4.10 APPENDIX: PROOFS

Proof of Proposition 4.3.2: Assume that the null hypothesis is true. The test statistic T_λ is a function of the sample. Let $T_\lambda(\mathcal{X})$ denote the value of the test statistic T_λ for the sample \mathcal{X} . Let \mathcal{Y} be iid as \mathcal{X} . Then

$$p^*(\mathcal{X}) = \alpha_1 + \sup_{\lambda \in \mathcal{A}(\mathcal{X})} \text{pr}_{\mathcal{Y}|\mathcal{X}}\{T_\lambda(\mathcal{Y}) \geq T_\lambda(\mathcal{X}) \mid \theta = (0 : \lambda)\}.$$

Now,

$$\begin{aligned} \text{pr}_{\mathcal{X}}\{p^*(\mathcal{X}) \leq \alpha\} &= \text{pr}_{\mathcal{X}}\{p^*(\mathcal{X}) \leq \alpha \text{ and } \lambda_0 \in \mathcal{A}(\mathcal{X})\} + \\ &\quad \text{pr}\{p^*(\mathcal{X}) \leq \alpha \text{ and } \lambda_0 \notin \mathcal{A}(\mathcal{X})\} \\ &\leq \text{pr}\{p^*(\mathcal{X}) \leq \alpha \text{ and } \lambda_0 \in \mathcal{A}(\mathcal{X})\} + \text{pr}\{\lambda_0 \notin \mathcal{A}(\mathcal{X})\}. \end{aligned} \quad (4.119)$$

Let F denote the cumulative distribution function of $T_{\lambda_0}(\mathcal{X})$. Then $F\{T_{\lambda_0}(\mathcal{X})\}$ has the uniform distribution on $(0, 1)$, and hence $\text{pr}_{\mathcal{Y}|\mathcal{X}}\{T_{\lambda_0}(\mathcal{Y}) \geq T_{\lambda_0}(\mathcal{X})\}$, which is equal to $[1 - F\{T_{\lambda_0}(\mathcal{X})\}]$, has the uniform distribution on $(0, 1)$. Now,

$$\begin{aligned} \text{pr}_{\mathcal{X}}\{p^*(\mathcal{X}) \leq \alpha \text{ and } \lambda_0 \in \mathcal{A}(\mathcal{X})\} \\ = \text{pr}_{\mathcal{X}}[\alpha_1 + \sup_{\lambda \in \mathcal{A}(\mathcal{X})} \text{pr}_{\mathcal{Y}|\mathcal{X}}\{T_\lambda(\mathcal{Y}) \geq T_\lambda(\mathcal{X}) \mid \theta = (0 : \lambda)\}] \leq \alpha \text{ and } \lambda_0 \in \mathcal{A}(\mathcal{X})] \\ \leq \text{pr}_{\mathcal{X}}[\alpha_1 + \text{pr}_{\mathcal{Y}|\mathcal{X}}\{T_{\lambda_0}(\mathcal{Y}) \geq T_{\lambda_0}(\mathcal{X}) \mid \theta = (0 : \lambda_0)\}] \leq \alpha \\ = \alpha - \alpha_1. \end{aligned}$$

Now $\text{pr}(\text{Type I error}) = \text{pr}\{p^*(\mathcal{X}) \leq \alpha\} \leq \alpha$, by (4.119). ■

In the proof of Proposition 4.8.1 and that of other similar asymptotic results, it is often sufficient to restrict attention to neighborhoods of the true value that shrink at the rate of $n^{-1/2}$, more precisely, neighborhoods of the form $\{\theta : \sqrt{n}\|\theta - \theta_0\| \leq K\}$ for large K . The next lemma provides a justification.

Lemma 4.10.1 (1) Let $\Theta \subset \mathbb{R}^p$, $\theta_0 \in \Theta$, $n^{1/2}(T_n - \theta_0) = O_p(1)$ and $\Theta_{n,K} = \Theta \cap \{\theta : n^{1/2}\|\theta - \theta_0\| < K\}$. Then $\|T_n - \theta_0\| = \|T_n - \Theta_n - \theta_0\|$ with arbitrarily large probability for sufficiently large K and n . More precisely, given ϵ there exist $n_0(\epsilon)$ and $K_0(\epsilon)$ such that $\text{pr}\{\Pi(T_n, \Theta) = \Pi(T_n, \Theta_{n,K})\} > 1 - \epsilon$ for $n > n_0$ and $K > K_0$.

(2) Suppose that $\hat{\theta}$ is the maximizer of the objective function $R_n(\theta)$ over Θ . Suppose also that $\sqrt{n}(\hat{\theta} - \theta_0) = O_p(1)$. Then, given $\epsilon > 0$ there exists $n_0 > 0$ and $K_0 > 0$ such that

$$\sup_{\theta \in \Theta} R_n(\theta) = \sup_{\theta \in \Theta_{n,K}} R_n(\theta)$$

with arbitrarily large probability for $n > n_0$ and $K > K_0$.

Proof: Given $\epsilon > 0$ there exist $n_0(\epsilon)$ and $K_0(\epsilon)$ such that

$$\text{pr}(n^{1/2} \|T_n - \theta_0\| < K/2) > 1 - \epsilon, \text{ for } n > n_0 \text{ and } K > K_0.$$

Since $\theta_0 \in \Theta$, we have that

$$\|T_n - \theta_0\| \geq \|T_n - \Pi(T_n, \Theta)\|.$$

Therefore,

$$\|\Pi(T_n, \Theta) - \theta_0\| \leq \|\Pi(T_n, \Theta) - T_n\| + \|T_n - \theta_0\| \leq 2\|T_n - \theta_0\|.$$

Now, if $n^{1/2} \|T_n - \theta_0\| < K/2$ then $n^{1/2} \|\Pi(T_n, \Theta) - \theta_0\| < K$ and hence $\Pi(T_n, \Theta) = \Pi(T_n, \Theta_{n,K})$. Therefore,

$$\text{pr}\{\Pi(T_n, \Theta) = \Pi(T_n, \Theta_{n,K})\} \geq \text{pr}(n^{1/2} \|T_n - \theta_0\| < K/2) > 1 - \epsilon$$

for $n > n_0$ and $K > K_0$.
The proof of the second part is similar. ■

Proof of Proposition 4.8.1: First let us first establish a lemma that would be of independent interest.

Lemma 4.10.2 Suppose that $Y_n = O_p(1)$, $A \subset \mathbb{R}^p$, V is a positive definite matrix of order $p \times p$ and \hat{V} is a consistent estimator of V . Then

1. $\Pi_V(Y_n, A) = O_p(1)$.
2. $\|Y_n - A\|_{\hat{V}} = \|Y_n - A\|_V + o_p(1)$.
3. $\|\Pi_{\hat{V}}(Y_n, A)\|_{\hat{V}} = O_p(1)$.

Proof:

(1) The first part follows by triangular inequality. Let P be a fixed point in A , O be the origin, $Y_n = OQ$, and R be the point in A that is V -closest to Q . Then

$$\begin{aligned} \|Y_n - A\|_V &= QR \leq QP \leq QO + OP = O_p(1) \\ \|\Pi(Y_n, A)\|_V &= OR \leq OQ + QR = O_p(1). \end{aligned}$$

(2) For a given matrix A , let $\|A\|$ denote the matrix norm that is equal to the largest

eigenvalue of AA^T ; then an upper bound for $\|A\|$ is $\sum \sum |a_{ij}|$. Here we shall use the inequality

$$\|Ax\|_W \leq \|A\| \|x\|_W \leq \sum \sum |a_{ij}| \|x\|_W$$

for any x and any positive definite matrix W (for example, see Pryce (1973, p 101)). From this and the Cauchy-Schwartz inequality, we have $x^T Ax \leq \sum \sum |a_{ij}| \|x\|^2$. Since \hat{V} is consistent, it follows that $\|V - \hat{V}\| = o_p(1)$. Let

$$f_n(\theta) = (Y_n - \theta)^T \hat{V}^{-1} (Y_n - \theta) \quad \text{and} \quad g_n(\theta) = (Y_n - \theta)^T V^{-1} (Y_n - \theta).$$

Then

$$\begin{aligned} \sup_{\theta \in A} |f_n(\theta) - g_n(\theta)| &\leq \sup_{\theta \in A} (\|Y_n - \theta\|^2) (\|\hat{V}^{-1} - V^{-1}\|) \\ &\leq (\|Y_n - A\|_F^2) (\|\hat{V}^{-1} - V^{-1}\|) = O_p(1) o_p(1) = o_p(1). \end{aligned}$$

Note that, since f_n and g_n are nonnegative, we have

$$\inf f_n(\theta) \leq |f_n(\theta)| \leq |f_n(\theta) - g_n(\theta) + g_n(\theta)| \leq |f_n(\theta) - g_n(\theta)| + g_n(\theta).$$

Therefore, $\inf f_n(\theta) \leq \sup |f_n(\theta) - g_n(\theta)| + g_n(\theta)$, and hence

$$\inf f_n(\theta) \leq \sup |f_n(\theta) - g_n(\theta)| + \inf g_n(\theta).$$

$$\text{Thus, } \inf f_n(\theta) - \inf g_n(\theta) \leq \sup |f_n(\theta) - g_n(\theta)|. \text{ By symmetry}$$

$$\inf g_n(\theta) - \inf f_n(\theta) \leq \sup |f_n(\theta) - g_n(\theta)|.$$

Therefore, we have

$$\inf_{\theta \in A} f_n(\theta) - \inf_{\theta \in A} g_n(\theta) \leq \sup_{\theta \in A} |f_n(\theta) - g_n(\theta)| = o_p(1).$$

This completes the proof of part (2). Part (3) follows from the first two. ■

Proof of Proposition 4.7.3: The proof of this uses several results concerning Constraint Qualification (CQ) in nonlinear programming and optimization. The main steps are indicated here without attempting to define terms and concepts in nonlinear optimization that are used in this proof. We shall use the same notation as in Bazaraa et al. (1993, section 5.3); for convenience, the relevant notations therein are stated below:

$$\begin{aligned} T &= \{\mathbf{d} : \mathbf{d} = \lim_{k \rightarrow \infty} \lambda_k (\theta_k - \theta_0), \theta_k \in \Theta_1, \theta_k \rightarrow \theta_0\}; \\ G' &= \{\mathbf{d} : \mathbf{d}^T \nabla h_i(\theta_0) \geq 0 \quad \text{for } i \in J(\theta_0)\}, \\ H_0 &= \{\mathbf{d} : \mathbf{d}^T \nabla h_i(\theta_0) = 0 \quad \text{for } i = 1, \dots, \ell\}, \\ A &= \{\mathbf{d} : \exists \delta > 0, \exists f : \mathbb{R} \rightarrow \mathbb{R}^p \text{ such that } f(\lambda) \in \Theta_1 \text{ for } \lambda \in (0, \delta), \\ &\quad \mathbf{f}(\mathbf{0}) = \theta_0, \text{ and } \lim_{\lambda \rightarrow 0} \lambda^{-1}[f(\lambda) - \mathbf{f}(\mathbf{0})] = \mathbf{d}\}. \end{aligned}$$

Thus, T is the cone of tangents and A is the derived tangent cone. For these cones the following holds (Bazaraa et al. (1993, p 193)):

1. Closure of $A \subset T \subset G' \cap H_0$.
2. Kuhn-Tucker CQ : Closure of $A = G' \cap H_0$.
3. Cottle CQ \Rightarrow Kuhn-Tucker (KT) CQ
4. Cottle CQ is equivalent to Mangasarian-Fromowitz (MF) CQ provided Θ is open.

Since Θ is open and the MF-CQ is satisfied, it follows that Cottle-CQ is satisfied, and hence KT-CQ is satisfied. Therefore, it follows that Closure of $A = T = G' \cap H_0$, and hence Θ_1 is Chernoff regular and the tangent cone is $G' \cap H_0$. ■

Proof of Proposition 4.7.4 using Hausdorff distance definition

Let $\epsilon > 0$ be given and $\|\mathbf{y} - \theta_0\| < \epsilon/2$. We shall refer to Figure 4.9 for this proof. Let $\mathbf{y} = OA$ and $\Pi(A, \Theta)$ denote the point in Θ that is V -closest (i.e., closest with respect to $\|\cdot\|_V$) to A . Let $C = \Pi(A, \mathcal{A})$, $D = \Pi(A, \Theta)$, $G = \Pi(C, \Theta)$ and $F = \Pi(D, \mathcal{A})$. First, note that since $C \in \mathcal{A}$ and $D \in \Theta$,

$$CD \leq h\{\Theta \cap B_\epsilon, \mathcal{A} \cap B_\epsilon\} = o(\epsilon).$$

By triangle inequality, $AC \leq AD + CD$; squaring both sides and subtracting AD^2 , we have

$$AC^2 - AD^2 \leq CD^2 + 2AD \cdot CD \leq o(\epsilon^2) + 2 \cdot OA \cdot CD = o(\epsilon^2) + \epsilon o(\epsilon) = o(\epsilon^2).$$

Similarly, by considering $AD \leq AC + CD$ we have $AD^2 - AC^2 \leq o(\epsilon^2)$. Therefore, $AD^2 - AC^2 = o(\epsilon^2)$. ■

Proof of Proposition 4.7.4 without using Hausdorff distance

Let $\mathbf{y} = OA$ and $\Pi(A, \Theta)$ denote the point in Θ that is V -closest (i.e., closest with respect to $\|\cdot\|_V$) to A . Let $C = \Pi(A, \mathcal{A})$, $D = \Pi(A, \Theta)$, $G = \Pi(C, \Theta)$ and $F = \Pi(D, \mathcal{A})$.

Case 1: $AC > AD$: see the diagram on the right in Fig.4.9. Since $AD \leq OA$, we have that $OD/OA \leq (OA + AD)/OA \leq (OA + OA)/OA = 2$. Now by making use of the definition of “closest point”, we have that $0 \leq (AC - AD)/OA \leq 2(DF/OD)$. $OD/OA = DF/OA = (DF/OD)(OD/OA) \leq$

Case 2: $AC < AD$; see the diagram on the left in Fig.4.9. Since OAC is a right-angle triangle, we have $OC \leq OA$. Now, $0 \leq (AD - AC)/OA \leq (AG - AC)/OA \leq (AC + CG - AC)/OA \leq CG/OA \leq CG/OC$. From cases (1) and (2), we have that

$$|(AC - AD)|/OA \leq \max\{2(DF/OD), CG/OC\}.$$

Since $AC \leq OA$ and $AD \leq OA$, it follows that $(AC + AD)/OA \leq 2$. Now, $|(AC^2 - AD^2)|/OA^2 = \{|(AC - AD)|/OA\}\{|(AC + AD)|/OA\} \leq 2|(AC - AD)|/OA \leq 2 \max\{2DF/OD, CG/OC\}$.

Now, as $OA \rightarrow 0$, we have $OD \rightarrow 0$ and $OC \rightarrow 0$ because $OD \leq OA + AD \leq 2OA$. As $OD \rightarrow 0$ we have $DF/OD \rightarrow 0$ by condition (a) of the definition of approximating cone; as $OC \rightarrow 0$ we have $CG/OC \rightarrow 0$ by condition (b) of the definition of approximating cone. Therefore, as $OA \rightarrow 0$, we have that $\max\{2DF/OD, CG/OC\} \rightarrow 0$, $(|AC^2 - AD^2|)/OA^2 \rightarrow 0$ and this establishes the claim of the proposition. ■