

SIMPLIFICATION SEARCHES

6.1 Simplification for Conditional Prediction 208
 6.2 Causally Constrained Conditional Predictions 214
 6.3 Simplification for Control 217
 6.4 Conclusion 223

In the two previous chapters we have considered simplification searches that are intended to introduce into a data analysis uncertain prior information. Hypothesis-testing searches arise when more than one model or hypothesis receive positive a priori probability. Interpretive searches involve prior density functions that, although allocating zero probability to all but one hypothesis, do concentrate the prior probability in certain regions of the parameter space. In the case of hypothesis-testing searches the statistical testing selects among a set of hypotheses with no presumption that in a large sample one of the hypotheses will be favored. In contrast, interpretive searches recognize that in a sufficiently large sample the most general hypothesis will necessarily be favored. The intent is not to select among legitimately competing models but rather to "improve" the estimate of the parameters by using an a priori estimate when the data evidence is too weak to yield a reliable sample estimate.

In this chapter we discuss a variety of search that has yet another motivation: simplification. The most general models appropriate for inference with nonexperimental data are usually so cluttered with variables of an incidental nature that they are nearly impossible to comprehend directly. It is thus incumbent on the researcher to find vehicles for communication of his results. He might, for example, focus his discussion on a particular parameter of special interest or perhaps on a linear combination of parameters. Alternatively, the researcher might seek from the data an indication

of the "important" variables. We call this a simplification search. Thus the function of simplification search is not to ask if a restricted specification is true, nor to ask if a restricted specification might lead to better parameter estimates, but rather to ask if a restricted specification that is undeniably simpler and more easily understood is not also "significantly" inferior to the more general model for some hypothetical or real decisions. If it is, we reject the hypothesis that the benefits of the restriction outweigh the costs.

Formal analysis of simplification problems requires a precise definition of the costs and benefits of simplicity. The costs of simplicity may be assessed in the context of some hypothetical decision problems, but the benefits are likely to elude precise definition. Consequently we concentrate our formal attention on the cost side, but we first comment informally on the likely benefits from simplification.

Justifications for simplicity can usefully be divided into two categories. The first makes a "metaphysical" reference to the inherent simplicity of Nature, or at least to man's belief in such. The second category of justifications accepts a complex Nature but rests simplicity on the finiteness and fallibility of Man's perceptive and reasoning faculties. Briefly simplicity is preferred because "Nature is simple" or because "Man is simple."

The "Nature is simple" hypothesis has, I think, little support among philosophers and statisticians. Jeffreys' is a widely cited exception. He writes [1961, p. 4] "It is asserted, for instance, that the choice of the simplest law is purely a matter of economy of description or thought, and has nothing to do with any reason for believing the law... I say, on the contrary, that the simplest law is chosen because it is the most likely to give correct predictions; that the choice is based on a reasonable degree of belief;..."

Jeffreys is asserting not only that constrained hypotheses should be assigned positive probability but also that they ought to be assigned greater prior probability than any alternative, more complex hypotheses. Such a preference for simple models might be inductively derived. Simple hypotheses could usually yield better predictions. But it is not enough to observe merely that people act as if simple models had a greater degree of believability. Any observed preference for simple models may derive not from the inherent superior believability of parsimonious models but rather from the undeniable difficulties encountered in working with complex descriptions of reality. Nor do I know of any proper empirical evidence to support the assertion that simpler models generally yield better predictions. There is the oft-told story of overfitting in which a naive researcher fits a polynomial of degree $T-1$ given T pairs of observations (y_i, x_i) . This

yield inferior predictions relative to a polynomial of degree n . But that can be fully remedied by assigning a proper weight to the parameters of the higher-degree polynomial. I can use as an illustration of the illogic of using a prior that is not supported by the data evidence when the data evidence simply contradicts it. It is hardly evidence in favor of simpler models. I discuss in the chapter on hypothesis testing that I know of few models that would assign positive probability to a restricted (simple) hypothesis. I can find nothing that compels me to favor the simpler hypothesis of probability. It is the other set of reasons for finding a hypothesis persuasive: Simplicity is desirable because it is easier to transmit and accumulation of knowledge. It greatly facilitates communication between and among observers and theorists. A theory that might take years to filter accurately to other observers is transmitted rapidly (but inaccurately) if it is simplified. At the very least, it is likely to be superior knowledge for their own use. It is likely to engage in this kind of marketing activity. Many models are likely to be used to advertise it.

I have argued in various ways that simplicity encourages the use of simpler models (1972) because they are more easily understood. This would be true if the simple model were assigned a higher probability, but if such a model is derived from a more complex hypothesis, it is likely to be rejected. It is apparent that falsifying evidence can be taken to mean only that a hypothesis does not work under all conditions. In that case, it is likely to be rejected. On the other hand, protection of a hypothesis from potential falsification is an essential feature of a theory (according to Kuhn (1962)). Filling in the details of a theory that all its implications requires a vast amount of tedious work would hardly be performed by doubters or even agnostics. Nature is simple" nor the "Man is simple" hypothesis. Ambiguous definitions or methods of measuring the benefits of a number of uncertain parameters is a possible mechanical simplicity, but it cannot be generally satisfactory. If we take (as a consequence of man's and society's shortcomings, the complexity necessarily changes from social milieu to social milieu) as impossible and even undesirable to define simplicity, we instead must content ourselves with the satisfaction that in any social information process can know themselves and what it is not.

The prototypical example of this is the construction of a map (Polan 1964). We may take as a theory of the world an enormously detailed globe which identifies every object down to the smallest grain of sand. The complexity of this theory effectively prevents us from using it for a purpose whatsoever. Instead, we simplify it in the form of a set of maps. We use one map to find my way to the subway station, another to select the station at which to depart. The pilot of the airplane uses yet another map to navigate from Boston to Washington. Each map is a greatly simplified version of the theory of the world; each is designed for some class of decisions and works relatively poorly for others.

The construction of a language is another good example of a simplification problem. The number of aurally and visually discernible words and word patterns is absolutely enormous, perhaps limitless. With as few as 26 characters as are in our alphabet we could form $26^5 > 10^7$ distinct five-letter words. Such a vocabulary would be beyond the reach of even the most verbally talented, and the mistaken use of words used infrequently would greatly distort intended communications. A highly limited vocabulary likewise distorts communications by not distinguishing one communication from another, for example, the American overuse of the word "nice" to describe a wide variety of generally pleasing responses to environmental stimuli. An optimal vocabulary ideally solves the trade-off between miscommunications from too few words and miscommunications from too many.

Incidentally, there is a great danger that a simple language is not only a vehicle for communication but that it also creates an impoverished reality of its own. The art of communication forces an awareness of reality, and the more subtle is the language, the more practice one obtains in distinguishing subtleties. Conversely, a coarse language creates no situation for exercising one's capacities to distinguish subtleties, and those faculties may atrophy like any unused muscle. We may, in fact, be unable anymore to distinguish the great variety of sensations we refer to as "nice." This may also be the case in the communication of scientific theories. We may come erroneously to believe in the simplicity of Nature because that is the way scientific theories are communicated.

I do not think it is possible to define simplicity, which is to say in the language of decision theory that it is difficult to compute precise benefits or precise costs from any simplification. In this chapter the cost of simplification is measured in the context of several simple decision problems, but the benefits are not quantified at all. We hope that what we learn can have implications for more complex and more realistic decisions.

One thing that is important to understand is that simplification is a decision problem which uses as an *input* the current information about the parameters. When a current sample is available, simplification logical

These three examples illustrate, respectively, a prediction problem, a control problem, and an inference problem. Relative to a model of the form $y = z\gamma + w\delta + u$ they ask if we may act as if δ were zero (1) if we wanted to predict y , (2) if we wanted to control y , or (3) if we wanted to make inferences about γ . The inference problem is distinguished from the others only in that more data is to be gathered before decisions are to be made. The actual, ultimate decision may, in fact, be either a prediction or a control problem. This is called a presimplification problem, referring to the fact that simplification occurs prior to observation. We make much use of the presimplification notion in Chapter 9, when we discuss postdata model construction.

It is easy to demonstrate the inappropriateness of classical hypothesis testing at a fixed level of significance for the simplification problem. Suppose the prior distribution were diffuse. The only information conveyed by the fact that the hypothesis $\gamma = 0$ is or is not rejected at the 5% level of significance is the information that the posterior 95% credible interval includes or does not include the point $\gamma = 0$. Thus you may reject the hypothesis $\gamma = 0$ even though with near certainty γ is infinitesimal (Figure 6.1a). And you may accept the hypothesis even though with high probability γ is enormous. (Figure 6.1b) It is thus important to distinguish the words "statistically significant" from the words "economically signif-

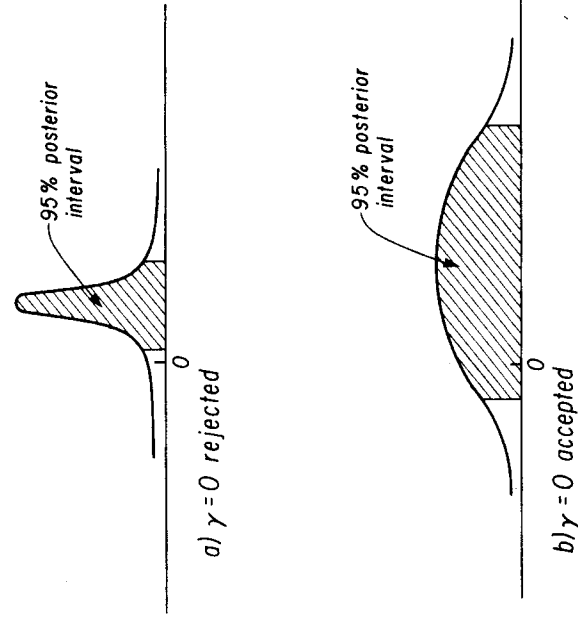


Fig. 6.1 Posterior distributions.

ATION SEARCHES

and is confused with the inferential process, at great peril of a statistical analysis. I would recommend making as far as possible between inference and decision, by discussions of a research report first the inferential question of distributions are influenced by the data and second the of how given various posterior distributions the model l. Incidentally, since a data set is taken as given, any nts reported in this chapter are necessarily conditional ; thus notationally convenient to suppress the data when al probability statements, and it is hoped that this will not For example, the statement $E(\beta) = (X'X)^{-1}X'Y$ implies nal mean of β , $E(\beta|Y, X)$, is equal to the least-squares $X'Y$.

merits repeating is that a simplified model that might y for some decision-making circumstances will be unam- ptable in others. It is, therefore, essential to identify problem that is considered. Three examples suggest the diversity of decision problems.

regate consumption of apples C_a and aggregate consump- C_b depend on aggregate GNP Y through the functions and $C_b = \alpha_a + \beta_a Y$. If we wish to predict future levels of apples and bananas, may we without great detriment to nstrain the marginal propensities to consume to equal b_a and therefore "remember" only the marginal propensity rather than separate propensities for each fruit?

IP Y is thought to depend on the government deficit G ock M , $Y = \alpha + \beta G + \gamma M$. If we wish GNP to attain some re effectively assure that goal by selecting an appropriate rument deficit G^* while treating money M as if it had no onversely, might we better control money M and act as if

constant-elasticity-of-substitution production function ex- ; a function of capital and labor inputs. If the elasticity of equal to one, the investment function assuming profit-maxi- is a function of one explanatory variable rather than two. ation generated by observing the production process, may ces about the investment process acting as if certain of its on special values?

measures the amount of information in the data; the size of the coefficient in the context of some decision joint is that classical tests have built into them rather unwarranted assumptions about the behavior of the variables. Consider again the model $y = z\gamma + w\delta + u$ with δ and w exactly, and with the explanatory variables z and w in a linear relationship $w = rz + \varepsilon$, where u and ε are independent variables and r is known. The hypothesis $\delta = 0$ can be used to test, yielding either $H_0: y = z\gamma$ or $H'_0: y = z(\gamma + r\delta)$, where r allows the included variable to play partly the role of the

excluded variable. The only goal, the hypothesis H'_0 is unambiguously more restrictive than H_0 and yields a lower expected loss. But for other reasons H_0 may be preferred. A simplification is intended to facilitate communication. It may be difficult to communicate, since it seems to say "we are testing $\delta = 0$ on y is $(\gamma + r\delta)$ when, in fact, it is only γ . It is *least* to distinguish the hypotheses "we may act as if $\delta = 0$ " from "we may act as if w_T were zero" from the hypothesis "we are testing $\delta = 0$ for not observing w_T ," the former pair implying the latter and the latter implying H'_0 . Classical hypothesis testing is a special case of H_0 with r implicitly estimated in a special case of the general form of simplification is

discussed in Section 6.2. The other form of simplification is discussed in this chapter consists of three sections and a conclusion. We report Lindley's (1968) formal decision-theoretic approach to the prediction-simplification problem. Among the lessons to be learned is the great importance of assumptions about the process that generates the explanatory variables. In fact, the simplification problem is not more on the process that generates these variables than on the regression process linking the dependent variable to the explanatory variables. That observation is used in Section 6.2 to argue in favor of simplification that makes fewer demands on our knowledge of the explanatory variable process and that also communicates the dependence of the decision process on the decision process. The third section emphasizes the dependence of the decision process on the decision process under consideration by using the (1968) analysis of a control-simplification problem and a solution with the prediction-simplification problem.

Simplification for Conditional Prediction

In a conditional-prediction problem, consider the two-variable model

$$y = \alpha + \gamma z + w_T \delta + u \tag{6.1}$$

where α , δ and γ are unobservable scalar parameters, u_t is an unobservable error, and y_t , z_t , and w_t are observable variables. Suppose that u_t ($t = 0, \dots, T$) is a sequence of independent normal random variables with zero means and known variance σ^2 . Let a set of T previous observations of the process be (Y, z, w) , which together with a multivariate prior distribution for the parameters (α, γ, δ) imply a multivariate posterior distribution with mean

$$E([\alpha, \gamma, \delta] | Y, z, w) = [\bar{\alpha}, \bar{\gamma}, \bar{\delta}]$$

In making a conditional prediction of the next outcome, say, Y_T , we assume that both the explanatory variables z_T and w_T are potentially observable prior to the announcement of the prediction, hence the adjective "conditional" modifying prediction. It is perhaps obvious, but it is demonstrated here that if the penalty for prediction error is quadratic, the optimal prediction given both z_T and w_T is

$$\hat{y}_T = \bar{\alpha} + z_T \bar{\gamma} + w_T \bar{\delta} \tag{6.2}$$

where $\bar{\gamma}$ and $\bar{\delta}$ are the posterior means of γ and δ . There will, of course, be prediction errors, partly because of the residual error process u_t and partly because the actual values of the parameters α , γ , and δ are not known.

Suppose, now, that we wished to determine if it is worth the expense to observe the second variable w_T . If w_T is not observed, we must estimate it by say, \hat{w}_T , and predict y_T as a function of z_T only as

$$\hat{y}_T^* = \bar{\alpha} + z_T \bar{\gamma} + \hat{w}_T \bar{\delta} \tag{6.3}$$

The squared discrepancy between Equations (6.2) and (6.3) is a measure of the error induced by not observing w_T :

$$(\hat{y}_T - \hat{y}_T^*)^2 = (w_T - \hat{w}_T)^2 \bar{\delta}^2 \tag{6.4}$$

Note especially that this error depends on the mean of $\bar{\delta}$ but not on its variance. Note also that in testing the hypothesis $\delta = 0$ in the sense of this chapter, that is, by computing numbers like (6.4), we are partly asking the question "is δ small?" but more importantly we are asking also "how well can we forecast w_T ?" To answer the latter question, we must model the process that generates w_T and z_T —there is no way the simplification question can be answered without such a model.

One model (that should, I think, be of little interest to an economist operating with time-series data) is the multivariate random model, in which the explanatory variables are treated as if they were drawn randomly from a population with fixed mean vector and covariance matrix. In particular

matrix

$$\mathbf{V} = \begin{bmatrix} v_{11} & v_{12} \\ v_{21} & v_{22} \end{bmatrix}.$$

arameters of this distribution, a prediction of w_T would be conditional regression

$$\hat{w}_T = E(w_T | z_T) = \mu_w + v_{21}v_{11}^{-1}(z_T - \mu_z)$$

variance

$$E[(w_T - \hat{w}_T)^2 | z_T] = v_{22} - v_{21}v_{11}^{-1}v_{12}.$$

arts of these unknown parameters can be used if the prior diffuse and if the number of observations is large.¹ We

$$E(w_T | z_T, \mathbf{z}, \mathbf{w}) = \bar{w} + (\mathbf{w}'\mathbf{Mz})(\mathbf{z}'\mathbf{Mz})^{-1}(z_T - \bar{z}) \quad (6.5)$$

$$-\hat{w}_T)^2 | z_T, \mathbf{z}, \mathbf{w}] = \frac{\mathbf{w}'\mathbf{Mw} - (\mathbf{w}'\mathbf{Mz})(\mathbf{z}'\mathbf{Mz})^{-1}(\mathbf{z}'\mathbf{Mw})}{T} \quad (6.6)$$

are the sample means of \mathbf{w} and \mathbf{z} and \mathbf{M} is the matrix that $\mathbf{M} = \mathbf{I} - \mathbf{1}_T\mathbf{1}'_T$. The predicting equations (6.3) and error (6.4) thus become

$$\bar{w} - (\mathbf{w}'\mathbf{Mz})(\mathbf{z}'\mathbf{Mz})^{-1}\bar{z} + z_T(\bar{y} + (\mathbf{w}'\mathbf{Mz})(\mathbf{z}'\mathbf{Mz})^{-1}\bar{\delta}) \quad (6.7)$$

$$= E(w_T - \hat{w}_T)^2 \delta^2 = \frac{[\mathbf{w}'\mathbf{Mw} - \mathbf{w}'\mathbf{Mz}(\mathbf{z}'\mathbf{Mz})^{-1}\mathbf{z}'\mathbf{Mw}]\delta^2}{T} \quad (6.8)$$

ns may now be made. If the prior for γ and δ were diffuse, eans $\bar{\alpha}$, $\bar{\delta}$, and \bar{y} would be just the least-squares estimates, b_z . The coefficient of z_T in Equation (6.7) would then be b_z , which is just the estimated coefficient of a regression alone. Furthermore, the penalty (6.8) can be written as χ^2 is the χ -square value for testing the restriction $\delta = 0$,

$$\chi^2 = \frac{b_z^2 [\mathbf{w}'\mathbf{Mw} - \mathbf{w}'\mathbf{Mz}(\mathbf{z}'\mathbf{Mz})^{-1}\mathbf{z}'\mathbf{Mw}]}{\sigma^2}.$$

ature just described measures the increase in the expected r when w_T is not observed in terms of the usual χ^2 variable 0. As is discussed in detail subsequently, it differs from thesis testing in implicitly defining the significance level as a rial from Section 3.4, and the diffuse prior assumption with $T^* = 0$, $\mathbf{S}^* = \mathbf{0}$, riance of w_T given z_T is not (6.6) but rather (6.6) times the adjustment

decreasing function of the sample size. What is perhaps more important is the fact that the decision theory logic makes unambiguous the otherwise implicit assumptions about the process that generates the explanatory variables. In particular, classical tests are appropriate only if the explanatory variable vectors are independently drawn from the same population.

Let us now repeat this logic for a general model and for general linear restrictions. Write the linear regression process as

$$\begin{bmatrix} \mathbf{Y} \\ y_T \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{x}'_T \end{bmatrix} \boldsymbol{\beta} + \begin{bmatrix} \mathbf{u} \\ u_T \end{bmatrix}$$

where \mathbf{Y} and \mathbf{X} are $(T \times 1)$ and $(T \times k)$ matrices and are already observed, where y_T is a future outcome of the process and \mathbf{x}_T is a $k \times 1$ vector of future explanatory variables, and where $[\mathbf{u}', u'_T]$ is a $(1 \times (T+1))$ vector of errors with mean zero and covariance $\boldsymbol{\Sigma}$. We are asked to predict y_T given \mathbf{Y} , \mathbf{X} , and \mathbf{x}_T and in particular to minimize squared prediction error $[y_T - \hat{y}(\mathbf{Y}, \mathbf{X}, \mathbf{x}_T)]^2$ with prediction \hat{y} . The expected prediction error can be written as

$$E[y_T - \hat{y}(\mathbf{Y}, \mathbf{X}, \mathbf{x}_T)]^2 = E\left\{E\left([y_T - \hat{y}(\mathbf{Y}, \mathbf{X}, \mathbf{x}_T)]^2 | \mathbf{Y}, \mathbf{X}, \mathbf{x}_T\right)\right\},$$

where the expression in the internal brackets is straightforwardly minimized for every value of $(\mathbf{Y}, \mathbf{X}, \mathbf{x}_T)$ by setting

$$\begin{aligned} \hat{y}(\mathbf{Y}, \mathbf{X}, \mathbf{x}_T) &= E(y_T | \mathbf{Y}, \mathbf{X}, \mathbf{x}_T) = \mathbf{x}'_T E(\boldsymbol{\beta} | \mathbf{Y}, \mathbf{X}, \mathbf{x}_T) + E(u_T | \mathbf{Y}, \mathbf{X}, \mathbf{x}_T) \\ &= \mathbf{x}'_T E(\boldsymbol{\beta} | \mathbf{Y}, \mathbf{X}) + E(u_T | \mathbf{Y}, \mathbf{X}) \end{aligned}$$

which is a linear function of \mathbf{x}_T . If some part of \mathbf{x}_T is not observed, we assume that the complete vector is predicted as a linear function of that which is observed. That is, letting $\mathbf{x}'_T = (\mathbf{x}'_T, \mathbf{x}'_T')$, we assume that $E(\mathbf{x}_T | \mathbf{Y}, \mathbf{X}, \mathbf{x}'_T) = \mathbf{A}\mathbf{x}'_T$, and thus the optimal predicting equation becomes

$$\hat{y}(\mathbf{Y}, \mathbf{X}, \mathbf{x}'_T) = \mathbf{x}'_T \mathbf{A}' E(\boldsymbol{\beta} | \mathbf{Y}, \mathbf{X}) + E(u_T | \mathbf{Y}, \mathbf{X}).$$

Or, to make a long story short, we wish to restrict our attention to predictions linear in \mathbf{x}_T

$$\hat{y}(\mathbf{Y}, \mathbf{X}, \mathbf{x}_T) = \mathbf{x}'_T \boldsymbol{\theta}(\mathbf{Y}, \mathbf{X}) \quad (6.9)$$

where the function $\boldsymbol{\theta}$ may be completely free, in which case it is just the posterior mean of $\boldsymbol{\beta}$, or it may be constrained to have certain elements zero to reflect the fact that certain elements of the vector \mathbf{x}_T are not observed prior to the prediction of y_T . Incidentally, Equation 6.9 implicitly includes the $E(u_T | \mathbf{Y}, \mathbf{X})$ term, since \mathbf{x}'_T is assumed to have one element equal to one.

For ease of notation we write the conditional expected value operator $E(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{X})$ henceforth as just $E(\cdot)$. If \mathbf{Y} and \mathbf{X} are given, $\boldsymbol{\theta}$ is just a vector

the expected loss can be written as

$$\begin{aligned} & \beta'x_T + u_T - \theta'x_T)^2 \\ & r(\beta - \theta)(\beta' - \theta')x_T \\ & r(\beta - E\beta + E\beta - \theta)(\beta - E\beta + E\beta - \theta)'x_T \\ & r(\beta - E\beta)(\beta - E\beta)'x_T + x_T'(E\beta - \theta)(E\beta - \theta)'x_T] \end{aligned} \quad (6.10)$$

$(x_T)V(\beta) + (E\beta - \theta)'S(x_T)(E\beta - \theta)$ written $S(x_T) = Ex_Tx_T'$. The three terms in the last line of the irreducible mean-square error Euw_T^2 , a penalty for and an additional penalty for $\theta \neq E\beta$, this last term of the uncertainty in β .

pected loss if θ lies in the linear subspace $R\theta = r$ is ed posterior loss (6.10) minimized over that linear sub- nization is a simple Lagrangian problem requiring the

$$= (E\beta - \theta)'S(x_T)(E\beta - \theta) + 2\lambda'(R\theta - r)$$

That is,

$$\frac{\partial f}{\partial \lambda} = R\theta - r = 0 \quad (6.11)$$

$$\frac{\partial f}{\partial \theta} = -S(x_T)(E\beta - \theta) + R'\lambda = 0. \quad (6.12)$$

ed by premultiplying (6.12) by $RS^{-1}(x_T)$ and calculating

$$\lambda = (RS^{-1}(x_T)R')^{-1}(RE(\beta) - r) \quad (6.13)$$

the mean-square error (6.10) becomes

$$(E\beta - \theta) = (RE(\beta) - r)'(RS^{-1}(x_T)R')^{-1}(RE(\beta) - r). \quad (6.14)$$

om the positive definiteness of the third term in (6.10) ected posterior loss requires $\theta = E\beta$, or by (6.14) that a es expected loss unless $R\theta = r$. A simplification thus ases expected prediction accuracy. We assume that a benefits also, and in the absence of any clear quantita- those benefits, a reasonable number to report is the

percentage increase in the expected posterior loss due to the restriction

$R\theta = r$:

$$L^2(R, r) = \frac{(RE(\beta) - r)'(RS^{-1}(x_T)R')^{-1}(RE(\beta) - r)}{Euw_T^2 + \text{tr}S(x_T)V(\beta)}. \quad (6.15)$$

With suitable definitions of prior vagueness we have simply the least-squares results (remembering that the expected value operator is conditional on X and Y)

$$\begin{aligned} E(\beta|Y, X) &= (X'X)^{-1}X'Y \\ V(\beta|Y, X) &= \sigma^2(X'X)^{-1}. \end{aligned}$$

Further, if the explanatory variables are independent observations from a multivariate process, we would have the x_T moment matrix be approximately (see Section 3.4)

$$S(x_T) = E(x_Tx_T') = \frac{X'X}{T}.$$

Using these in (6.13), θ is seen to be simply the constrained least-squares estimate subject to $R\theta = r$. Inserting them into (6.14), we obtain the increase in the posterior expected loss to be T^{-1} times a factor that is well known to be the increase in the error-sum squares due to the restriction. The summary L^2 becomes

$$L^2(R, r) = \frac{T^{-1}\Delta ESS}{\sigma^2\left(1 + \frac{k}{T}\right)} \quad (6.16)$$

where ΔESS is the increase in the error sum of squares, k is the number of coefficients, and T is the number of observations. This contrasts with the classical summary statistic $\Delta ESS/\sigma^2$, which is compared with $\chi_p^2(\alpha)$ where p is the rank of R and α is the significance level. Thus the classical counterpart of (6.16) is the ratio $\Delta ESS/\sigma^2\chi_p^2(\alpha)$. In addition to the nonoccurrence of the factor T^{-1} (which for large T necessitates a "significant" finding), the classical summary differs from the subjectivist summary in depending on p , the number of restrictions. The measure (6.16), incidentally, is just the difference in the multiple correlation coefficients of the two models times a factor that tends to a constant as sample size grows, $L^2(R, r) = (R^2 - R_0^2)(Y'MY/T\sigma^2)/(1 + kT^{-1})$. Thus if a restriction does not greatly affect the R^2 of an equation, it will not greatly increase the expected squared prediction error.

This rough coincidence of approaches usefully highlights the assumptions that are implicit in the use of classical tests to simplify models for

the dependent variable. Indirect effects depend on other unspecified causal linkages.

An example is in order to make clear these relatively obscure notions. Suppose the equation of motion of a body falling from rest is

$$\frac{d^2y}{(dt)^2} = g(1 - \beta zt)$$

where z measures the wind resistance and t the time since departure. The parameter g is acceleration in a vacuum, and terminal velocity is reached at time $t = 1/\beta z$. Suppose, further, that observations on a set of falling bodies are used to estimate the equation of location

$$y = \frac{\hat{g}t^2}{2} - \frac{\hat{\beta}\hat{g}zt^3}{6}, \quad R^2 = .98$$

where the circumflexes indicate estimated parameters. For this particular sample of falling bodies (including feathers and bowling balls) the following auxiliary regression is also calculated

$$(zt^3) = \hat{r}t^2 - \hat{\alpha}.$$

The model may be simplified to exclude the wind resistance variable. Two alternative simpler models are

$$y = \hat{g}t^2/2 \quad R^2 = .70 \quad (6.17)$$

$$y = \frac{(\hat{g} - \hat{\beta}\hat{g}\hat{r}/3)t^2}{2} + \frac{\hat{\alpha}\hat{\beta}\hat{g}}{6} \quad R^2 = .95 \quad (6.18)$$

It is my contention that the first of these equations is the one that should be used to discuss simplification. It asserts that in a vacuum the estimate rate of acceleration is \hat{g} and that for the class of bodies and for the time periods considered, we ought not to think of the experiment as if it were conducted in a vacuum, since one's ability to predict the location of falling bodies is seriously affected by that assumption. (The R^2 drops from .98 to .7.) Contrast that perfectly clear statement with the statement appropriate for the second equation. "Wind resistance is 'negligible' since by adjusting the rate of acceleration to $\hat{g} - \hat{\beta}\hat{g}\hat{r}/3$ and by acting as if the initial location of the body were $\hat{\alpha}\hat{\beta}\hat{g}/6$ rather than zero, we can track the position of this class of falling bodies almost as well as we would if we actually observed the wind resistance."

In fact, wind resistance is not negligible; rather, it can be compensated for. At the very least we ought to make clear the distinction between the two statements. For reasons I have explained, I think simplification more appropriately interpreted as the problem of neglecting variables

CATION SEARCHES

course, there is the diffuse prior assumption. But more vectors of explanatory variables are assumed to be $T+1$ indications of a multivariate process. Autocorrelation and similar are assumed away. Few economists would find that also worth stating explicitly that the variance in the t statistic does not measure the uncertainty in the rather the inverse of the conditional variance of the explanatory. From (6.10), it is seen that uncertainty in the β) does not influence choice of restrictions $\theta \neq E(\beta)$.

Causally Constrained Conditional Predictions

spect of the solution discussed in the previous section is variables are used to forecast correlated unobserved variables; assumption that the correlation structure is maintained. effect (the coefficient) of an observed variable thus includes an estimated coefficient but also a part due to the effect of variables assumed to be correlated with it. Interpreted in terms of R^2 due to a restriction is calculated. restricted equation with a reestimated set of coefficients. may make good sense if we intend the test to determine the of the restriction, it makes less sense for the simplification we really mean to say that an effect of an explanatory variable when it can be predicted well from observation of a variable? This is the question implicit in a classical statistical application of Webster suggests that a variable is negligible if we can neglect it without substantial loss. means not bothering to predict it or otherwise to make it not observing it. As will be shown, this is the question of statistical beta coefficients and variants thereof.

from semantics to metaphysics, we can find another same argument. To the extent that the full unconstrained theory our beliefs about the causal nature of the world, the of the coefficients implicit in hypothesis testing constitutes a notion of that causality. That is, since included variables play role of dropped variables, the constrained equation is causally less the included variables do, in fact, cause the excluded they do not, the resulting equation is causally inaccurate. An ide toward the causality within the explanatory variable set is reporting the original estimates of the coefficients of the variables calculated in the context of the unconstrained equation. described as the direct effects of the included variables on

Given the assumption of diffuse priors, and supposing that δ is a scalar, the criterion (6.19) is just the square of the least-squares coefficient times the sample variance of the variable. If this were divided by the square of the sample variance of the dependent variable, the resulting number would be just the beta coefficient, which can be computed as least squares with variables standardized to have unit variance. Although standardized coefficients are used in other disciplines, in the econometrics literature they are rarely even mentioned. Goldberger (1964, pp. 197-198) is an exception.

To conclude, criterion (6.20), which is equivalent to (6.16) under diffuseness assumptions, ranks variables considered individually for discarding in the same way as traditional t tests. Criterion (6.19), however, provides a ranking identical to the ranking implied by classical beta coefficients. It seems to me, therefore, that the rarely used beta coefficients could be usefully resurrected as indicators of significance when models are being simplified, although the variance of the explanatory variables ought at a minimum be trend and autocorrelated adjusted.

6.3 Simplification for Control

A point that may be obvious is that simplification is problem specific, and, for example, simplification for prediction may be quite different from simplification for control. The one-period control problem of Lindley (1968) illustrates this fact. Suppose a scalar variable y_T is determined by the linear-regression process

$$y_T = \alpha + \gamma'z_T + \delta'w_T + u_T \tag{6.21}$$

where γ and δ are vector parameters, α is a scalar parameter, u_T is a residual error with mean zero and variance σ^2 , and z_T and w_T are vectors of explanatory variables. The control problem is to select z_T and w_T in such a way that y_T is likely to be close to some target t . In particular, let us choose the explanatory variables to minimize expected loss where loss is quadratic

$$L(y_T, t) = (y_T - t)^2.$$

Writing the regression process as

$$y_T = \alpha + \beta'x_T + u_T \tag{6.22}$$

where $\beta' = [\gamma, \delta']$ and $x_T = [z_T, w_T']$, the expected loss can be written as a function of x_T as

$$E(L(y_T, t)|x_T) = E([\alpha + \beta'x_T + u_T - t]^2|x_T).$$

Setting the derivatives of this expression to zero to obtain the minimizing

CAUTION SEARCHES

problem of compensating for their effects. There is first the point that if simplification were intended to compensate for neglect certain secondary influences, we might expect a more appropriate adjective than "negligible." Second, for a secondary influence, the theory may be fundamentally distorted. Consider the gravity example. If we assume we assert what is completely true: a body falling in vacuum accelerates at the constant rate \hat{g} . Contrast that with results when wind resistance is compensated for: a body at the time of departure instantaneously falls to a height initial position, attaining thereby absolutely no velocity, falls, accelerating at the constant rate $\hat{g} - \beta\hat{g}/3$. The Newtonian mechanics is obvious and absurd.

tests with unrecomputed coefficients can be calculated formulas as tests with recomputed coefficients provided the constraint matrices R and r appropriately. If we write $= \alpha + z_T'\gamma + w_T'\delta + u_T$, a simplification hypothesis is $\delta = 0$. recomputation of the coefficients on the z variables by constraint that the coefficients must equal their posterior as a causally constrained simplification is implied by the

$$R = \begin{bmatrix} 0 & I & 0 \\ 0 & 0 & I \end{bmatrix} \quad r = \begin{bmatrix} E(\gamma) \\ 0 \end{bmatrix}$$

column of R is a vector of zeroes multiplying the constant α restriction matrices the mean-square-error penalty (6.14)

$$RE(\beta) - r'(RS^{-1}(x_T)R)^{-1}(RE(\beta) - r) \\ E(\delta)]'V(w_T)[E(\delta)]. \tag{6.19}$$

bles without the causal constraint requires constraint

$$R = \begin{bmatrix} 0 & 0 & I \\ 0 & 0 & 0 \end{bmatrix}, \quad r = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \\ \text{quare-error penalty (6.14) becomes} \\ [E(\delta)]'V(w_T|z_T)[E(\delta)] \tag{6.20}$$

is the conditional variance of w_T , given z_T . Penalty (6.20) penalty (6.19) depending on the correlation between the included variables, because the included variables are used ided variables.

$$\mathbf{0} = E[2\beta\beta'x_T + 2\beta(\alpha - t)]$$

$$x_T = (E\beta\beta')^{-1}E\beta(t - \alpha).$$

alve of x_T into the expected loss, we obtain the minimum

$$\begin{aligned} & \text{in } E[L(y_T, t)|x_T] \\ & + E(\alpha - t)^2 - E(t - \alpha)\beta'(E\beta\beta')^{-1}E\beta(t - \alpha). \end{aligned} \quad (6.23)$$

for the expected loss simplifies nicely in the case when α and β derives only from observation of the regression y . Letting Y be the T -dimensional vector of previous ne process and X be the matrix of observations of the bles, the posterior moments are

$$\begin{aligned} E(\alpha) &= \bar{Y} - \bar{X}'(X'MX)^{-1}X'MY = b_0 \\ E(\beta) &= (X'MX)^{-1}X'MY = \mathbf{b} \end{aligned}$$

dimensional vector of ones and $M = I - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}'$, $\bar{X} =$ Also, the variance matrix can be written as

$$\begin{bmatrix} \mathbf{1}'\mathbf{1} & \mathbf{1}'X \\ X'\mathbf{1} & X'X \end{bmatrix}^{-1} = \begin{bmatrix} T^{-1}(1 + \mathbf{1}'X(X'MX)^{-1}\bar{X}) & -\bar{X}'(X'MX)^{-1} \\ -(\bar{X}'MX)^{-1}\bar{X} & (X'MX)^{-1} \end{bmatrix}.$$

$\bar{Y} = b_0 + \bar{X}'\mathbf{b}$ we may write the regression process as

$$\begin{aligned} &= \alpha + \beta'x_T - b_0 - \mathbf{b}'\bar{X} + u_T \\ &= (\alpha - b_0) + (\beta - \mathbf{b})'\bar{X} + \beta'(x_T - \bar{X}) + u_T \\ &\equiv \alpha^* + \beta'x_T^* + u_T \end{aligned} \quad (6.24)$$

r and if α and β were known to equal $E\alpha$ and $E\beta$, then the instrument $E\alpha)/E\beta$, which is called the certainty equivalence control rule. Assum- ent, the optimal rule can be written in terms of the certainty equivalence $+ E^2(\beta\beta')^{-1}E(\beta)(t - E(\alpha)) = (1 + t\beta^{-2})^{-1}x_T^*$, where $t\beta^2 = E^2(\beta)/V(\beta)$. is more conservative than the certainty equivalence rule in the sense ble is not turned on as far. The shrinkage factor $(1 + t\beta^{-2})^{-1}$ is a function

where

$$\begin{aligned} \alpha^* &= (\alpha - b_0) + (\beta - \mathbf{b})'\bar{X} \\ x_T^* &= x_T - \bar{X}. \end{aligned}$$

Controlling y_T at t is equivalent to controlling $y_T^* = y_T - \bar{Y}$ at $t^* = t - \bar{Y}$, where y_T^* is generated by the process described in (6.24). The expected loss (6.23) attains a simple form since $E(\beta\alpha^*) = \mathbf{0}$

$$\begin{aligned} L_1 &= E(\alpha^* - t^*)^2 - t^{*2}\mathbf{b}'(\mathbf{b}\mathbf{b}' + \sigma^2(X'MX)^{-1})^{-1}\mathbf{b} + \sigma^2 \\ &= E\alpha^{*2} + t^{*2}\left(1 - \mathbf{b}'(\mathbf{b}\mathbf{b}' + \sigma^2(X'MX)^{-1})^{-1}\mathbf{b}\right) + \sigma^2 \\ &= \frac{\sigma^2}{T} + \sigma^2 + \frac{t^{*2}}{1 + \mathbf{b}'X'MX\mathbf{b}/\sigma^2} \\ &= \sigma^2(1 + T^{-1}) + \frac{t^{*2}}{1 + \chi^2} \end{aligned} \quad (6.25)$$

where we have used the inverse formula $(xx' + A)^{-1} = A^{-1} - A^{-1}x(x'Ax' + A^{-1})^{-1}x'A^{-1}$.

Thus the minimum expected loss is a quadratic function of the deviation of the target from the historical level of the process, $(t - \bar{Y})^2 = t^{*2}$. The coefficient multiplying this term is $(1 + \chi^2)^{-1}$ where χ^2 is the value of the chi-square statistic for testing $\beta = \mathbf{0}$. A large χ^2 statistic thus implies that y_T can be pushed from its historical mean without incurring great expected loss. The part of the expected loss independent of the target is just the variance of y_T assuming that x_T is set to its historical level \bar{X} ,

$$V(y_T|x_T = \bar{X}) = E\alpha^{*2} + \sigma^2 = \sigma^2(1 + T^{-1}).$$

Next consider the possibility that none of the variables is controlled. To compute expected control error it is then necessary to "guess" what the explanatory variables will be. This means modeling the process that generates the explanatory variables. For our purposes it is enough to know the first two moments of x_T , since the expected loss can be written as

$$\begin{aligned} E(y_T - t)^2 &= E(\alpha + \beta'x_T + u_T - t)^2 \\ &= \sigma^2 + E(\alpha - t)^2 + 2E(\alpha - t)\beta'x_T + E\beta'x_Tx_T'\beta. \end{aligned}$$

Taking as we did in the previous section the assumption of an independent multivariate process for the explanatory variables, we have approximately $E\mathbf{x}_T = \bar{X} = X'\mathbf{1}/T$, $V\mathbf{x}_T = X'MX/T$, where $M = I - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}'$. These together with the least squares moments for α and β impute in the above

$$\begin{aligned}
 t)^2 &= E(y_T - \bar{Y} + \bar{Y} - t)^2 \\
 &= E(y_T - \bar{Y})^2 + (\bar{Y} - t)^2 \\
 &= \sigma^2 + E(\alpha^* + \beta'x_T^*)^2 + t^{*2} \\
 &= \sigma^2 + E\alpha^{*2} + E \operatorname{tr}(x_T^*x_T^{*'}\beta\beta') + t^{*2} \\
 &= \sigma^2 + \frac{\sigma^2}{T} + \frac{k\sigma^2}{T} + \frac{\chi^2\sigma^2}{T} + t^{*2} \\
 &= \sigma^2 \left(1 + \frac{k+1+\chi^2}{T} \right) + t^{*2}
 \end{aligned} \tag{6.26}$$

nessionality of β and χ^2 is the chi-square value for testing

), the expected loss with no control, is to be contrasted

6.25), the expected loss with optimal control. Their dif-

ue to decontrolling x_T is then

$$\frac{\sigma^2(k + \chi^2)}{T} + \frac{t^{*2}\chi^2}{1 + \chi^2}, \tag{6.27}$$

minimum of $k/(T+1)$ when $\chi^2=0$. We are thus led to

with $T+1$ to determine if decontrolling x_T could be

ase expected losses substantially.

r hand, it is desired to control y_T at some value far from

in, the second term in (6.27) dominates the expected loss.

increase in expected loss would then be just χ^2 , and we

mpare χ^2 with the number one to decide if controlling x_T

examined the extreme cases in which either all or none of

he vector $x_T^*=(z_T', w_T')$ is under control. The intermediate

control affects only z_T is more difficult, since it requires a

; how z_T affects the distribution of w_T or, more accurately,

he conditional distribution $f(y_T|z_T)$. Both the prediction

ion 6.1 and the control problem of this section are most

the conditional moments of y_T :

$$\begin{aligned}
 E(y_T|z_T) &= Ey_T + g'(z_T - Ez_T) \\
 V(y_T|z_T) &= a + (z_T - Ez_T)'A(z_T - Ez_T).
 \end{aligned} \tag{6.29}$$

These assumptions—that the mean is a linear function and that the variance is a quadratic function of z_T —are implicit in the foregoing discussion. The prediction problem of minimizing $E(y_T - \hat{y})^2$ where \hat{y} is a function of z_T is straightforwardly solved by letting $\hat{y} = E(y_T|z_T)$ with resultant expected loss $E(y_T - \hat{y})^2 = E[y_T - E(y_T|z_T)]^2 = EV(y_T|z_T) = a + \operatorname{tr}AV(z_T)$.

The control problem is equally trivial. We wish to choose z_T to minimize

$$\begin{aligned}
 \min_{z_T} E[(y_T - t)^2|z_T] &= \min_{z_T} E\{[(y_T - E(y_T|z_T))]^2|z_T\} + [t - E(y_T|z_T)]^2 \\
 &= \min_{z_T} V(y_T|z_T) + [t - E(y_T|z_T)]^2.
 \end{aligned}$$

With the foregoing moments the derivatives of this expression with respect to z_T are

$$2A(z_T - Ez_T) - 2g(t - Ey_T - g'[z_T - Ez_T])$$

which when set to zero yields the optimizing value of z_T

$$z_T^* = Ez_T + (A + gg')^{-1}g(t - Ey_T).$$

The resulting expected loss is

$$\begin{aligned}
 E[(y_T - t)^2|z_T = z_T^*] &= a + (t - Ey_T)g'(A + gg')^{-1}A(A + gg')^{-1}g(t - Ey_T) \\
 &\quad + [t - Ey_T - g'(A + gg')^{-1}g(t - Ey_T)]^2 \\
 &= a + (t - Ey_T)^2[1 - g'(A + gg')^{-1}g] \\
 &= a + \frac{(t - Ey_T)^2}{1 + g'A^{-1}g}
 \end{aligned} \tag{6.30}$$

Note that this is a quadratic function of $(t - Ey_T)$, the discrepancy between the target and the expected value of y_T .

To be specific, let us again work with the diffuse prior assumption. After some minor manipulation, we may obtain for the constants in the moments (6.29) the following

$$\begin{aligned}
 E(y_T) &= \bar{Y}, & E(z_T) &= \bar{Z} \\
 g &= (Z'MZ)^{-1}(Z'MY) & & \text{(the regression of } Z \text{ on } Y) \\
 A &= \sigma^2(Z'MZ)^{-1} & &
 \end{aligned}$$

$Y'MW - W'MZ)(Z'MZ)^{-1}Z'MW)b_w/\sigma^2$. Notice that the chi-square statistic for testing $\gamma=0$, given that $\delta=0$, is then

$$\min_z E[(y_T - t)^2 | z_T] = \frac{\chi_\delta^2 \sigma^2}{T} + \frac{k_w \sigma^2}{T} + \frac{t^2}{1 + \chi_{\gamma|\delta=0}^2}$$

and the historical mean, $t^{*2}=0$, the percentage increase in $\chi_{\gamma|\delta=0}^2$ is not controlled is thus

$$\frac{T^{-1}(\chi_\delta^2 + k_w)}{1 + T^{-1}} = \frac{\chi_\delta^2 + k_w}{1 + T}$$

compare $\chi_\delta^2 + k_w$ with $(1 + T)$ to determine if w_T can be little increase in expected error. from the historical mean the percentage increase in $\chi_{\gamma|\delta=0}^2$ is due to decontrolling w_T is

$$\frac{(1 + \chi_{\gamma|\delta=0}^2)^{-1} - (1 + \chi_\delta^2)^{-1}}{(1 + \chi_\delta^2)^{-1}}$$

$$= \frac{\chi_\delta^2 - \chi_{\gamma|\delta=0}^2}{1 + \chi_{\gamma|\delta=0}^2} = \frac{\chi_\delta^2}{1 + \chi_{\gamma|\delta=0}^2}$$

compare χ_δ^2 with $1 + \chi_{\gamma|\delta=0}^2$. repeated that these results involve the unlikely assumption that z_T we do not alter the process that generates the variables (in the sense that the conditional distribution is altered). The assumption of known σ^2 can be altered by prior mean where relevant. Mathematically more appropriate treat the vector (y_T, x'_T) as coming from a multivariate normal with unknown mean and unknown variance matrix. A prior or the uncertain parameters implies that the marginal distribution (y_T, x'_T) is a multivariate Student distribution with means μ and Σ (6.28) and (6.29). We leave to the tenacious reader the calculation.

6.4 Conclusion

To conclude we may restate, first, the more important formal results of this chapter and then reiterate the more important informal lessons to be learned.

The results of this chapter listed in Table 6.2 make use of the assumptions listed in Table 6.1. If a variable y is generated by a linear regression process with explanatory variables w and z , if w and z themselves come from a multivariate normal process, and if priors for the various parameters are appropriately diffuse, then: (1) for a conditional prediction problem, we need not observe w , if the χ^2 statistic for testing whether w can be omitted (χ_δ^2) is small relative to $(T+k)$ where T is the number of observations and k is the dimension of $x'=(w',z')$; (2) for control with target equal to the historical mean of y , w may be decontrolled if $\chi_\delta^2 + k_w$ is small relative to $(1+T)$, where k_w is the dimension of w ; (3) for control far from the historical mean, w may be decontrolled if χ_δ^2 is small relative to $1 + \chi_{\gamma|\delta=0}^2$, one plus the χ^2 value for testing if z belongs in the equation given that w does not.

The principal caveat that has been repeated ad nauseam is that these results involve a very specific and often unwarranted assumption about the

Table 6.1

Assumptions for Simplification Analysis

Model

$$y_t = \alpha + z'_t \gamma + w'_t \delta + u_t \\ \equiv \alpha + x'_t \beta + u_t, \quad t=0, 1, \dots, T \\ u_t \sim N(0, \sigma^2), \quad \sigma^2 \text{ known} \\ x_t \sim N(\mu, \Sigma) \\ \mu, \Sigma, \alpha, \beta \text{ have diffuse priors}$$

Observations

$$Y(T \times 1), Z(T \times k_z), W(T \times k_w), \\ X=(Z, W)(T \times k_x)$$

Statistics

$$b=(X'MX)^{-1}X'MY, \quad M=I-1T^{-1}1' \\ \bar{Y}=1'Y/T, \bar{X}=1'X/T \\ b_0=\bar{Y}-\bar{X}b \\ g=(Z'MZ)^{-1}Z'MY, \quad \hat{b}_0=\bar{Y}-\bar{Z}'g \\ \chi_\delta^2 = b'_w [W'MW - W'MZ(Z'MZ)^{-1}Z'MW]b_w / \sigma^2 \\ \chi_{\gamma|\delta}^2 = b'X'MXb / \sigma^2 \\ \chi_{\gamma|\delta}^2 = g'Z'MZg / \sigma^2$$

process that generates the explanatory variables. No simplification decisions can be made without either an implicit or explicit study of the behavior of the explanatory variables, and we hardly need say that it seems clear that an explicit study of their behavior is highly desirable.

For both prediction and control problems the effects of the excluded variables have been compensated for by adjustment of the included variables, and we have argued at length that it may be desirable not to adjust in this way. Semantically, adjustment is undesirable, because rather than asking if a variable can be neglected, in fact, we ask if it can be compensated for. Metaphysically, adjustment is undesirable, since it implies a causal link between the included and excluded variables. Statistically, the predictions and control that result may be quite inferior if anything happens to change the historical correlations between the variables. Control, especially, is likely to alter those correlations.

Decision Rules		Expected Losses	
D_1 : unconstrained		$L_0 - L_1$	
Prediction	$y_T = b_0 + z_T b_z + w_T b_w$	$\sigma^2(1 + T^{-1}k_x)$	$T^{-1}\sigma^2 k_x^2$
Causally Constrained	same as above	same as above	$T^{-1}b_w^2 W^2 M^2 b_w^2$
Prediction	$x_T = \bar{X} - (bb' + \sigma^2[X'MX]^{-1})^{-1}b_1^* z_T = \bar{z} + (g'g + \sigma^2[Z'MZ]^{-1})^{-1}g_1^*$	$\sigma^2(1 + T^{-1})$	$\sigma^2(\chi_z^2 + k_w)T^{-1}$
Control	\bar{X}	$\sigma^2(1 + T^{-1})$	$\sigma^2(\chi_z^2 T^{-1} + k_w)$
Control (approx.)	same as two lines above	same as two lines above	
t^* very large			
Control (approx.)	same as two lines above	$t^{*2}(1 + \chi_z^2)^{-1}$	$t^{*2}[(1 + \chi_z^2)^{-1} - (1 + \chi_z^2)^{-1}]$