



Model Selection Strategies

• In this lecture we will address assumptions (A1)-(A5). In particular, how do we propose and select a model (a DGP)?

• Potentially, we have a huge number of possible models (different functional form, *f(.)*, and explanatory variables, **X**). Say, we have

Model 1 $y = X\beta + \varepsilon$ Model 2 $y = Z\gamma + \xi$ Model 3 $y = (W\gamma)^{\lambda} + \eta$ Model 4 $y = \exp(Z D \delta) + \epsilon$

• We want to select the best model, the one that is closest to the DGP. In practice, we aim for a good model.

Model Selection Strategies

• A model is a simplification. Many approaches:

• "Pre-eminence of theory." Economic theory should drive a model. Data is only used to quantify theory. Econometric methods offer sophisticated ways 'to bring data into line' with a particular theory.

• Purely data driven models. Success of ARIMA models (late 60s – early 70s). No theory, only exploiting the time-series characteristics of the data to build models.

• Modern (LSE) view. A compromise: theory and the characteristics of the data are used to build a model.

Model Selection Strategies

• Modern view: Theory and practice play a role in deriving a good model. David Hendry (2009) emphasizes:

"This implication is not a tract for mindless modeling of data in the absence of economic analysis, but instead suggests formulating more general initial models that embed the available economic theory as a special case, consistent with our knowledge of the institutional framework, historical record, and the data properties. ... Applied econometrics cannot be conducted without an economic theoretical framework to guide its endeavours and help interpret its findings. Nevertheless, since economic theory is not complete, correct, and immutable, and never will be, one also cannot justify an insistence on deriving empirical models from theory alone."

Model Selection Strategies According to David Hendry, a good model should be:

- Data admissible -i.e., modeled and observed **y** should have the same properties.
- Theory consistent -our model should "make sense"
- Predictive valid -we should expect out-of-sample validation
- Data coherent -all information should be in the model.
 - Nothing left in the errors (*white noise errors*).
- Encompassing -our model should explain earlier models.

• That is, we are searching for a statistical model that can generate the observed data (\mathbf{y}, \mathbf{X}) , this is usually referred as *statistical adequacy*, makes theoretical sense and can explain other findings.

Model Selection Strategies

- FAQ in practice:
 - Should I include all the variables in the database in my model?
 - How many explanatory variables do I need in my model?
 - How many models do I need to estimate?
 - What functional form should I be using?
 - Should the model allow for structural breaks?
 - Should I include dummies & interactive dummies ?
 - Which regression model will work best and how do I arrive at it?

Model Selection Strategies: Some Concepts

• *Diagnostic testing:* We test assumptions behind the model. In our case, assumptions (A1)-(A5) in the CLM.

Example: Test $E[\mathbf{\epsilon} | \mathbf{X}] = 0$ -i.e., the residuals are zero-mean, white noise distributed errors.

• Parameter testing: We test economic H_0 's. **Example**: Test $\beta_k = 0$ -say, there is no size effect on the expected return equation.

- There are several *model-selection methods*. We will consider two:
 - Specific to General
 - General to Specific

Model Selection Strategies: Two Methods

• There are several *model-selection methods*. We will consider two:

- Specific to General
- General to Specific

• Specific to General. Start with a small "restricted model," do some testing and make model bigger model in the direction indicated by the tests (for example, add variable x_k when test reject H₀: $\beta_k=0$).

• General to Specific. Start with a big "general unrestricted model," do some testing and reduce model in the direction indicated by the tests (for example, eliminate variable x_k when test cannot reject H₀: β_k =0).

• Begin with a small theoretical model	– for example, the CAPM			
$y = X\beta + \varepsilon.$				
• Estimate the model	– say, using OLS			
• Do some diagnostic testing	– are residuals white noise			
• If the assumptions do not hold, then us	se:			
- More advanced econometrics	- GLS instead of OLS?			
- A more general model	– APT? Lags?			
• Test economic H_0 on the parameters	- Is size significant?			
• This strategy is known as <i>specific to genera</i> (AER), and, in the machine learning litera	ıl, Average Economic Regression ature, forwards selection.			





Model Selection Strategies: Specific to General

• The specific-to-general method makes assumptions along the way Some remarks based on the previous example:

(1) Very likely the starting model is based on theory and experience (HML is not significant at the usual 5% level). Not clear how to proceed from there to a more general model.

(2) We tested for a January effect and then added to the model. However, we could have tested for a Dot.com effect or for an interactive Dot.com/January effect with the 3 FF factors. Not clear when to stop the search.

(3) Select a p-value to add variables to the model. In this case, we use the standard 5% for the tests.

Model Selection Strategies: Specific to General

• Note that in the previous example, we started with a model. What happens if are skeptical regarding models?

• A popular implementation of the specific-to-general model selection is the *stepwise regression*, where we start with only a set of potential explanatory variables and let the data, based on some criteria (R², AIC, etc.), determine which variables to keep.

Model Selection Strategies: Stepwise Regression

• Overall structure:

- The method begins with a k potential regressors.

- Do k one-variable regressions. Pick the one that shows the biggest tstat or maximizes a goodness of fit measure, say, Adjusted-R², \overline{R}^2 . Suppose x_j is selected.

- Then, do (k - 1)-variable regressions all with x_j . Select the regressor (in addition to x_j) that has the highest t-stat or that maximizes \overline{R}^2 .

- Continue. But, when we start adding regressors, we usually check if the added regressor(s) change the significance of previous steps. (Note: at each step, we remove or add a regressor(s) based on t- or F-tests.)

- Stop: Additional regressors do not have *significant* t-stats/increase \overline{R}^2 .

• Decisions: Selection of *k* variables, α for tests ($\alpha = 5\%$, 10%, 20%?) and goodness of fit statistic.

Model Selection Strategies: Stepwise Regression

• Overall structure:

- The method begins with a *k* potential regressors.

- Do k one-variable regressions. Pick the one that shows the biggest tstat or maximizes a goodness of fit measure, say, Adjusted-R², \overline{R}^2 . Suppose x_i is selected.

- Then, do (k - 1)-variable regressions all with x_j . Select the regressor (in addition to x_j) that has the highest t-stat or that maximizes \overline{R}^2 .

- Continue. But, when we start adding regressors, we usually check if the added regressor(s) change the significance of previous steps. (Note: at each step, we remove or add a regressor(s) based on t- or F-tests.)

- Stop: Additional regressors do not have *significant* t-stats/increase \overline{R}^2 .

• <u>Decisions</u>: Selection of *k* variables, α for tests ($\alpha = 5\%$, 10%, 20%?) and goodness of fit statistic.

Model Selection Strategies: Stepwise Regression

• Decisions: Selection of k initial variables, α for tests ($\alpha = 5\%$, 10%, 30%?) and goodness of fit statistic.

<u>Remark</u>: Always keep in mind that the selected (final) model is not necessarily better than others. Type I and Type II errors are likely to occur, thus the final model may have irrelevant and/or omitted variables.

<u>Technical Note</u>: Though popular in practice, in general, selecting variables based on *p-values* is not advised, since the distribution of the OLS coefficients is affected. (Recall pre-testing.)

Example with the 5 ols_step_fo	: Step 5 FF fa rward_	wise regr actors as <i>p</i> in the <i>a</i>	ession s candidat olsrr pacl	trategy tes for sage, v	y to m IBM. vhich	odel II We us uses <i>p</i> -	BM returns. We se the function <i>values</i> to select:
library(olsrr) ff_step_data	<- data	.frame(Mkt	_RF, SMB.	HML, F	NMW, C	CMA)	
.bm_ff_mod	el <- lm	(ibm_x ~ .,	$data = ff_{ata}$	step_dat	a) ‡	≠ default	p-value (penter) is 0.2
ols_step_for	ward_p((ibm_ff_mo	del , details	s = TRU	Ъ) 7	≠ long fi	nal output
Pa	rameter	Estimates					
Pa model	rameter Beta	Estimates Std. Error	Std. Beta	t	Sig	lower	upper
Pa model (Intercept)	Beta -0.005	Estimates Std. Error 0.002	Std. Beta	t -1.999	Sig 0.046	lower -0.010	upper 0.000
Pa model (Intercept) Mkt_RF	rameter Beta -0.005 0.887	Estimates Std. Error 0.002 0.055	Std. Beta 	t -1.999 16.227	Sig 0.046 0.000	lower -0.010 0.780	upper 0.000 0.995
Pa model (Intercept) Mkt_RF SMB	Beta -0.005 0.887 -0.261	Estimates Std. Error 0.002 0.055 0.088	Std. Beta 0.574 -0.111	t -1.999 16.227 -2.960	Sig 0.046 0.000 0.003	-0.010 0.780 -0.435	upper 0.000 0.995 -0.088

© RS 2024 – Not to be posted/shared online without written consent from author

Exa	ample (continua	ation):				
		Selection	Summary				
	Variable		Adj.				
Step	Entered	R-Square	R-Square	C(p)	AIC	RMSE	
1	Mkt_RF	0.3087	0.3075	7.7108	-1665.5551	0.0594	
2	SMB	0.3174	0.3151	2.2117	-1671.0548	0.0590	
3	RMW	0.3188	0.3154	2.9552	-1670.3207	0.0590	



Model Selection Strategies: General to Specific

• General-to-Specific Method:

Step 1 - First ensure that the GUM does not suffer from any diagnostic problems. Check residuals in the GUM to ensure that they possess acceptable properties. (For example, test for heteroskedasticity, white noise, incorrect functional form, etc.).

Step 2 - Test the restrictions implied by the specific model against the general model – either by exclusion tests or other tests of linear restrictions.

Step 3 - If the restricted model is accepted, test its residuals to ensure that this more specific model is still acceptable on diagnostic grounds.

• This strategy is called *general to specifics* ("gets"), LSE, TTT (Test, test, test), and, in the ML literature, *backwards selection*. It was pioneered by Sargan (1964). The properties of gets are discussed in Hendry and Krolzig (2005, Economic Journal).



Model Selection Strategies: Properties

• A modeling strategy is *consistent* if its probability of finding the true model tends to 1 as *T* -the sample size- increases.

- Properties for strategies
- (1) Specific to General
- It is not consistent if the original model is incorrect.
- It need not be predictive valid, data coherent, & encompassing.
- No clear stopping point for an unordered search.
- (2) General to Specific
- It is consistent under some circumstances. But, it needs a large T.
- It uses data mining, which can lead to incorrect models for small T.

- The significance levels are incorrect. This is the problem of *mass significance*.



Model Selection Strategies: General to Specific							
Examp	le (con	tinuati	on):				
fit_ibm_gum <- lm (ibm_x ~ Mkt_RF + SMB + HML + Jan_1 + Mkt_RF_2 + SMB_2 + HML_2 + Mtt_HML + Mtt_SMB + SMB_HML + Mtt_Iaa + HML_Iaa + Mtt_Dat + SMB_Dat							
<pre>Mkt_HML + Mkt_SMD + SMD_HML + Mkt_Jan + HML_Jan + Mkt_Dot + HML_Dot + SMD_Dot) > summary(fit ibm_gum)</pre>							
Coefficients:							
	Estimate	Std. Error	r t value	e Pr(> t)			
(Intercept)	-0.007836	0.003063	-2.559	0.010772	*		
Mkt_RF	0.791866	0.090474	8.752	< 2e-16 *	k**		
SMB	-0.295790	0.110655	-2.673	0.007738	**		
HML	-0.233942	0.135146	-1.731	0.084004		\Rightarrow practice says	"keep it." Judgement call.
Jan_1	0.031769	0.009349	3.398	0.000727 *	***		
Mkt_RF_2	-0.433762	0.850899	-0.510	0.610417			
SMB_2	-0.927271	1.470645	-0.631	0.528615			
HML_2	2.707992	1.670366	1.621	0.105545		\Rightarrow almost 10%, 1	I keep it. Judgement call.
Mkt_HML	0.628721	1.557090	0.404	0.686531			
Mkt_SMB	0.791625	1.746939	0.453	0.650618			
SMB_HML	-1.044806	2.029091	-0.515	0.606819			
Mkt_Jan	-0.069413	0.189309	-0.367	0.714008			
HML_Jan	-0.259697	0.255484	-1.016	0.309841			







Model Selection Strategies: General to Specific

Example (continuation):

Step 2 - Further specification checks of Restricted GUM, for example,
perform a Ramsey's reset test (using the *resettest* in the lmtest library).
> resettest(fit_gum_r, type="fitted")

RESET test

data: **fit_ibm_gum_r** RESET = **1.0998**, df1 = 2, df2 = 561, p-value = **0.3337**

Step 3 - Test if Restricted GUM residuals are acceptable –i.e., do diagnostic tests (mainly, make sure they are white noise). If Restricted GUM passes all the diagnostic tests, it becomes the "final model."

<u>Note</u>: With the final model, we use it to justify/explain financial theory and features, and do forecasting.

Model Selection Strategies: General to Specific

• The general-to-specific method makes assumptions along the way. Some remarks based on the previous example:

(1) Select a *p-value* for the tests of significance in the discovery stage (we use 10%). Given that we performed 15 *t-tests*, we should not be surprised we rejected the GUM, since we had an overall significance, $\alpha^* = .79 [= 1 - (1 - .10)^{15}]$. *Mass significance* is an issue.

(2) Judgement calls are also made.

(3) The reduction of the GUM involves "*pre-testing*" –i.e., data mining. We are likely rejecting a true H_0 (false positives) & not rejecting a true H_1 (false negatives), along the way. This increases the probability that the final model is not a good approximation. It is common to ignore (or not even acknowledge) pre-testing issues.

Model Selection Strategies: Best Subset

• Begin with a big model, with *k* regressors:

$$y = X\beta + \varepsilon$$

The idea is to select the "best" subset of the k regressors in **X**, where "best" is defined by the researcher, say MSE, Adjusted-R², etc.

- In theory, it requires 2^k regressions. It can take a while if k is big (k < 40 is no problem).
- Many tricks are used to reduce the number of regressions.

• In practice, we use best subset to reduce the number of models to consider. For example, from the regressions with one-variable, keep the best one-variable model, from the regression with two-variables, keep the best two-variable model, etc.



Issues: Pre-testing • A special case of omitted variables. - First, a researcher starts with an unrestricted model (U): $\mathbf{v} = \mathbf{X}\mathbf{\beta} + \mathbf{\varepsilon}.$ (U)- Then, based on ("preliminary") tests --say, an F-test- a researcher decides to use restricted estimator. That is, $\mathbf{y} = \mathbf{X}\mathbf{\beta} + \mathbf{\varepsilon}.$ s.t. $\mathbf{R}\boldsymbol{\beta} = \mathbf{q}$ (R) - We can think of the estimator we get from estimating R as: $\mathbf{b}_{\rm PT} = \mathbf{I}_{\{0,c\}}(F) \ \mathbf{b}^* + \mathbf{I}_{\{c,\infty\}}(F) \ \mathbf{b},$ where $I_{\{0,c\}}$ is an indicator function: $I_{\{0,c\}}(F) = 1$, if *F*-stat not in the rejection region – say, F < c – $I_{\{c,\infty\}}(F) = 0$, otherwise. *c*: critical value chosen for testing H_0 : $R\beta = q$, using the *F*-stat.

Issues: Pre-testing

• The *pre-test estimator* is a rule which chooses between the restricted estimator, **b***, or the OLS estimator, **b**:

 $\mathbf{b}_{\mathrm{PT}} = \mathrm{I}_{\{0,\mathbf{c}\}}(F) \ \mathbf{b}^* + \mathrm{I}_{\{\mathbf{c},\infty\}}(F) \ \mathbf{b}.$

where $b^* = b - (X'X)^{-1}R'[R(X'X)^{-1}R']^{-1}(Rb-q)$

• Two "negative" situations:

(1) H_0 : **R** β = **q** is true. The *F*-test will incorrectly reject $H_0 \alpha\%$ of the time. That is, in $\alpha\%$ of the repeated samples, OLS **b** \Rightarrow No bias, inefficient estimator.

(2) $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{q}$ is false. The *F-test* will correctly reject H_0 a % of times equal to the power π of the test. That is, $(100 - \pi)$ % of the time, $\mathbf{R}\boldsymbol{\beta}=\mathbf{q}$ will be incorrectly imposed, \mathbf{b}^* will be used \Rightarrow bias!

Issues: Pre-testing

• The failure of the OLS estimator to have the properties under correct specification is called *pre-test bias*.

• Pre-testing (also called *sequential estimation, data mining*) is common in practice. In general, it is ignored –and not even acknowledged.

• Main argument to ignore pre-testing: We need some assumptions to decide which variables are included in a model. Is the probability that pre-testing yields an incorrect set of *X* greater than the probability of selecting the "correct" assumption?

• The LSE methodology does not see pre-testing in the discovery stage as a problem. For the LSE method, pre-testing at that stage is part of the process of discovery.

Issues: Pre-testing

 \bullet Checking the MSE of $b_{\mbox{\scriptsize PT}},\ b^*$ and b helps to evaluate the problem

• <u>Practical advise</u>: Be aware of the problem. Do not rely solely on stats to select a model –use economic theory as well.

• Do not use same sample evidence to generate an H₀ and to test it!

Example: The Fama-French factors have been "discovered" using the CRSP/Compustast database for a long, long time. Thus, testing the Fama-French factors using the CRSP/Compustat is not advisable! (You can test them with another dataset, for example, get international data.)

Issues: Mass significance

• Two cases (1) Independent tests $\alpha^* = 1 - (1 - \alpha)^k$ & $\alpha = 1 - (1 - \alpha^*)^{1/k}$ (2) Dependent tests: $\alpha^* \le k\alpha$ & $\alpha \ge \alpha^*/k$ \Rightarrow close to the "independent" values for small α , but can differ for large α . **Example**: $\alpha = 0.05$ and $k=5 \Rightarrow \alpha^*(\text{Indep}) = .23$ & $\alpha^*(\text{Dep}) = .25$ $\alpha = 0.05$ and $k=20 \Rightarrow \alpha^*(\text{Indep}) = .64$ & $\alpha^*(\text{Dep}) = 1$ $\alpha^* = 0.05$ and $k=5 \Rightarrow \alpha(\text{Indep}) = .0102$ & $\alpha(\text{Dep}) = .01$ $\alpha^* = 0.20$ and $k=5 \Rightarrow \alpha(\text{Indep}) = .044$ & $\alpha(\text{Dep}) = .04$ $\alpha^* = 0.20$ and $k=20 \Rightarrow \alpha(\text{Indep}) = .011$ & $\alpha(\text{Dep}) = .01$

Issues: Mass significance In repeated *parametric testing* (overall level 5%): Only accept variables as important when their *p-values* are less than 0.001, preferably smaller Maybe look for other ways of choosing variables, say IC. In repeated *diagnostic testing* (overall level 20%), we should only accept there is no misspecification if All *p-values* are greater than 0.05, *or*Most *p-values* are greater than 0.10 with a few in the range 0.02 to 0.10

Modeling Strategies: Information Criteria • IC's are equal to the estimated variance or the log-likelihood function plus a penalty factor, that depends on k. Many IC's: • Theil Information Criterion (Adjusted R²) \overline{R}^2) $\overline{R}^2 = 1 - [(T-1)/(T-k)](1 - R^2) = 1 - [(T-1)/(T-k)]$ RSS/TSS \Rightarrow maximizing Adjusted R² \Leftrightarrow minimizing s² • Akaike Information Criterion (AIC) AIC = $-2/T(\ln L - k) = -2 \ln L/T + 2 k/T$ \Rightarrow if normality AIC = $\ln(e^*e/T) + (2/T) k$ (+constants) • Bayes-Schwarz Information Criterion (BIC) BIC = $-(2/T \ln L - [\ln(T)/T] k)$ \Rightarrow if normality AIC = $\ln(e^*e/T) + [\ln(T)/T] k$ (+constants)

Modeling Strategies: Information Criteria

• The goal of these criteria is to provide us with an easy way of comparing alternative model specifications, by ranking them.

<u>General Rule</u>: The lower the IC, the better the model. For the previous IC's, then choose model to minimize $s_p^2 AIC_p$, or BIC_p .

- Some remarks about IC's:
- They are used for ranking. The raw value tends to be ignored.
- They have two components: a *goodness of fit* component –based on *lnL* and a model complexity component –the penalty based on *k*.
- Different penalties, different IC's.

- Some authors do not scale the IC's by *T*, like we do above. If raw values are irrelevant, this is not an issue.

Modeling Strategies: Information Criteria

• IC's are not test statistics. They do not test a model. But, they are statistics –i.e., they are functions of RVs- with sampling distributions.

• We would like these statistics –i.e., the IC's– to have good properties. For example, if the true model is being considered among many, we'd want the information criteria to select it. This can be done on average (unbiased) or as *T* increases (consistent).

• Usually, inconsistency is a fatal flaw for a statistics. But, in model selection, it is very likely that the true DGP is not among the models considered. That is, inconsistency may not matter in these cases.

• Information? It refers to Kullback and Leibler's (1951) *information discrepancy* measure, used in information theory (in telecom literature).

Modeling Strategies: IC - K-L divergence

• Kullback and Leibler's (1951) *information discrepancy* measure is also called *information divergence*.

• Information divergence measures the difference between two probability distributions P and Q; where P represents the true DGP. Here, we look at the difference between the expected values of *Y* when *Y* is determined by: (i) P and (ii) some Q model.

• Minimizing the K-L divergence, when considering several Q models, gets us close to the true DGP.

• But, expected values are unobservable, they need to be estimated. The information associated with Y is given by L-i.e., the joint pdf. The AIC uses ln L evaluated at the estimated parameter values.

Modeling Strategies: IC - AIC and BIC

• Some results regarding AIC and BIC.

- AIC and Adjusted R² are not consistent.

- AIC is conservative –i.e., it tends to over-fit; that is, choose too large models.

- AIC selects the model that minimizes the leave-one-out crossvalidation MSE for cross-sectional data. In time series, it selects the model that minimizes the out-of-sample one-step ahead forecast MSE.

- BIC is more parsimonious than AIC. It penalizes the inclusion of parameters more ($k_{\text{BIC}} \le k_{\text{AIC}}$).

- BIC is consistent in hierarchical (gets) autoregressive models.

Modeling Strategies: IC - AIC and BIC

• There are several *small sample corrections* of IC's. But, asymptotically they have no impact. Because of this feature, using corrected IC's is not a bad choice.

• Comparing models based on IC's can be expensive

• In 'unstructured problems' (natural order to the hypotheses to be tested), there is a huge number of potential combinations to investigate: 2^m possible models for *m* candidate variables.

• For the Lovell (1983) database, that would be $2^{40} \approx 10^{12}$ models. Even at a USD 0.001 per model, that would cost USD 1 billion.

Modeling Strategies: Other Criteria

• A related criteria is Mallows' (1973) C_p statistic (Notation: p = k): $C_p = RSS(k)/s^2 - T + 2 * k$

where RSS(k) is the RSS for the model with k regressors.

It can be shown that the C_p -statistic estimates the size of the bias that is introduced into the predicted responses by having omitted variables –i.e., an underspecified model.

It has useful properties for selection of regressors:

- For a model that fits the data "adequately" $\Rightarrow E[C_p] \approx k$
- For the full model (no bias), with *k* parameters $\Rightarrow E[C_p] = k$.

• Other popular statistics: RIC (Risk Inflation Criteria), FPE, OOS R².

Modeling Strategies: Model Validation

• Cross validation, as in Lecture 5, can be used to select a model. For example, *K*-fold cross-validation. We have already done this in combination with Best subset.

Example: Suppose using best subsets to model IBM excess returns, using the k=3 Fama-French factors, we selected three model: CAPM (M1); Mkt_RF & SMB (M2); and the 3-factor F-F Model (M3).

Now, we use *K*-fold cross-validation, with *K*=5. CV₅ M1: 0.003542756 CV₅ M2: **0.003505873** CV₅ M3: 0.003556918

<u>Note</u>: Models look very similar. Practitioners compute a SE for CV_K and use a one SE rule. If within one SE, keep simplest model (M1).

Testing Model Specification: Non-Nested Models Example:

Model 1	$y = X\beta + W\delta + \varepsilon$
Model 2	$y = X\beta + Z\gamma + \xi$

• If the dependent variable is the same in both models (as is the case here), we can simply use Adjusted- R^2 to rank the models and select the one with the largest Adjusted- R^2 .

- We can also use AIC and/or BIC to rank the models.
- But, we can also use more sophisticated, testing-based, methods.

• Testing-based Method 1: Encompassing

(1) Form a composite or *encompassing* model that nests both rival models –Model 1 & Model 2. This is the **unrestricted Model**, ME.

(2) Test the relevant restrictions of each rival model against ME. We do two F-tests:

(i) Test ME (Unrestricted Model) against Model 1 (Restricted Model)(ii) Test ME (Unrestricted Model) against Model 2 (Restricted Model)

• If we reject the restrictions against one Model, say Model 1, and we cannot reject the restrictions against the other, Model 2, we are done: We select the Model that the F test do not reject restrictions (Model 2).

Assuming the restrictions cannot be rejected, we prefer the model with the lower F statistic for the test of restrictions.

Example: We have: Model 1 $y = X\beta + W\delta + \varepsilon$ Model 2 $y = X\beta + Z\gamma + \xi$ Then, the **Encompassing Model (ME)** is: ME: $y = X\beta + W\delta + Z\gamma + \varepsilon$ Now test, separately, the hypotheses (1) $\delta = 0$ and (2) $\gamma = 0$. That is, F-test for H₀: $\gamma = 0$: ME (U Model) vs Model 1 (R Model). F-test for H₀: $\delta = 0$: ME (U Model) vs Model 2 (R Model). If we reject H₀: $\gamma = 0 \Rightarrow$ We have evidence against Model 1 If we reject H₀: $\delta = 0 \Rightarrow$ We have evidence against Model 1. If we reject H₀: $\delta = 0 \Rightarrow$ We have evidence against Model 2. <u>Note</u>: We test a hybrid model, a combination of two models. Also, multicollinearity may appear.

Non-nested Models and Tests: Encompassing

Non-nested Models and Tests: IFE or PPP?

• Two of the main theories to explain the behaviour of exchange rates, S_t, are the **International Fisher Effect (IFE)** and the **Purchasing Power Parity (PPP)**. We use the direct notation for S_t, that is, units of Domestic Currency per 1 unit of Foreign currency.

• IFE states that, in equilibrium, changes in exchange rates (e) are driven by the interest rates differential between the domestic currency, i_d , and the foreign currency, i_f . A DGP consistent with IFE is:

$$\mathbf{e} = \alpha^1 + \beta^1 \left(\mathbf{i_d} - \mathbf{i_f} \right) + \mathbf{\epsilon}^2$$

• Relative PPP states that that, in equilibrium, **e** are driven by the inflation rates differential between the domestic Inflation rate, I_d , and the foreign Inflation rate, I_f A GDP consistent with IFE is:

$$\mathbf{e} = \boldsymbol{\alpha}^2 + \beta^2 \left(\mathbf{I}_{\mathbf{d}} - \mathbf{I}_{\mathbf{f}} \right) + \boldsymbol{\epsilon}$$

• Theories are non-nested, use non-nested methods to pick a model.

Non-nested Models and Tests: *J*-test • Testing-based Method 1: Davidson-MacKinnon (1981)'s *J*-test. We start with two non-nested models. Say, <u>Model 1</u>: $y = X\beta + \varepsilon$ <u>Model 2</u>: $y = Z\gamma + \xi$ <u>Idea</u>: If Model 2 is true, then the fitted values from the Model 1, when added to the 2nd equation, should be insignificant. • Steps: (1) Estimate Model 1 \Rightarrow obtain fitted values: **Xb**. (2) Add **Xb** to the list of regressors in Model 2 $\Rightarrow y = Z\gamma + \lambda Xb + \xi$ (3) Do a *t-test* on λ . A significant *t*-value would be evidence against Model 2 and in favour of Model 1.

Non-nested Models: *J*-test (4) Repeat the procedure for the models the other way round. (4.1) Estimate Model 2 ⇒ obtain fitted values: Zc. (4.2) Add Zc to the list of regressors in Model 1: ⇒ y = Xβ + λ Zc + ε (4.3) Do a *t*-test on λ. A significant *t*-value would be evidence against Model 1 and in favour of Model 2. (5) Rank the models on the basis of this test. It is possible that we cannot reject both models. This is possible in small samples, even if one model, say Model 2, is true. It is also possible that both *t*-tests reject H₀ (λ ≠ 0 & λ ≠ 0). This is not unusual. McAleer's (1995), in a survey, reports that out of 120 applications all models were rejected 43 times.

Non-nested Models: J-test

<u>Technical Note</u>: As some of the regressors in step (3) are stochastic, Davidson and MacKinnon (1981) show that the *t-test* is *asymptotically* valid.

• One would also want to examine the diagnostic test results when choosing between two models.

Non-nested Models: *J*-test – IFE or PPP? **Example**: Now, we test Model 1 vs Model 2, using the *J*-test. $\mathbf{e} = \mathbf{\alpha}^1 + \beta^1 \left(\mathbf{i}_{\mathbf{d}} - \mathbf{i}_{\mathbf{f}} \right) + \mathbf{\epsilon}^1$ Model 1 (IFE): Model 2 (PPP): $\mathbf{e} = \alpha^2 + \beta^2 (\mathbf{I}_d - \mathbf{I}_f) + \mathbf{\epsilon}^2$ y <- lr_usdgbp fit_m1 <- $lm(y \sim int_dif)$ summary(fit_m1) y_hat1 <- fitted(fit_m1)</pre> fit_J1 <- $lm(y \sim inf_dif + y_hat1)$ summary(fit_J1) fit_m2 <- $lm(y \sim inf_dif)$ summary(fit_m2) y_hat2 <- fitted(fit_m2)</pre> fit_J2 <- $lm(y \sim int_dif + y_hat2)$ summary(fit_J2)

Non-nested Models: J-test – Application • We want to test $\mathbf{H}_{\mathbf{0}}: \boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}_{0}$ (additive) vs $\mathbf{H}_{1}:\ln \mathbf{y}=(\ln \mathbf{X})\mathbf{\gamma}+\mathbf{\varepsilon}_{1}$ (multiplicative) • We look at the *J*-test Step 1: OLS on H₁: get $\hat{\gamma}$ OLS $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\lambda}_1 \exp\{\ln(\mathbf{X})\,\hat{\boldsymbol{\gamma}}\} + \boldsymbol{\varepsilon} \implies t\text{-test on }\boldsymbol{\lambda}_1$ Step 2: OLS on **H**₀: get **b** OLS $\ln y = (\ln X) \gamma + \lambda_0 Xb + \varepsilon \implies t$ -test on λ_0 • Situations: (1) Both OK: $\lambda_1 = 0$ and $\lambda_0 = 0 \implies$ get more data (2) Only 1 is OK: $\lambda_1 \neq 0$ and $\lambda_0 = 0$ (multiplicative is OK); $\lambda_0 \neq 0$ and $\lambda_1 = 0$ (additive is OK) (3) Both rejected: $\lambda_1 \neq 0$ and $\lambda_0 \neq 0 \implies$ new model is needed.

Non-nested Models: J-test - Considerations

• The *J*-test was designed to test non-nested models (one model is the true model, the other is the false model), not for choosing competing models –the usual use of the test.

• The *J*-test is likely to over reject the true (model) hypothesis when one or more of the following features is present:

i) A poor fit of the true model

ii) A low/moderate correlation between the regressors of the 2 models

iii) The false model includes more regressors than the correct model.

Davidson and MacKinnon (2004) state that the *J*-test will over-reject, *often quite severely* in finite samples when the sample size is small or where conditions (i) or (iii) above are obtained.

Modeling Strategies: Significance level, a

• So far, we have assumed that the distribution of the test statistic –say the *F*-statistic-- under H_0 is known exactly, so that we have what is called an *exact test*.

• Technically, the *size of a test* is the supremum of the rejection probability over all DGPs that satisfy H_0 . For an exact test, the size equals the *nominal level*, α –i.e., the Prob[Type I error] = α .

• Usually, the distribution of a test is known only approximately *(asymptotically)*. In this case, we need to draw a distinction between the nominal level *(nominal size)* of the test and the actual *rejection probability (empirical size)*, which may differ greatly from the nominal level.

• Simulations are needed to gauge the empirical size of tests.

Modeling Strategies: A word about a

• Ronald Fisher, before computers, tabulated distributions. He used a .10, .05, and .01 percentiles. These tables were easy to use and, thus, those percentile became the de-facto standard α for testing H₀.

• "It is usual and convenient for experimenters to take 5% as a standard level of significance." –Fisher (1934).

• Given that computers are powerful and common, why is p = 0.051 unacceptable, but p = 0.049 is great? There is no published work that provides a theoretical basis for the standard thresholds.

• Rosnow and Rosenthal (1989): " ... surely God loves .06 nearly as much as .05."

Modeling Strategies: A word about a

<u>Practical advise</u>: In the usual Fisher's null hypothesis (significance) testing, significance levels, α , are arbitrary. Make sure you pick one, say 5%, and stick to it throughout your analysis or paper.

• Report *p-values*, along with CI's. Search for *economic significance*.

• Q: .10, .05, or .01 significance?

Many tables will show *, **, and *** to show .10, .05, and .01 significance levels. Throughout the paper, the authors will point out the different significance levels. In these papers, it is not clear what α is the paper using for inference.

• In a Neyman-Pearson world, we can think of these stars (or *p*-values) as ways of giving weights to H_0 relative to H_1 .

Modeling Strategies: A word about H₀

• In applied work, we only learn when we reject H_0 . Failing to reject H_0 provides almost no information about the state of the world.

• Thus, failing to reject H_0 does not rule out an infinite number of other competing research hypotheses.

• Null hypothesis significance testing is asymmetric: if the test statistic is "too large" for a given H_0 then H_0 is rejected; but if the test statistic is not "too large" then H_0 is not automatically accepted.

• It is dangerous to "accept" the conclusion from a non-rejected H_0 . But, it is common. Eight of the twenty (40%) articles in the *American Political Science Review* Volume 91 (1997), that used a H_0 , drew substantive conclusions from a fail to reject decision.

Modeling Strategies: A word about H₀

• In applied work, we only learn when we reject H_0 ; say, when the *p*-value< α . But, rejections are of two types:

- Correct ones, driven by the power of the test,

- Incorrect ones, driven by Type I Error ("statistical accident," luck).

• It is important to realize that, however small the *p*-value, there is always a finite chance that the result is a pure accident. At the 5% level, there is 1 in 20 chances that the rejection of H_0 is just luck.

• Since negative results are difficult to publish (*publication bias*), there is an unknown but possibly large number of false claims taken as truths.

Example (from Lecture 4): If $\alpha = 0.05$, proportion of false H₀=10%, and $\pi = .50$, 47.4% of rejections are true H₀ -i.e., "false positives."

Model Selection Methods: Summary

• Eight literature strands can be delineated:

(1) *Specific-to-general*: Anderson (1962), Hendry and Mizon (1978), and Hendry (1979), for critiques;

(2) Retaining the general model: Yancey and Judge (1976), and Judge and Bock (1978);

(3) *Testing Theory-based models*: Hall (1978), criticized by Davidson and Hendry (1981), and Hendry and Mizon (2000); Stigum (1990) proposes a formal approach;

(4) Other 'rules' for model selection, such as:

- step-wise regression: Learner (1983a), for a critical appraisal

- 'optimal' regression: algorithm to maximize the Adj-R² with a specified set of regressors. See Coen, Gomme and Kendall (1969);

Model Selection Methods: Summary

• Eight literature strands can be delineated (continuation): (5) *Model comparisons*, often based on non-nested hypothesis tests or encompassing: Cox (1961, 1962), Pesaran (1974), and the survey in Hendry and Richard (1989);

(6) *Model selection by information criteria*: Schwarz (1978), Hannan and Quinn (1979), Amemiya (1980);

(7) *Bayesian model comparisons*: Learner (1978) and Clayton, Geisser and Jennings (1986);

(8) General-to-specifics (gets): Anderson (1962), Sargan (1973, 1981), Hendry (1979), and White (1990).

Criteria for Model Selection: Judgement Call

• In the end, judgment must be used in weighing up various criteria:

- The Economic Criterion –are the estimated parameters plausible? (Economic Significance)

- The First Order Statistical Criterion –does the model provide a good fit (in-sample) with statistically significant parameter estimates?

- The Second Order Statistical Criterion —is the model generally free of misspecification problems — as evidenced in the diagnostic tests?

- The Out of Sample Predictive Criterion –does the model provide good out of sample predictions? Model validation, with the different flavours, can be used here.

Model Selection: Causality and Identification

• In empirical work, we are interested in identifying causal relations, say from **X** to **y**, as implied in the DGP of the CLM: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$.

• Suppose we have two correlated variables: $Y_t \& X_t$. The co-variation in $Y_t \& X_t$ can be driven by (not mutually exclusive):

- Causation from Y_t to X_t : Changes in $Y_t \implies$ changes in X_t
- Causation from X_t to Y_t : Changes in $X_t \implies$ changes in Y_t
- Correlated through a 3rd variable, W_t : changes in $W_t \Rightarrow$ changes $X_t \& Y_t$

• In practice, it is not easy to say what generates variation in $Y_t \& X_t$. The third case, especially when W_t is an unobservable variable, creates a lot complications.

Example: Y_t : earnings; X_t : schooling; W_t : ability.

Model Selection: Causality and Identification

• There are four approaches for identification (of variation):

- Experiments. The researcher generates the variation in the variables.

- *Natural Experiments*. A known exogenous event generates the variation in the variables.

- Instrumental variables. A variable provides variation.

- *Econometric Identification*. We use econometric assumptions for identification.

• In time series, there is the concept of *Granger causality*, where past changes in one variable affect the present values of another variable. This is not, strictly speaking, the causation we discuss here.

Model Selection: Causality and Identification

• To be precise, the identification problem in econometrics refers to the problem of identifying and estimating one or more coefficients of a system of simultaneous equations.

Model Selection: Causality and Identification

• Experiments

Experiments are popular in the sciences (say, biology, physics). For example, we want to test a new treatment. Then,

A sample is divided randomly in two similar groups: *treated* group
 control group. (A *randomized* study: Only difference is the treatment!)
 Look for differences in both groups.

⇒ Rare in economics and finance; they can be very expensive or unethical (say, exposing people to a "poverty shock"). Some work in small communities and small units in some businesses.

<u>Problem</u>: Not easy to randomize these man-made experiments that involve humans.

Model Selection: Causality and Identification

• Natural experiments

An exogenous (historical) event (not necessarily a nature event) provides a situation where groups can be reasonably randomized in a *treated* (affected by the natural event) and a *control group* (not affected by the natural event).

In the absence of experiments, natural experiments give us a very good way to identify causation.

Examples: Changes in tax code and regulations; changes in accounting standards, shocks (Covid-19, stock market crisis), disasters (earthquakes, floods, etc.), laws or rules that impose thresholds (*discontinuity*) for behaviors, etc. (More on Lectures 8 & 15.)

Problem: Not easy to generalize, not clear how robust results are.

Model Selection: Causality and Identification

• Instrumental Variables

Suppose we want to study the effect of networking on CEO compensation. Since CEO compensation and networking may be affected by the unobserved natural ability of an individual (W_t) , a simple regression will be biased (omitted variables problem).

Suppose we have a variable, Z, correlated with networking, but not with natural ability (*ethnicity*?, *age*?, *number of childhood friends*?) –i.e., Z induces variation in X *unrelated* to W_r . Then, we use Z to study the effect of networking on CEO compensation.

We call Z an *instrument*. Usually, we can relate Z to a *natural experiment*.

Problem: As we will see later, in Lecture 8, finding Z is not easy.

Model Selection: Causality and Identification

• Econometric Identification

We think that networking is correlated with ability, then we model it. Actually, we model everything. Very transparent in the assumptions.

We end up with a Simultaneous Equations Models (SEM), which we will study later in Lecture 16.

Problem: They tend to be (very) complicated.