

Nonlinear Time Series in Financial Forecasting*

Gloria González-Rivera
Department of Economics
University of California, Riverside
Riverside, CA 92521-0427
E-mail: gloria.gonzalez@ucr.edu
phone: +1 951-827-1470
fax +1 951-827-5685

Tae-Hwy Lee
Department of Economics
University of California, Riverside
Riverside, CA 92521-0427
E-mail: taelee@ucr.edu
phone: +1 951-827-1509
fax +1 951-827-5685

January 18, 2008

*[Entry ID: 00221] “Financial Forecasting, Nonlinear Time Series in”, for *Encyclopedia of Complexity and System Science*. Publisher contact: Julia.Koerting@springer.com, Kerstin.Kindler@springer.com. Editor: Bruce Mizrach <mizrach@economics.rutgers.edu>. Project website: <http://refworks.springer.com/complexity>

Article Outline

Glossary

1. Definitions
2. Introduction
3. Nonlinear models for the conditional mean
4. Nonlinear forecasting models for the conditional variance of returns
5. Forecasting quantiles, directions, durations, density, intervals, and etc.
6. Conclusions
7. Future Directions
8. Bibliography

Glossary

ARCH: autoregressive conditional heteroskedasticity

Artificial neural network: is a nonlinear flexible functional form, connecting inputs to outputs, being capable of approximating a measurable function to any desired level of accuracy provided that sufficient complexity (in terms of number of hidden units) is permitted.

Bagging: short for *bootstrap aggregating*. Bagging is a method of smoothing the instabilities by averaging the predictors over bootstrap predictors and thus lowering the sensitivity of the predictors to training samples. A predictor is said to be unstable if perturbing the training sample can cause significant changes in the predictor.

Functional coefficient model: a model with time-varying and state-dependent coefficients. The number of states can be infinite.

Linearity in mean conditional mean: the process $\{y_t\}$ is linear in mean conditional on X_t if

$$\Pr [\mathbb{E}(y_t|X_t) = X_t'\theta^*] = 1 \text{ for some } \theta^* \in \mathbb{R}^k.$$

Markov-switching model: features parameters changing in regime, but in contrast with the threshold models the change is dictated by a non-observable state variable that is modelled as a Markov chain. a.k.a., hidden Markov chain model.

Sieves: the sieves or approximating spaces are approximations to an unknown function, that are dense in the original function space. Sieves can be constructed using linear spans of power series, e.g., Fourier series, splines, or many other basis functions such as artificial neural network (ANN), and various polynomials (Hermite, Laguerre, etc.).

Smooth transition models: threshold model with the indicator function replaced by a smooth monotonically increasing differentiable function such as a probability distribution function.

Threshold model: a nonlinear model with time-varying coefficients specified using an indicator which takes a non-zero value when a state variable falls on a specified partition of a set of states, and zero otherwise. The number of partitions is finite.

Volatility: Volatility in financial economics is often measured by the conditional variance (e.g., ARCH) or the conditional range. It is import for any decision making under uncertainty such as portfolio allocation, option pricing, risk management.

1 Definitions

1.1 Financial forecasting

Financial forecasting is concerned with the prediction of prices of financial assets such as stocks, bonds, options, interest rates, exchange rates, etc. The importance of financial forecasting derives primarily from the role of financial markets within the macro economy. The development of financial instruments and financial institutions contribute to the growth and stability of the overall economy. Because of this interconnection between financial markets and the real economy, financial forecasting is also intimately linked to macroeconomic forecasting, which is concerned with the prediction of macroeconomic aggregates such as growth of the gross domestic product, consumption growth, inflation rates, commodities prices, etc. Financial forecasting and macroeconomic forecasting share many of the techniques and statistical models that will be explained in detail in this article.

In financial forecasting a major object of study is the return to a financial asset, mostly calculated as the continuously compounded return, i.e., $y_t = \log p_t - \log p_{t-1}$ where p_t is the price of the asset at time t . Nowadays financial forecasters use sophisticated techniques that combine the advances in modern finance theory, pioneered by Markowitz (1959), with the advances in time series econometrics, in particular the development of nonlinear models for conditional moments and quantiles of asset returns.

The aim of finance theory is to provide models for expected returns taking into account the uncertainty of the future asset payoffs. In general, financial models are concerned with investors' decisions under uncertainty. For instance the portfolio allocation problem deals with the allocation of wealth among different assets that carry different levels of risk. The implementation of these theories relies on econometric techniques that aim to estimate financial models and testing them against the data. Financial econometrics is the branch of econometrics that provides model-based statistical inference for financial variables, and therefore financial forecasting will provide their corresponding model-based predictions. However there are also econometric developments that inform the construction of *ad hoc* time series models that are valuable on describing the stylized facts of financial data.

Since returns $\{y_t\}$ are random variables, the aim of financial forecasting will be to forecast

any conditional moments, quantiles, and eventually the conditional distribution of these variables. Most of the time our interest will be centered on expected returns and expected volatility as these two moments are crucial components on portfolio allocation problems, option valuation, and risk management, but it is also possible to forecast quantiles of a random variable, and therefore to forecast the expected probability density function. Density forecast is the most complete forecast as it embeds all the information on the financial variable of interest. Financial forecasting is also concerned with other financial variables like durations between trades and directions of price changes. In these cases, it is also possible to construct conditional duration models and conditional probit models that can be used to forecasting durations and timing the markets.

The earliest characterization of financial prices has its roots in the games of chance that are also associated with the beginnings of probability theory in the XVI century. Borrowing from the concept of fair game, financial prices are said to enjoy the *martingale property* if tomorrow's price is expected to be equal to today's price given some information set; in other words tomorrow's price has an equal chance to either move up or move down, and thus the best forecast must be the current price. The martingale property is written as

$$\mathbb{E}(p_{t+1}|\mathcal{F}_t) = p_t$$

where \mathbb{E} is the expectation operator and the information set $\mathcal{F}_t \equiv \{p_t, p_{t-1}, p_{t-2}, \dots\}$ is the collection of past and current prices. From a forecasting point of view, the martingale model implies that changes in financial prices ($p_{t+1} - p_t$) are not predictable.

A most restrictive form of the martingale property, proposed by Bachelier (1900) in his theory of speculation is the model (in logarithms)

$$\log p_{t+1} = \mu_t + \log p_t + \varepsilon_{t+1},$$

where $\mu_t = \mu$ is a constant drift and ε_{t+1} is an identically and independently distributed (i.i.d.) error that is assumed to be normally distributed with zero mean and constant variance σ^2 . Since the return is the percentage change in prices, i.e. $y_t = \log p_t - \log p_{t-1}$, an equivalent model for asset returns is

$$y_{t+1} = \mu_t + \varepsilon_{t+1}.$$

Taking a conditional expectation, we have that $\mathbb{E}(y_{t+1}|\mathcal{F}_t) = \mu_t$. If the conditional mean return is not time-varying and fixed a constant ($\mu_t = \mu$), then the returns are not forecastable based on past price information. In addition and given the assumptions on the error term, returns are independent and identically distributed random variables. These two properties that the drift is a constant and the error term is i.i.d. is too restrictive and it rules out the possibility of any predictability in asset returns. A less restrictive and more practically plausible version is obtained when these restrictions are relaxed. The error term may be heteroscedastic so that returns have different (unconditional or conditional) variances and consequently they are not identically distributed, and/or the error term, though uncorrelated, may exhibit dependence in higher moments and in this case the returns are not independent random variables.

Modern finance theory quantifies the trade-off between expected returns and risk. Investors are willing to hold risky assets when they are adequately compensated by the amount of risk they bear. Arguably, the two most important asset pricing models in modern finance theory are the Capital Asset Pricing Model (CAPM) proposed by Sharpe (1964) and Lintner (1965) and the Arbitrage Pricing Theory (APT) proposed by Ross (1976). Both models claim that the expected return to an asset is a linear function of risk; in CAPM risk is related to the covariance of the asset return with the return to the market portfolio, and in APT risk is measured as exposure to a set of factors, which may include the market portfolio among others. The CAPM model in its original version is written as

$$\mathbb{E}(y_i) = y_f + \beta_{im} [\mathbb{E}(y_m) - y_f],$$

where y_f is the risk-free rate, y_m is the return to the market portfolio, and β_{im} is the risk of asset i defined as

$$\beta_{im} = \frac{\text{cov}(y_i, y_m)}{\text{var}(y_m)} = \frac{\sigma_{im}}{\sigma_m^2}.$$

This model has a time series version known as the conditional CAPM (Bollerslev, Engle, and Wooldridge, 1988) that it may be useful for forecasting purposes. For asset i and given an information set as $\mathcal{F}_t = \{y_{i,t}, y_{i,t-1}, \dots; y_{m,t}, y_{m,t-1}, \dots\}$, the expected return is a linear function of a time-varying beta

$$\mathbb{E}(y_{i,t+1}|\mathcal{F}_t) = y_f + \beta_{im,t} [\mathbb{E}(y_{m,t+1}|\mathcal{F}_t) - y_f]$$

where $\beta_{im,t} = \frac{\text{cov}(y_{i,t+1}, y_{m,t+1} | \mathcal{F}_t)}{\text{var}(y_{m,t+1} | \mathcal{F}_t)} = \frac{\sigma_{im,t}}{\sigma_{m,t}^2}$. From this type of models is evident that we need to model the conditional second moments of returns jointly with the conditional mean. A general finding of this type of models is that when there is high volatility expected returns are high, and hence forecasting volatility becomes important for the forecasting of expected returns. In the same spirit, the APT models have also conditional versions that exploit the information contained in past returns. A K -factor APT model is written as

$$y_t = c + B' f_t + \varepsilon_t,$$

where f_t is a $K \times 1$ vector of factors and B is a $K \times 1$ vector of sensitivities to the factors. If the factors have time-varying second moments, it is possible to specify an APT model with a factor structure in the time-varying covariance matrix of asset returns (Engle, Ng, and Rothschild, 1990), which in turn can be exploited for forecasting purposes.

The conditional CAPM and conditional APT models are fine examples on how finance theory provides a base to specify time-series models for financial returns. However there are other time series specifications, more *ad hoc* in nature, that claim that financial prices are nonlinear functions of the information set and by that they impose some departures from the martingale property. In this case it is possible to observe some predictability in asset prices. This is the subject of nonlinear financial forecasting. We begin with defining the concepts of linearity and nonlinearity.

1.2 Linearity and nonlinearity

We begin with defining the concepts of linearity and nonlinearity. For nonlinearity, it is relevant to ask first how this concept should be defined. It is important to be precise about the meaning of the word ‘linearity’. Lee, White, and Granger (LWG, 1993) is the first who clarify it. Let $\{Z_t\}$ be a stochastic process, and partition Z_t as $Z_t = (y_t \ X_t')'$, where (for simplicity) y_t is a scalar and X_t is a $k \times 1$ vector. X_t may (but need not necessarily) contain a constant and lagged values of y_t . LWG define that the process $\{y_t\}$ is *linear in mean conditional on X_t* if

$$\Pr [\mathbb{E}(y_t | X_t) = X_t' \theta^*] = 1 \text{ for some } \theta^* \in \mathbb{R}^k.$$

In the context of forecasting, Granger and Lee (1999) define linearity as follows. Define $\mu_{t+h} = \mathbb{E}(y_{t+h}|\mathcal{F}_t)$ being the optimum least squares h -step forecast of y_{t+h} made at time t . μ_{t+h} will generally be a nonlinear function of the contents of \mathcal{F}_t . Denote m_{t+h} to be the optimum *linear* forecast of y_{t+h} made at time t , being the best forecast that is constrained to be a linear combination of the contents of $X_t \in \mathcal{F}_t$. Granger and Lee (1999) define that $\{y_t\}$ is said to be *linear in conditional mean* if μ_{t+h} is linear in X_t , i.e., $\Pr[\mu_{t+h} = m_{t+h}] = 1$ for all t and for all h . Typically, for simplicity, the interest is only with $h = 1$ in this definition, as considered in Granger and Lee (1999) and many others. Under this definition the focus is the conditional mean and thus a process exhibiting autoregressive conditional heteroskedasticity (ARCH) of Engle (1982) may nevertheless exhibit linearity of this sort because ARCH does not refer to the conditional mean. This is appropriate whenever we are concerned with the adequacy of linear models for forecasting the conditional mean returns. See White (2006, Section 2) for more rigorous treatments on the definitions of linearity and nonlinearity.

This definition may be extended with some caution to the concepts for linearity in a higher moment or quantiles, but the definition may depend on the focus or interest of the researcher. Let $\varepsilon_{t+h} = y_{t+h} - \mu_{t+h}$ and $\sigma_{t+h}^2 = \mathbb{E}(\varepsilon_{t+h}^2|\mathcal{F}_t)$. If we may like to consider the ARCH and GARCH as linear models, we may say $\{\sigma_{t+h}^2\}$ is linear in conditional variance if σ_{t+h}^2 is a linear function of lagged ε_{t-j}^2 and σ_{t-j}^2 for some h or for all h . Alternatively, $\sigma_{t+h}^2 = \mathbb{E}(\varepsilon_{t+h}^2|\mathcal{F}_t)$ is said to be linear in conditional variance if σ_{t+h}^2 is a linear function of $x_t \in \mathcal{F}_t$ for some h or for all h . Similarly, we can have consider linearity in conditional quantiles. It seems that the issue of linearity vs nonlinearity is most relevant for the conditional mean, and it may be a more relevant matter whether a certain specification is correct or incorrect (rather than linear or nonlinear) for higher order conditional moments or quantiles.

2 Introduction

There exists a nontrivial gap between martingale difference and serial uncorrelatedness. The former implies the latter, but not vice versa. Consider a stationary time series $\{y_t\}$. Often, serial dependence of $\{y_t\}$ is described by its autocorrelation function $\rho(j)$, or by its

standardized spectral density

$$h(\omega) = \frac{1}{2\pi} \sum_{j=-\infty}^{\infty} \rho(j)e^{-ij\omega}, \omega \in [-\pi, \pi].$$

Both $h(\omega)$ and $\rho(j)$ are the Fourier transform of each other, containing the same information of serial correlations of $\{y_t\}$. A problem of using $h(\omega)$ and $\rho(j)$ is that they can not capture nonlinear time series that have zero autocorrelation but are not serially independence. Nonlinear MA and Bilinear series are examples:

$$\begin{aligned} \text{Nonlinear MA} & : Y_t = be_{t-1}e_{t-2} + e_t, \\ \text{Bilinear} & : Y_t = be_{t-1}Y_{t-2} + e_t. \end{aligned}$$

These processes are serially uncorrelated, but they are predictable using the past information. Hong and Lee (2003a) note that autocorrelation, variance ratio, and power spectrum can easily miss these processes. Misleading conclusions in favor of the martingale hypothesis could be reached when these test statistics are insignificant. It is therefore important and interesting to explore whether there exists a gap between serial uncorrelatedness and martingale difference for financial forecasting, and if so, whether the neglected nonlinearity in conditional mean can be explored to forecast financial asset returns.

In the forthcoming sections, we will present, without being exhaustive, nonlinear time series models for financial returns, which are the basis for nonlinear forecasting. In Section 3, we review nonlinear models for the conditional mean of returns. A general representation is $y_{t+1} = \mu(y_t, y_{t-1}, \dots) + \varepsilon_{t+1}$ with $\mu(\cdot)$ a nonlinear function of the information set. If $\mathbb{E}(y_{t+1}|y_t, y_{t-1}, \dots) = \mu(y_t, y_{t-1}, \dots)$, then there is a departure from the martingale hypothesis, and past price information will be relevant to predict tomorrow's return. In Section 4, we review models for the conditional variance of returns. For instance, a model like $y_{t+1} = \mu + u_{t+1}\sigma_{t+1}$ with time-varying conditional variance $\sigma_{t+1}^2 = \mathbb{E}((y_{t+1} - \mu)^2|\mathcal{F}_t)$ and i.i.d. error u_{t+1} , is still a martingale-difference for returns but it represents a departure from the independence assumption. The conditional mean return may not be predictable but the conditional variance of the return will be. In addition, as we have seen modeling time-varying variances and covariances will be very useful for the implementation of conditional CAPM and APT models.

3 Nonlinear models for the conditional mean

To forecast the changes in prices of financial assets such as stocks, bonds, options, interest rates, exchange rates, etc., we consider models for $\mu_{t+h} = \mathbb{E}(y_{t+h}|\mathcal{F}_t)$. In considering the conditional mean in this section, we are restricting the loss function of the forecast error to be the mean squared forecast error (MSFE). Other loss functions will lead to forecasting different aspect of the forecast density. For example, the loss function of the mean absolute error is associated with the conditional median. Some evidence suggests that μ_{t+h} is time-varying but is of a complicated form. In particular, it cannot be modeled simply by autoregressive polynomials in lagged changes of financial asset prices. Various parametric and nonparametric models can be used. Examples of parametric models are autoregressive bilinear and threshold models. Examples of nonparametric models are artificial neural network, kernel and nearest neighbor regression models.

In this section we will consider a brief discussion on a small set of nonlinear specifications. However, it is impossible to cover an exhaustive set of all nonlinear models. In the mean time, some of the nonlinear models are universal approximators as discussed in White (2006) and Chen (2006). Thus, based on many theoretical results that provide good and computable approximations to an unknown function we can also consider various *sieves*. For example, the sieves or approximating spaces can be constructed using linear spans of power series, Fourier series, splines, or many other basis functions such as artificial neural network (ANN), Hermite polynomials as used in e.g., Gallant and Nychka (1987) for modelling seminonparametric density, and Laguerre polynomials used in Nelson and Siegel (1987) in modelling the yield curve. Diebold and Li (2006) and Huang, Lee, and Li (2007) use the Nelson-Siegel model in forecasting yields and inflation. We consider various parametric nonlinear models (threshold model, smooth transition model, Markov switching model, nonlinear models based on random fields), nonparametric models (e.g., local linear, local polynomial, local exponential, functional coefficient models), and nonlinear models based on sieves (e.g., ANN, various polynomials). Some other particular parametric nonlinear models that are not included below may be found from other books on nonlinear time series models such as Fan and Yao (2003), Gao (2007), Tsay (2005). We begin with a very simple nonlinear model.

3.1 A simple nonlinear model with dummy variables

Goyal and Welch (2006) forecast the equity premium namely, the S&P 500 index return minus T-bill rate using many predictors such as stock-related variables (e.g., dividend-yield, earning-price ratio, book-to-market ratio, corporate issuing activity, etc.), interest-rate-related variables (e.g., treasury bills, long-term yield, corporate bond returns, inflation, investment to capital ratio), and ex ante consumption, wealth, income ratio (modified from Lettau and Ludvigson 2001). They test the out-of-sample performance conditional on observed IS significance. They find the predictors have better performance in bad times, such as Great Depression (1930-33), oil-shock period (1973-75), and bubble-crash period (1999-2001). Also, they argue that reasonable truncation on equity premium because no investor is interested in negative premium.

Campbell and Thompson (2007), inspired by out-of-sample forecasting of Goyal and Welch (2006), argue that if we impose some restrictions on the signs of predictors' coefficients and excess return forecasts, predictive power of some predictors can beat that of historical average equity premium. Similar to Goyal and Welch (2006), they also use rich forecasting variables – valuation ratios (e.g., dividend price ratio, earning price ratio, and book to market ratio), real return on equity, nominal interest rates and inflation, and equity share of new issues and consumption-wealth ratio. They apply two restrictions –the first one is to restrict predictors' coefficients to have the theoretically expected sign and to set wrong-signed coefficients to zero, and the second one is to rule out negative equity premium forecast. They show the effectiveness of these theoretically-inspired restrictions almost always improve the out-of sample performance of predictive regressions. This is an example where shrinkage works by reducing the forecast error variance at the cost of higher forecast bias but with a smaller mean squared forecast error (the sum of variance and the squared bias).

The results from Goyal and Welch (2006) and Campbell and Thompson (2007) support simple form of nonlinearity that can be generalized to threshold models or time-varying coefficient models, which we consider next.

3.2 Threshold models

Many financial and macroeconomic time series exhibit different characteristics over time depending upon the state of the economy. For instance, we observe bull and bear stock markets, high volatility versus low volatility periods, recessions versus expansions, credit crunch versus excess liquidity, etc. If these different regimes are present in economic time series data, econometric specifications should go beyond linear models as these assume that there is only a single structure or regime over time. Nonlinear time series specifications that allow for the possibility of different regimes, also known as state-dependent models, include several types of models: threshold, smooth transition, and regime-switching models.

Threshold autoregressive (TAR) models (Tong, 1983, 1990) assume that the dynamics of the process is explained by an autoregression in each of the n regimes dictated by a conditioning or threshold variable. For a process $\{y_t\}$, a general specification of a TAR model is

$$y_t = \sum_{j=1}^n \left[\phi_o^{(j)} + \sum_{i=1}^{p_j} \phi_i^{(j)} y_{t-i} + \varepsilon_t^{(j)} \right] \mathbf{1}(r_{j-1} < x_t \leq r_j).$$

There are n regimes, in each one there is an autoregressive process of order p_j with different autoregressive parameters $\phi_i^{(j)}$, the threshold variable is x_t with r_j thresholds and $r_o = -\infty$ and $r_n = +\infty$, and the error term is assumed i.i.d. with zero mean and different variance across regimes $\varepsilon_t^{(j)} \sim \text{i.i.d. } (0, \sigma_j^2)$, or more generally $\varepsilon_t^{(j)}$ is assumed to be a martingale difference. When the threshold variable is the lagged dependent variable itself y_{t-d} , the model is known as self-exciting threshold autoregressive (SETAR) model. The SETAR model has been applied to the modelling of exchange rates, industrial production indexes, and gross national product (GNP) growth, among other economic data sets. The most popular specifications within economic time series tend to find two, at most three regimes. For instance, Boero and Marrocu (2004) compare a two and three-regime SETAR models with a linear AR with GARCH disturbances for the euro exchange rates. On the overall forecasting sample, the linear model performs better than the SETAR models but there is some improvement in the predictive performance of the SETAR model when conditioning on the regime.

3.3 Smooth transition models

In the SETAR specification, the number of regimes is discrete and finite. It is also possible to model a *continuum* of regimes as in the Smooth Transition Autoregressive (STAR) models, (Teräsvirta, 1994). A typical specification is

$$y_t = \phi_0 + \sum_{i=1}^p \phi_i y_{t-i} + (\theta_0 + \sum_{i=1}^p \theta_i y_{t-i}) F(y_{t-d}) + \varepsilon_t$$

where $F(y_{t-d})$ is the transition function that is continuous and in most cases is either a logistic function or an exponential,

$$\begin{aligned} F(y_{t-d}) &= [1 + \exp(-\gamma(y_{t-d} - r))]^{-1} \\ F(y_{t-d}) &= 1 - [\exp(-\gamma(y_{t-d} - r))]^2 \end{aligned}$$

This model can be understood as many autoregressive regimes dictated by the values of the function $F(y_{t-d})$, or alternatively as an autoregression where the autoregressive parameters change smoothly over time. When $F(y_{t-d})$ is logistic and $\gamma \rightarrow \infty$, the STAR model collapses to a threshold model SETAR with two regimes. One important characteristic of these models, SETAR and STAR, is that the process can be stationary within some regimes and non-stationary within others moving between explosive and contractionary stages.

Since the estimation of these models can be demanding, the first question to solve is whether the nonlinearity is granted by the data. A test for linearity is imperative before engaging in the estimation of nonlinear specifications. An LM test that has power against the two alternatives specifications SETAR and STAR is proposed by Luukkonen et al (1988) and it consists of running two regressions: under the null hypothesis of linearity, a linear autoregression of order p is estimated in order to calculate the sum of squared residuals, SSE_0 ; the second is an auxiliary regression

$$y_t = \beta_0 + \sum_{i=1}^p \beta_i y_{t-i} + \sum_{i=1}^p \sum_{j=1}^p \psi_{ij} y_{t-i} y_{t-j} + \sum_{i=1}^p \sum_{j=1}^p \zeta_{ij} y_{t-i} y_{t-j}^2 + \sum_{i=1}^p \sum_{j=1}^p \xi_{ij} y_{t-i} y_{t-j}^3 + u_t$$

from which we calculate the sum of squared residuals, SSE_1 . The test is constructed as $\chi^2 = T(SSE_0 - SSE_1)/SSE_0$ that under the null hypothesis of linearity is chi-squared

distributed with $p(p+1)/2+2p^2$ degrees of freedom. There are other tests in the literature, for instance Hansen (1996) proposes a likelihood ratio test that has a non-standard distribution, which is approximated by implementing a bootstrap procedure. Tsay (1998) proposes a test based on arranged regressions with respect to the increasing order of the threshold variable and by doing this the testing problem is transformed into a change-point problem.

If linearity is rejected, we proceed with the estimation of the nonlinear specification. In the case of the SETAR model, if we fix the values of the delay parameter d and the thresholds r_j , the model reduces to n linear regressions for which least squares estimation is straightforward. Tsay (1998) proposes a conditional least squares (CLS) estimator. For simplicity of exposition suppose that there are two regimes in the data and the model to estimate is

$$y_t = \left[\phi_o^{(1)} + \sum_{i=1}^{p_1} \phi_i^{(1)} y_{t-i} \right] \mathbf{1}(y_{t-d} \leq r) + \left[\phi_o^{(2)} + \sum_{i=1}^{p_2} \phi_i^{(2)} y_{t-i} \right] \mathbf{1}(y_{t-d} > r) + \varepsilon_t$$

Since r and d are fixed, we can apply least squares estimation to the model and to obtain the LS estimates for the parameters ϕ_i 's. With the LS residual $\hat{\varepsilon}_t$, we obtain the total sum of squares $S(r, d) = \sum_t \hat{\varepsilon}_t^2$. The CLS estimates of r and d are obtained from $(\hat{r}, \hat{d}) = \arg \min S(r, d)$.

For the STAR model, it is also necessary to specify *a priori* the functional form of $F(y_{t-d})$. Teräsvirta (1994) proposes a modeling cycle consisting of three stages: specification, estimation, and evaluation. In general, the specification stage consists of sequence of null hypothesis to be tested within a linearized version of the STAR model. Parameter estimation is carried out by nonlinear least squares or maximum likelihood. The evaluation stage mainly consists of testing for no error autocorrelation, no remaining nonlinearity, and parameter constancy, among other tests.

Teräsvirta and Anderson (1992) find strong nonlinearity in the industrial production indexes of most of the OECD countries. The preferred model is the logistic STAR with two regimes, recessions and expansions. The dynamics in each regime are country dependent. For instance, in USA they find that the economy tends to move from recessions into expansions very aggressively but it will take a large negative shock to move rapidly from an expansion into a recession.

For forecasting with STAR models, see Lundbergh and Teräsvirta (2002). It is easy to construct the one-step-ahead forecast but the multi-step-ahead forecast is a complex problem. For instance, for the 2-regime threshold model, the one-step-ahead forecast is constructed as the conditional mean of the process given some information set

$$\mathbb{E}(y_{t+1}|\mathcal{F}_t; \theta) = \left[\phi_o^{(1)} + \sum_{i=1}^{p_1} \phi_i^{(1)} y_{t+1-i} \right] \mathbf{1}(y_{t+1-d} \leq r) + \left[\phi_o^{(2)} + \sum_{i=1}^{p_2} \phi_i^{(2)} y_{t+1-i} \right] \mathbf{1}(y_{t+1-d} > r)$$

provided that $y_{t+1-i}, y_{t+1-d} \in \mathcal{F}_t$. However, a multi-step-ahead forecast will be a function of variables that being dated at a future date do not belong to the information set; in this case the solution requires the use of numerical integration techniques or simulation/bootstrap procedures. See Granger and Teräsvirta (1993, Chapter 9) and Teräsvirta (2006) for more details on numerical methods for multi-step forecasts.

3.4 Markov-switching models

A Markov-switching (MS) model (Hamilton 1989, 1996) also features changes in regime, but in contrast with the SETAR models the change is dictated by a non-observable state variable that is modelled as a Markov chain. For instance, a first order autoregressive Markov switching model is specified as

$$y_t = c_{s_t} + \phi_{s_t} y_{t-1} + \varepsilon_t$$

where $s_t = 1, 2, \dots, N$ is the unobserved state variable that is modelled as an N -state Markov chain with transition probabilities $p_{ij} = P(s_t = j | s_{t-1} = i)$, and $\varepsilon_t \sim \text{i.i.d. } N(0, \sigma^2)$ or more generally ε_t is a martingale difference. Conditioning in a given state and an information set \mathcal{F}_t , the process $\{y_t\}$ is linear but unconditionally the process is nonlinear. The conditional forecast is $\mathbb{E}(y_{t+1} | s_{t+1} = j, \mathcal{F}_t; \theta) = c_j + \phi_j y_t$ and the unconditional forecast based on observable variables is the sum of the conditional forecasts for each state weighted by the probability of being in that state,

$$\mathbb{E}(y_{t+1} | \mathcal{F}_t; \theta) = \sum_{j=1}^N P(s_{t+1} = j | \mathcal{F}_t; \theta) \mathbb{E}(y_{t+1} | s_{t+1} = j, \mathcal{F}_t; \theta).$$

The parameter vector $\theta = (c_1 \dots c_N, \phi_1 \dots \phi_N, \sigma^2)'$ as well as the transition probabilities p_{ij} can be estimated by maximum likelihood.

MS models have been applying to the modeling of foreign exchange rates with mixed success. Marsh (2000) estimates a two-state MS for the Deutschemark, the Pound Sterling, and the Japanese Yen. Though the model approximates the characteristics of the data well, the forecasting performance is poor when measured by the profit/losses generated by a set of trading rules based on the predictions of the MS model. On the contrary, Dueker and Neely (2007) find that for the same exchange rate a MS model with three states variables – in the scale factor of the variance of a Student-t error, in the kurtosis of the error, and in the expected return– produces out-of-sample excess returns that are slightly superior to those generated by common trading rules. For stock returns, there is evidence that MS models perform relatively well on describing two states in the mean (high/low returns) and two states in the variance (stable/volatile periods) of returns (Maheu and McCurdy, 2000). In addition, Perez-Quiros and Timmermann (2001) propose that the error term should be modelled as a mixture of Gaussian and Student-t distributions to capture the outliers commonly found in stock returns. This model provides some gains in predictive accuracy mainly for small firms returns. For interest rates in USA, Germany, and United Kingdom, Ang and Bekaert (2002) find that a two-state MS model that incorporates information on international short rate and on term spread is able to predict better than an univariate MS model. Additionally they find that in USA the classification of regimes correlates well with the business cycles.

SETAR, STAR, and MS models are successful specifications to approximate the characteristics of financial and macroeconomic data. However, good in-sample performance does not imply necessarily a good out-of-sample performance, mainly when compared to simple linear ARMA models. The success of nonlinear models depends on how prominent the nonlinearity is in the data. We should not expect a nonlinear model to perform better than a linear model when the contribution of the nonlinearity to the overall specification of the model is very small. As it is argued in Granger and Teräsvirta (1993), the prediction errors generated by a nonlinear model will be smaller only when the nonlinear feature modelled in-sample is also present in the forecasting sample.

3.5 A state dependent mixture model based on cross-sectional ranks

In the previous section, we have dealt with nonlinear time series models that only incorporate time series information. González-Rivera, Lee, and Mishra (2008) propose a nonlinear model that combines time series with cross sectional information. They propose the modelling of expected returns based on the joint dynamics of a sharp jump in the cross-sectional rank and the realized returns. They analyze the marginal probability distribution of a jump in the cross-sectional rank within the context of a duration model, and the probability of the asset return conditional on a jump specifying different dynamics depending on whether or not a jump has taken place. The resulting model for expected returns is a mixture of normal distributions weighted by the probability of jumping.

Let $y_{i,t}$ be the return of firm i at time t , and $\{y_{i,t}\}_{i=1}^M$ be the collection of asset returns of the M firms that constitute the *market* at time t . For each time t , the asset returns are ordered from the smallest to the largest, and define $z_{i,t}$, the *Varying Cross-sectional Rank* (VCR) of firm i within the market, as the proportion of firms that have a return less than or equal to the return of firm i . We write

$$z_{i,t} \equiv M^{-1} \sum_{j=1}^M \mathbf{1}(y_{j,t} \leq y_{i,t}), \quad (1)$$

where $\mathbf{1}(\cdot)$ is the indicator function, and for M large, $z_{i,t} \in (0, 1]$. Since the rank is a highly dependent variable, it is assumed that small movements in the asset ranking will not contain significant information and that most likely large movements in ranking will be the result of news in the overall market and/or of news concerning a particular asset. Focusing on large rank movements, we define, at time t , a sharp jump as a binary variable that takes the value one when there is a minimum (upward or downward) movement of 0.5 in the ranking of asset i , and zero otherwise:

$$J_{i,t} \equiv \mathbf{1}(|z_{i,t} - z_{i,t-1}| \geq 0.5). \quad (2)$$

A jump of this magnitude brings the asset return above or below the median of the cross-sectional distribution of returns. Note that this notion of jumps differs from the more traditional meaning of the word in the context of continuous-time modelling of the univariate

return process. A jump in the cross-sectional rank implicitly depends on numerous univariate return processes.

The analytical problem now consists in modeling the joint distribution of the return $y_{i,t}$ and the jump $J_{i,t}$, i.e. $f(y_{i,t}, J_{i,t}|\mathcal{F}_{t-1})$ where \mathcal{F}_{t-1} is the information set up to time $t-1$. Since $f(y_{i,t}, J_{i,t}|\mathcal{F}_{t-1}) = f_1(J_{i,t}|\mathcal{F}_{t-1})f_2(y_{i,t}|J_{i,t}, \mathcal{F}_{t-1})$, the analysis focuses first on the modelling of the marginal distribution of the jump, and subsequently on the modelling of the conditional distribution of the return.

Since $J_{i,t}$ is a Bernoulli variable, the marginal distribution of the jump is $f_1(J_{i,t}|\mathcal{F}_{t-1}) = p_{i,t}^{J_{i,t}}(1 - p_{i,t})^{(1-J_{i,t})}$ where $p_{i,t} \equiv \Pr(J_{i,t} = 1|\mathcal{F}_{t-1})$ is the conditional probability of a jump in the cross-sectional ranks. The modelling of $p_{i,t}$ is performed within the context of a dynamic duration model specified in calendar time as in Hamilton and Jordà (2002). The calendar time approach is necessary because asset returns are reported in calendar time (days, weeks, etc.) and it has the advantage of incorporating any other available information also reported in calendar time.

It is easy to see that the probability of jumping and duration must have an inverse relationship. If the probability of jumping is high, the expected duration must be short, and vice versa. Let $\Psi_{N(t)}$ be the expected duration. The expected duration until the next jump in the cross-sectional rank is given by $\Psi_{N(t)} = \sum_{j=1}^{\infty} j(1 - p_t)^{j-1}p_t = p_t^{-1}$. Note that $\sum_{j=0}^{\infty}(1 - p_t)^j = p_t^{-1}$. Differentiating with respect to p_t yields $\sum_{j=0}^{\infty} -j(1 - p_t)^{j-1} = -p_t^{-2}$. Multiplying by $-p_t$ gives $\sum_{j=0}^{\infty} j(1 - p_t)^{j-1}p_t = p_t^{-1}$ and thus $\sum_{j=1}^{\infty} j(1 - p_t)^{j-1}p_t = p_t^{-1}$. Consequently, to model $p_{i,t}$, it suffices to model the expected duration and compute its inverse. Following Hamilton and Jordà (2002), we specify an autoregressive conditional hazard (ACH) model. The ACH model is a calendar-time version of the autoregressive conditional duration (ACD) of Engle and Russell (1998). In both ACD and ACH models, the expected duration is a linear function of lag durations. However as the ACD model is set up in event time, there are some difficulties on how to introduce information that arrives between events. This is not the case in the ACH model because the set-up is in calendar time. In the ACD model, the forecasting object is the expected time between events; in the ACH model, the objective is to forecast the probability that the event will happen tomorrow

given the information known up to today. A general ACH model is specified as

$$\Psi_{N(t)} = \sum_{j=1}^m \alpha_j D_{N(t)-j} + \sum_{j=1}^r \beta_j \Psi_{N(t)-j}. \quad (3)$$

Since p_t is a probability, it must be bounded between zero and one. This implies that the conditional duration must have a lower bound of one. Furthermore, as we mentioned above, working in calendar time has the advantage that we can incorporate information that becomes available between jumps and can affect the probability of a jump in future periods. We specify the conditional hazard rate as

$$p_t = [\Psi_{N(t-1)} + \delta' X_{t-1}]^{-1}, \quad (4)$$

where X_{t-1} is a vector of relevant calendar time variables such as past VCRs and past returns. This completes the marginal distribution of the jump $f_1(J_{i,t}|\mathcal{F}_{t-1}) = p_{i,t}^{J_{i,t}}(1 - p_{i,t})^{(1-J_{i,t})}$.

On modelling $f_2(y_t|J_t, \mathcal{F}_{t-1}; \theta_2)$, it is assumed that the return to asset i may behave differently depending upon the occurrence of a jump. The modelling of two potential different states (whether a jump has occurred or not) will permit to differentiate whether the conditional expected return is driven by active or/and passive movements in the asset ranking in conjunction with its own return dynamics. *A priori*, different dynamics are possible in these two states. A general specification is

$$f_2(y_t|J_t, \mathcal{F}_{t-1}; \theta_2) = \begin{cases} N(\mu_{1,t}, \sigma_{1,t}^2) & \text{if } J_t = 1 \\ N(\mu_{0,t}, \sigma_{0,t}^2) & \text{if } J_t = 0, \end{cases} \quad (5)$$

where $\mu_{j,t}$ is the conditional mean and $\sigma_{j,t}^2$ the conditional variance in each state ($j = 1, 0$). Whether these two states are present in the data is an empirical question and it should be answered through statistical testing.

Combining the models for the marginal density of the jump and the conditional density of the returns, the estimation can be conducted with maximum likelihood techniques. For a sample $\{y_t, J_t\}_{t=1}^T$, the joint log-likelihood function is

$$\sum_{t=1}^T \ln f(y_t, J_t|\mathcal{F}_{t-1}; \theta) = \sum_{t=1}^T \ln f_1(J_t|\mathcal{F}_{t-1}; \theta_1) + \sum_{t=1}^T \ln f_2(y_t|J_t, \mathcal{F}_{t-1}; \theta_2).$$

Let us call $\mathcal{L}_1(\theta_1) = \sum_{t=1}^T \ln f_1(J_t|\mathcal{F}_{t-1}; \theta_1)$ and $\mathcal{L}_2(\theta_2) = \sum_{t=1}^T \ln f_2(y_t|J_t, \mathcal{F}_{t-1}; \theta_2)$. The maximization of the joint log-likelihood function can be achieved by maximizing $\mathcal{L}_1(\theta_1)$ and $\mathcal{L}_2(\theta_2)$ separately without loss of efficiency by assuming that the parameter vectors θ_1 and θ_2 are “variation free” in the sense of Engle *et al* (1983).

The log-likelihood function $\mathcal{L}_1(\theta_1) = \sum_{t=1}^T \ln f_1(J_t|\mathcal{F}_{t-1}; \theta_1)$ is

$$\mathcal{L}_1(\theta_1) = \sum_{t=1}^T [J_t \ln p_t(\theta_1) + (1 - J_t) \ln(1 - p_t(\theta_1))], \quad (6)$$

where θ_1 includes all parameters in the conditional duration model.

The log-likelihood function $\mathcal{L}_2(\theta_2) = \sum_{t=1}^T \ln f_2(y_t|J_t, \mathcal{F}_{t-1}; \theta_2)$ is

$$\mathcal{L}_2(\theta_2) = \sum_{t=1}^T \ln \left[\frac{J_t}{\sqrt{2\pi\sigma_{1,t}^2}} \exp \left\{ -\frac{1}{2} \left(\frac{y_t - \mu_{1,t}}{\sigma_{1,t}} \right)^2 \right\} + \frac{1 - J_t}{\sqrt{2\pi\sigma_{0,t}^2}} \exp \left\{ -\frac{1}{2} \left(\frac{y_t - \mu_{0,t}}{\sigma_{0,t}} \right)^2 \right\} \right],$$

where θ_2 includes all parameters in the conditional means and conditional variances under both regimes.

If the two proposed states are granted in the data, the marginal density function of the asset return must be a mixture of two normal density functions where the mixture weights are given by the probability of jumping p_t :

$$\begin{aligned} g(y_t|\mathcal{F}_{t-1}; \theta) &\equiv \sum_{J_t=0}^1 f(y_t, J_t|\mathcal{F}_{t-1}; \theta) \\ &= \sum_{J_t=0}^1 f_1(J_t|\mathcal{F}_{t-1}; \theta_1) f_2(y_t|J_t, \mathcal{F}_{t-1}; \theta_2) \\ &= p_t \cdot f_2(y_t|J_t = 1, \mathcal{F}_{t-1}; \theta_2) + (1 - p_t) \cdot f_2(y_t|J_t = 0, \mathcal{F}_{t-1}; \theta_2), \end{aligned} \quad (7)$$

as $f_1(J_t|\mathcal{F}_{t-1}; \theta_1) = p_t^{J_t} (1 - p_t)^{(1-J_t)}$. Therefore, the one-step ahead forecast of the return is

$$\mathbb{E}(y_{t+1}|\mathcal{F}_t; \theta) = \int y_{t+1} \cdot g(y_{t+1}|\mathcal{F}_t; \theta) dy_{t+1} = p_{t+1}(\theta_1) \cdot \mu_{1,t+1}(\theta_2) + (1 - p_{t+1}(\theta_1)) \cdot \mu_{0,t+1}(\theta_2). \quad (8)$$

The expected return is a function of the probability of jumping p_t , which is a nonlinear function of the information set as shown in (4). Hence the expected returns are nonlinear functions of the information set, even in a simple case where $\mu_{1,t}$ and $\mu_{0,t}$ are linear.

This model was estimated for the returns of the constituents of the SP500 index from 1990 to 2000, and its performance was assessed in an out-of-sample exercise from 2001 to 2005 within the context of several trading strategies. Based on the one-step-ahead forecast of the mixture model, a proposed trading strategy called VCR-Mixture Trading Rule is shown to be a superior rule because of its ability to generate large risk-adjusted mean returns when compared to other technical and model-based trading rules. The VCR-Mixture Trading Rule is implemented by computing for each firm in the SP500 index the one-step ahead forecast of the return as in (8). Based on the forecasted returns $\{\hat{y}_{i,t+1}(\hat{\theta}_t)\}_{t=R}^{T-1}$, the investor predicts the VCR of all assets in relation to the overall market, that is,

$$\hat{z}_{i,t+1} = M^{-1} \sum_{j=1}^M \mathbf{1}(\hat{y}_{j,t+1} \leq \hat{y}_{i,t+1}), \quad t = R, \dots, T-1, \quad (9)$$

and buys the top K performing assets if their forecasted return is above the risk-free rate. In every subsequent out-of-sample period ($t = R, \dots, T-1$), the investor revises her portfolio, selling the assets that fall out of the top performers and buying the ones that rise to the top, and she computes the one-period portfolio return

$$\pi_{t+1} = K^{-1} \sum_{j=1}^M y_{j,t+1} \cdot \mathbf{1}(\hat{z}_{j,t+1} \geq z_{t+1}^K), \quad t = R, \dots, T-1, \quad (10)$$

where z_{t+1}^K is the cutoff cross-sectional rank to select the K best performing stocks such that $\sum_{j=1}^M \mathbf{1}(\hat{z}_{j,t+1} \geq z_{t+1}^K) = K$. In the analysis of González-Rivera, Lee, and Mishra (2008) a portfolio is formed with the top 1% ($K = 5$ stocks) performers in the SP500 index. Every asset in the portfolio is weighted equally. The evaluation criterion is to compute the “mean trading return” over the forecasting period

$$MTR = P^{-1} \sum_{t=R}^{T-1} \pi_{t+1}.$$

It is also possible to correct MTR according to the level of risk of the chosen portfolio. For instance, the traditional Sharpe ratio will provide the excess return per unit of risk measured by the standard deviation of the selected portfolio

$$SR = P^{-1} \sum_{t=R}^{T-1} \frac{(\pi_{t+1} - r_{f,t+1})}{\sigma_{t+1}^{\pi}(\hat{\theta}_t)},$$

where $r_{f,t+1}$ is the risk free rate. The VCR-Mixture Trading Rule produces a weekly *MTR* of 0.243% (63.295% cumulative return over 260 weeks), equivalent to a yearly compounded return of 13.45%, that is significantly more than the next most favorable rule, which is the Buy-and-Hold-the-Market Trading Rule with a weekly mean return of -0.019% , equivalent to a yearly return of -1.00% . To assess the return-risk trade off, we implement the Sharpe ratio. The largest *SR* (mean return per unit of standard deviation) is provided by the VCR-Mixture rule with a weekly return of 0.151% (8.11% yearly compounded return per unit of standard deviation), which is lower than the mean return provided by the same rule under the *MTR* criterion, but still a dominant return when compared to the mean returns provided by the Buy-and-Hold-the-Market Trading Rule.

3.6 Random fields

Hamilton (2001) proposed a flexible parametric regression model where the conditional mean has a linear parametric component and a potential nonlinear component represented by an isotropic Gaussian random field. The model has a nonparametric flavor because no functional form is assumed but, nevertheless, the estimation is fully parametric.

A scalar random field is defined as a function $m(\omega, x) : \Omega \times A \rightarrow R$ such that $m(\omega, x)$ is a random variable for each $x \in A$ where $A \subseteq R^k$. A random field is also denoted as $m(x)$. If $m(x)$ is a system of random variables with finite dimensional Gaussian distributions, then the scalar random field is said to be Gaussian and it is completely determined by its mean function $\mu(x) = \mathbb{E}[m(x)]$ and its covariance function with typical element $C(x, z) = \mathbb{E}[(m(x) - \mu(x))(m(z) - \mu(z))]$ for any $x, z \in A$. The random field is said to be homogeneous or stationary if $\mu(x) = \mu$ and the covariance function depends only on the difference vector $x - z$ and we should write $C(x, z) = C(x - z)$. Furthermore, the random field is said to be isotropic if the covariance function depends on $d(x, z)$, where $d(\cdot)$ is a scalar measure of distance. In this situation we write $C(x, z) = C(d(x, z))$.

The specification suggested by Hamilton (2001) can be represented as

$$y_t = \beta_0 + x_t' \beta_1 + \lambda m(g \odot x_t) + \epsilon_t, \quad (11)$$

for $y_t \in R$ and $x_t \in R^k$, both stationary and ergodic processes. The conditional mean has

a linear component given by $\beta_0 + x_t'\beta_1$ and a nonlinear component given by $\lambda m(g \odot x_t)$, where $m(z)$, for any choice of z , represents a realization of a Gaussian and homogenous random field with a moving average representation; x_t could be predetermined or exogenous and is independent of $m(\cdot)$, and ϵ_t is a sequence of independent and identically distributed $N(0, \sigma^2)$ variates independent of both $m(\cdot)$ and x_t as well as of lagged values of x_t . The scalar parameter λ represents the contribution of the nonlinear part to the conditional mean, the vector $g \in R_{0,+}^k$ drives the curvature of the conditional mean, and the symbol \odot denotes element-by-element multiplication.

Let H_k be the covariance (correlation) function of the random field $m(\cdot)$ with typical element defined as $H_k(x, z) = \mathbb{E}[m(x)m(z)]$. Hamilton (2001) proved that the covariance function depends solely upon the Euclidean distance between x and z , rendering the random field isotropic. For any x and $z \in R^k$, the correlation between $m(x)$ and $m(z)$ is given by the ratio of the volume of the overlap of k -dimensional unit spheroids centered at x and z to the volume of a single k -dimensional unit spheroid. If the Euclidean distance between x and z is greater than two, the correlation between $m(x)$ and $m(z)$ will be equal to zero. The general expression of the correlation function is

$$H_k(h) = \begin{cases} G_{k-1}(h, 1)/G_{k-1}(0, 1) & \text{if } h \leq 1 \\ 0 & \text{if } h > 1 \end{cases}, \quad (12)$$

$$G_k(h, r) = \int_h^r (r^2 - w^2)^{k/2} dw,$$

where $h \equiv \frac{1}{2}d_{L_2}(x, z)$, and $d_{L_2}(x, z) \equiv [(x - z)'(x - z)]^{1/2}$ is the Euclidean distance between x and z .

Within the specification (11), Dahl and González-Rivera (2003a) provided alternative representations of the random field that permit the construction of Lagrange multiplier tests for neglected nonlinearity, which circumvent the problem of unidentified nuisance parameters under the null of linearity and, at the same time, they are robust to the specification of the covariance function associated with the random field. They modified the Hamilton framework in two directions. First, the random field is specified in the L_1 norm instead of the L_2 norm, and secondly they considered random fields that may not have a simple moving average representation. The advantage of the L_1 norm, which is exploited in the testing problem,

is that this distance measure is a linear function of the nuisance parameters, in contrast to the L_2 norm which is a nonlinear function. Logically, Dahl and González-Rivera proceeded in an opposite fashion to Hamilton. Whereas Hamilton first proposed a moving average representation of the random field, and secondly, he derived its corresponding covariance function, Dahl and González-Rivera first proposed a covariance function, and secondly they inquire whether there is a random field associated with it. The proposed covariance function is

$$C_k(h^*) = \begin{cases} (1 - h^*)^{2k} & \text{if } h^* \leq 1 \\ 0 & \text{if } h^* > 1 \end{cases}, \quad (13)$$

where $h^* \equiv \frac{1}{2}d_{L_1}(x, z) = \frac{1}{2}|x - z|'1$. The function (13) is a permissible covariance, that is, it satisfies the positive semidefiniteness condition, which is $q'C_kq \geq 0$ for all $q \neq 0_T$. Furthermore, there is a random field associated with it according to the Khinchin's theorem (1934) and Bochner's theorem (1959). The basic argument is that the class of functions which are covariance functions of homogenous random fields coincides with the class of positive semidefinite functions. Hence, (13) being a positive semidefinite function must be the covariance function of a homogenous random field.

The estimation of these models is carried out by maximum likelihood. From model (11), we can write $y \sim N(X\beta, \lambda^2 C_k + \sigma^2 I_T)$ where $y = (y_1, y_2, \dots, y_T)'$, $X_1 = (x'_1, x'_2, \dots, x'_T)'$, $X = (\mathbf{1} : X_1)$, $\beta = (\beta_0, \beta_1)'$, $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_T)'$ and σ^2 is the variance of ϵ_t . C_k is a generic covariance function associated with the random field, which could be equal to the Hamilton spherical covariance function in (12), or to the covariance in (13). The log-likelihood function corresponding to this model is

$$\begin{aligned} \ell(\beta, \lambda^2, g, \sigma^2) &= -\frac{T}{2} \log(2\pi) - \frac{1}{2} \log |\lambda^2 C_k + \sigma^2 I_T| \\ &\quad - \frac{1}{2} (y - X\beta)' (\lambda^2 C_k + \sigma^2 I_T)^{-1} (y - X\beta). \end{aligned} \quad (14)$$

The flexible regression model has been applied successfully to detect nonlinearity in the quarterly growth rate of the US real GNP (Dahl and González-Rivera 2003b) and in the Industrial Production Index of sixteen OECD countries (Dahl and González-Rivera 2003a). This technology is able to mimic the characteristics of the actual US business cycle. The cycle is dissected according to measures of duration, amplitude, cumulation and excess cumulation

of the contraction and expansion phases. In contrast to Harding and Pagan (2002) who find that nonlinear models are not uniformly superior to linear ones, the flexible regression model represents a clear improvement over linear models, and it seems to capture just the right shape of the expansion phase as opposed to Hamilton (1989) and Durland and McCurdy (1994) models, which tend to overestimate the cumulation measure in the expansion phase. It is found that the expansion phase must have at least two subphases: an aggressive early expansion after the trough, and a moderate/slow late expansion before the peak implying the existence of an inflexion point that we date approximately around one-third into the duration of the expansion phase. This shape lends support to parametric models of the growth rate that allow for three regimes (Sichel, 1994), as opposed to models with just two regimes (contractions and expansions). For the Industrial Production Index, testing for nonlinearity within the flexible regression framework brings similar conclusions to those in Teräsvirta and Anderson (1992), who propose parametric STAR models for industrial production data. However, the tests proposed in Dahl and González-Rivera (2003a), which have superior performance to detect smooth transition dynamics, seem to indicate that linearity cannot be rejected in the industrial production indexes of Japan, Austria, Belgium and Sweden as opposed to the findings of Teräsvirta and Anderson.

3.7 Nonlinear factor models

For the last ten years forecasting using vast data in data-rich environment has been one of the most researched topic in forecasting in economics and finance. See Stock and Watson (2002, 2006). These factor models in this literature are mostly linear models. Bai and Ng (BN, 2007) introduce a nonlinear factor model, and a quadratic principal component model as a special case. First consider a simple factor model

$$x_{it} = \lambda_i' F_t + e_{it}. \quad (15)$$

By the method of principal component, the elements of \mathbf{f}_t are linear combinations of elements of \mathbf{x}_t . The factors are estimated by minimizing the sum of squared residuals of the linear model, $x_{it} = \lambda_i' F_t + e_{it}$.

The factor model in (15) presupposes a linear link function between the predictor \mathbf{x}_t and the latent factors F_t . BN consider a more flexible approach by a nonlinear link function $g(\cdot)$ such that

$$g(x_{it}) = \phi_i' J_t + v_{it},$$

where J_t are the common factors, and ϕ_i is the vector of factor loadings. BN consider $g(x_{it})$ to be x_{it} augmented by some or all of the unique cross-products of the elements of $\{x_{it}\}_{i=1}^N$. The second-order factor model is then $x_{it}^* = \phi_i' J_t + v_{it}$ where x_{it}^* is an $N^* \times 1$ vector. Estimation of J_t then proceeds by the usual method of principal components. BN consider $x_{it}^* = \{x_{it} \ x_{it}^2\}_{i=1}^N$ with $N^* = 2N$, which they call the SPC (squared principal components).

Once the factors are estimated, the forecasting equation for y_{t+h} would be

$$y_{t+h} = (1 \ \hat{F}_t') \boldsymbol{\gamma} + \varepsilon_t.$$

The forecasting equation remains linear whatever the link function g is. An alternative way of capturing nonlinearity is to augment the forecasting equation to include functions of the factors

$$y_{t+h} = (1 \ \hat{F}_t') \boldsymbol{\gamma} + a(\hat{F}_t) + \varepsilon_t,$$

where $a(\cdot)$ is nonlinear. A simple case when $a(\cdot)$ is quadratic is referred to as PC2 (squared factors) in BN.

BN note that the PC2 is conceptually distinct from SPC. While the PC2 forecasting model allows the volatility of factors estimated by linear principal components to have predictive power for y , the SPC model allows the factors to be possibly nonlinear functions of the predictors while maintaining a linear relation between the factors and y . Ludvigson and Ng (2005) found that the square of the first factor estimated from a set of financial factors (i.e., volatility of the first factor) is significant in the regression model for the mean excess returns. In contrast, factors estimated from the second moment of data (i.e., volatility factors) are much weaker predictors of excess returns. This contrasts the predictive ability of the nonlinear models for the excess returns rather than usual interpretation of the risk premium.

3.8 Artificial neural network models

Consider an augmented single hidden layer feedforward neural network model $f(x_t, \theta)$ in which the network output y_t is determined given input x_t as

$$\begin{aligned} y_t &= f(x_t, \theta) + \varepsilon_t \\ &= x_t \beta + \sum_{j=1}^q \delta_j \psi(x_t \gamma_j) + \varepsilon_t \end{aligned}$$

where $\theta = (\beta' \ \gamma' \ \delta)'$, β is a conformable column vector of connection strength from the input layer to the output layer; γ_j is a conformable column vector of connection strength from the input layer to the hidden units, $j = 1, \dots, q$; δ_j is a (scalar) connection strength from the hidden unit j to the output unit, $j = 1, \dots, q$; and ψ is a squashing function (e.g., the logistic squasher) or a radial basis function. Input units x send signals to intermediate hidden units, then each of hidden unit produces an activation ψ that then sends signals toward the output unit. The integer q denotes the number of hidden units added to the affine (linear) network. When $q = 0$, we have a two layer *affine* network $y_t = x_t \beta + \varepsilon_t$. Hornick, Stinchcombe and White (1989) show that neural network is a nonlinear flexible functional form being capable of approximating any Borel measurable function to any desired level of accuracy provided sufficiently many hidden units are available. Stinchcombe and White (1998) show that this result holds for any $\psi(\cdot)$ belonging to the class of “generically comprehensively revealing” functions. These functions are “comprehensively revealing” in the sense that they can reveal arbitrary model misspecifications $\mathbb{E}(y_t|x_t) \neq f(x_t, \theta^*)$ with non-zero probability and they are “generic” in the sense that almost any choice for γ will reveal the misspecification.

We build an artificial neural network (ANN) model based on a test for neglected non-linearity likely to have power against a range of alternatives . See White (1989) and Lee, White, and Granger (1993) on the neural network test and its comparison with other specification tests. The neural network test is based on a test function $h(x_t)$ chosen as the activations of ‘phantom’ hidden units $\psi(x_t \Gamma_j)$, $j = 1, \dots, q$, where Γ_j are random column vectors independent of x_t . That is,

$$\mathbb{E}[\psi(x_t \Gamma_j) \varepsilon_t^* | \Gamma_j] = \mathbb{E}[\psi(x_t \Gamma_j) \varepsilon_t^*] = 0 \quad j = 1, \dots, q, \quad (16)$$

under H_0 , so that

$$\mathbb{E}(\Psi_t \varepsilon_t^*) = 0, \quad (17)$$

where $\Psi_t = (\psi(x_t \Gamma_1), \dots, \psi(x_t \Gamma_q))'$ is a phantom hidden unit activation vector. Evidence of correlation of ε_t^* with Ψ_t is evidence against the null hypothesis that y_t is linear in mean. If correlation exists, augmenting the linear network by including an additional hidden unit with activations $\psi(x_t \Gamma_j)$ would permit an improvement in network performance. Thus the tests are based on sample correlation of affine network errors with phantom hidden unit activations,

$$n^{-1} \sum_{t=1}^n \Psi_t \hat{\varepsilon}_t = n^{-1} \sum_{t=1}^n \Psi_t (y_t - x_t \hat{\beta}). \quad (18)$$

Under suitable regularity conditions it follows from the central limit theorem that $n^{-1/2} \sum_{t=1}^n \Psi_t \hat{\varepsilon}_t \xrightarrow{d} N(0, W^*)$ as $n \rightarrow \infty$, and if one has a consistent estimator for its asymptotic covariance matrix, say \hat{W}_n , then an asymptotic chi-square statistic can be formed as

$$(n^{-1/2} \sum_{t=1}^n \Psi_t \hat{\varepsilon}_t)' \hat{W}_n^{-1} (n^{-1/2} \sum_{t=1}^n \Psi_t \hat{\varepsilon}_t) \xrightarrow{d} \chi^2(q). \quad (19)$$

Elements of Ψ_t tend to be collinear with X_t and with themselves. Thus LWG conduct a test on $q^* < q$ principal components of Ψ_t not collinear with x_t , denoted Ψ_t^* . This test is to determine whether or not there exists some advantage to be gained by adding hidden units to the affine network. We can estimate \hat{W}_n robust to the conditional heteroskedasticity, or we may use with the empirical null distribution of the statistic computed by a bootstrap procedure that is robust to the conditional heteroskedasticity, e.g., wild bootstrap.

Estimation of an ANN model may be tedious and sometimes results in unreliable estimates. Recently, White (2006) proposes a simple algorithm called QuickNet, a form of “relaxed greedy algorithm” because QuickNet searches for a single best additional hidden unit based on a sequence of OLS regressions, that may be analogous to the least angular regressions (LARS) of Efron, Hastie, Johnstone, and Tibshirani (2004). The simplicity of the QuickNet algorithm achieves the benefits of using a forecasting model that is nonlinear in the predictors while mitigating the other computational challenges to the use of nonlinear forecasting methods. See White (2006, Section 5) for more details on QuickNet, and for

other issues of controlling for overfit and the selection of the random parameter vectors Γ_j independent of x_t .

Campbell, Lo, and MacKinlay (1997, Section 12.4) provide a review to start in the area. White (2006) is an excellent place to see how the research has reached. Trippi and Turban (1992) review the applications of ANNs to finance and investment.

3.9 Functional coefficient models

A functional coefficient model is introduced by Cai, Fan, and Yao (CFY, 2000), with time-varying and state-dependent coefficients. It can be viewed as a special case of Priestley’s (1980) state-dependent model, but it includes the models of Tong (1990), Chen and Tsay (1993) and regime-switching models as special cases. Let $\{(y_t, s_t)'\}_{t=1}^n$ be a stationary process, where y_t and s_t are scalar variables. Also let $X_t \equiv (1, y_{t-1}, \dots, y_{t-d})'$. We assume

$$\mathbb{E}(y_t | \mathcal{F}_{t-1}) = a_0(s_t) + \sum_{j=1}^d a_j(s_t) y_{t-j},$$

where the $\{a_j(s_t)\}$ are the autoregressive coefficients depending on s_t , which may be chosen as a function of X_t or something else. Intuitively, the functional coefficient model is an AR process with time-varying autoregressive coefficients. The coefficient functions $\{a_j(s_t)\}$ can be estimated by local linear regression. At each point s , we approximate $a_j(s_t)$ locally by a linear function $a_j(s_t) \approx a_j + b_j(s_t - s)$, $j = 0, 1, \dots, d$, for s_t near s , where a_j and b_j are constants. The local linear estimator at point s is then given by $\hat{a}_j(s) = \hat{a}_j$, where $\{(\hat{a}_j, \hat{b}_j)\}_{j=0}^d$ minimizes the sum of local weighted squares $\sum_{t=1}^n [y_t - \mathbb{E}(y_t | \mathcal{F}_{t-1})]^2 K_h(s_t - s)$, with $K_h(\cdot) \equiv K(\cdot/h)/h$ for a given kernel function $K(\cdot)$ and bandwidth $h \equiv h_n \rightarrow 0$ as $n \rightarrow \infty$. CFY (2000, p. 944) suggest to select h using a modified multi-fold “leave-one-out-type” cross-validation based on MSFE.

It is important to choose an appropriate smooth variable s_t . Knowledge on data or economic theory may be helpful. When no prior information is available, s_t may be chosen as a function of explanatory vector X_t or using such data-driven methods as AIC and cross-validation. See Fan, Yao and Cai (2003) for further discussion on the choice of s_t . For exchange rate changes, Hong and Lee (2003) choose s_t as the difference between the exchange

rate at time $t - 1$ and the moving average of the most recent L periods of exchange rates at time $t - 1$. The moving average is a proxy for the local trend at time $t - 1$. Intuitively, this choice of s_t is expected to reveal useful information on the direction of changes.

To justify the use of the functional coefficient model, CFY (2000) suggest a goodness-of-fit test for an AR(d) model against a functional coefficient model. The null hypothesis of AR(d) can be stated as

$$\mathbb{H}_0 : a_j(s_t) = \beta_j, \quad j = 0, 1, \dots, d,$$

where β_j is the autoregressive coefficient in AR(d). Under \mathbb{H}_0 , $\{y_t\}$ is linear in mean conditional on X_t . Under the alternative to \mathbb{H}_0 , the autoregressive coefficients depend on s_t and the AR(d) model suffers from “neglected nonlinearity”. To test \mathbb{H}_0 , CFY compares the residual sum of squares (RSS) under \mathbb{H}_0

$$RSS_0 \equiv \sum_{t=1}^n \hat{\varepsilon}_t^2 = \sum_{t=1}^n [Y_t - \hat{\beta}_0 - \sum_{j=1}^d \hat{\beta}_j Y_{t-j}]^2$$

with the RSS under the alternative

$$RSS_1 \equiv \sum_{t=1}^n \tilde{\varepsilon}_t^2 = \sum_{t=1}^n [Y_t - \hat{a}_0(s_t) - \sum_{j=1}^d \hat{a}_j(s_t) Y_{t-j}]^2.$$

The test statistic is $T_n = (RSS_0 - RSS_1)/RSS_1$. We reject \mathbb{H}_0 for large values of T_n . CFY suggest the following bootstrap method to obtain the p -value of T_n : (i) generate the bootstrap residuals $\{\varepsilon_t^b\}_{t=1}^n$ from the centered residuals $\tilde{\varepsilon}_t - \bar{\varepsilon}$ where $\bar{\varepsilon} \equiv n^{-1} \sum_{t=1}^n \tilde{\varepsilon}_t$ and define $y_t^b \equiv X_t' \hat{\beta} + \varepsilon_t^b$, where $\hat{\beta}$ is the OLS estimator for AR(d); (ii) calculate the bootstrap statistic T_n^b using the bootstrap sample $\{y_t^b, X_t', s_t\}_{t=1}^n$; (iii) repeat steps (i) and (ii) B times ($b = 1, \dots, B$) and approximate the bootstrap p -value of T_n by $B^{-1} \sum_{b=1}^B \mathbf{1}(T_n^b \geq T_n)$. See Hong and Lee (2003) for empirical application of the functional coefficient model to forecasting foreign exchange rates.

3.10 Nonparametric regression

Let $\{y_t, x_t\}, t = 1, \dots, n$, be stochastic processes, where y_t is a scalar and $x_t = (x_{t1}, \dots, x_{tk})$ is a $1 \times k$ vector which may contain the lagged values of y_t . Consider the regression model

$$y_t = m(x_t) + u_t$$

where $m(x_t) = \mathbb{E}(y_t|x_t)$ is the true but unknown regression function and u_t is the error term such that $\mathbb{E}(u_t|x_t) = 0$.

If $m(x_t) = g(x_t, \delta)$ is a correctly specified family of parametric regression functions then $y_t = g(x_t, \delta) + u_t$ is a correct model and, in this case, one can construct a consistent least squares (LS) estimator of $m(x_t)$ given by $g(x_t, \hat{\delta})$, where $\hat{\delta}$ is the LS estimator of the parameter δ .

In general, if the parametric regression $g(x_t, \delta)$ is incorrect or the form of $m(x_t)$ is unknown then $g(x_t, \hat{\delta})$ may not be a consistent estimator of $m(x_t)$. For this case, an alternative approach to estimate the unknown $m(x_t)$ is to use the consistent nonparametric kernel regression estimator which is essentially a local constant LS (LCLS) estimator. To obtain this estimator take Taylor series expansion of $m(x_t)$ around x so that

$$\begin{aligned} y_t &= m(x_t) + u_t \\ &= m(x) + e_t \end{aligned}$$

where $e_t = (x_t - x)m^{(1)}(x) + \frac{1}{2}(x_t - x)^2m^{(2)}(x) + \dots + u_t$ and $m^{(s)}(x)$ represents the s -th derivative of $m(x)$ at $x_t = x$. The LCLS estimator can then be derived by minimizing

$$\sum_{t=1}^n e_t^2 K_{tx} = \sum_{t=1}^n (y_t - m(x))^2 K_{tx}$$

with respect to constant $m(x)$, where $K_{tx} = K\left(\frac{x_t - x}{h}\right)$ is a decreasing function of the distances of the regressor vector x_t from the point $x = (x_1, \dots, x_k)$, and $h \rightarrow 0$ as $n \rightarrow \infty$ is the window width (smoothing parameter) which determines how rapidly the weights decrease as the distance of x_t from x increases. The LCLS estimator so estimated is

$$\hat{m}(x) = \frac{\sum_{t=1}^n y_t K_{tx}}{\sum_{t=1}^n K_{tx}} = (\mathbf{i}'\mathbf{K}(x)\mathbf{i})^{-1}\mathbf{i}'\mathbf{K}(x)\mathbf{y}$$

where $\mathbf{K}(x)$ is the $n \times n$ diagonal matrix with the diagonal elements K_{tx} ($t = 1, \dots, n$), \mathbf{i} is an $n \times 1$ column vector of unit elements, and \mathbf{y} is an $n \times 1$ vector with elements y_t ($t = 1, \dots, n$).

The estimator $\hat{m}(x)$ is due to Nadaraya (1964) and Watson (1964) (NW) who derived this in an alternative way. Generally $\hat{m}(x)$ is calculated at the data points x_t , in which case we can write the leave-one out estimator as

$$\hat{m}(x) = \frac{\sum_{t'=1, t' \neq t}^n y_{t'} K_{t't}}{\sum_{t'=1, t' \neq t}^n K_{t't}},$$

where $K_{t't} = K\left(\frac{x_{t'} - x_t}{h}\right)$. The assumption that $h \rightarrow 0$ as $n \rightarrow \infty$ gives $x_t - x = O(h) \rightarrow 0$ and hence $\mathbb{E}e_t \rightarrow 0$ as $n \rightarrow \infty$. Thus the estimator $\hat{m}(x)$ will be consistent under certain smoothing conditions on h, K , and $m(x)$. In small samples however $\mathbb{E}e_t \neq 0$ so $\hat{m}(x)$ will be a biased estimator, see Pagan and Ullah (1999) for details on asymptotic and small sample properties.

An estimator which has a better small sample bias and hence the mean square error (MSE) behavior is the local linear LS (LLS) estimator. In the LLS estimator we take first order Taylor-Series expansion of $m(x_t)$ around x so that

$$\begin{aligned} y_t &= m(x_t) + u_t = m(x) + (x_t - x)m^{(1)}(x) + v_t \\ &= \alpha(x) + x_t\beta(x) + v_t \\ &= X_t\delta(x) + v_t \end{aligned}$$

where $X_t = (1 \ x_t)$ and $\delta(x) = [\alpha(x) \ \beta(x)]'$ with $\alpha(x) = m(x) - x\beta(x)$ and $\beta(x) = m^{(1)}(x)$. The LLS estimator of $\delta(x)$ is then obtained by minimizing

$$\sum_{t=1}^n v_t^2 K_{tx} = \sum_{t=1}^n (y_t - X_t\delta(x))^2 K_{tx}$$

and it is given by

$$\tilde{\delta}(x) = (\mathbf{X}'\mathbf{K}(x)\mathbf{X})^{-1}\mathbf{X}'\mathbf{K}(x)\mathbf{y}. \quad (20)$$

where \mathbf{X} is an $n \times (k + 1)$ matrix with the t th row X_t ($t = 1, \dots, n$).

The LLS estimator of $\alpha(x)$ and $\beta(x)$ can be calculated as $\tilde{\alpha}(x) = (1 \ 0)\tilde{\delta}(x)$ and $\tilde{\beta}(x) = (0 \ 1)\tilde{\delta}(x)$. This gives

$$\tilde{m}(x) = (1 \ x)\tilde{\delta}(x) = \tilde{\alpha}(x) + x\tilde{\beta}(x).$$

Obviously when $X = \mathbf{i}$, $\tilde{\delta}(x)$ reduces to the NW's LCLS estimator $\hat{m}(x)$. An extension of the LLS is the local polynomial LS (LPLS) estimators, see Fan and Gijbels (1996).

In fact one can obtain the local estimators of a general nonlinear model $g(x_t, \delta)$ by minimizing

$$\sum_{t=1}^n [y_t - g(x_t, \delta(x))]^2 K_{tx}$$

with respect to $\delta(x)$. For $g(x_t, \delta(x)) = X_t \delta(x)$ we get the LLLS in (20). Further when $h = \infty, K_{tx} = K(0)$ is a constant so that the minimization of $K(0) \sum [y_t - g(x_t, \delta(x))]^2$ is the same as the minimization of $\sum [y_t - g(x_t, \delta)]^2$, that is the local LS becomes the global LS estimator $\hat{\delta}$.

The LLLS estimator in (20) can also be interpreted as the estimator of the functional coefficient (varying coefficient) linear regression model

$$\begin{aligned} y_t &= m(x_t) + u_t \\ &= X_t \delta(x_t) + u_t \end{aligned}$$

where $\delta(x_t)$ is approximated locally by a constant $\delta(x_t) \simeq \delta(x)$. The minimization of $\sum u_t^2 K_{tx}$ with respect to $\delta(x)$ then gives the LLLS estimator in (20), which can be interpreted as the LC varying coefficient estimator. An extension of this is to consider the linear approximation $\delta(x_t) \simeq \delta(x) + D(x)(x_t - x)'$ where $D(x) = \frac{\partial \delta(x_t)}{\partial x_t'}$ evaluated at $x_t = x$. In this case

$$\begin{aligned} y_t &= m(x_t) + u_t = X_t \delta(x_t) + u_t \\ &\simeq X_t \delta(x) + X_t D(x)(x_t - x)' + u_t \\ &= X_t \delta(x) + [(x_t - x) \otimes X_t] \text{vec} D(x) + u_t \\ &= X_t^x \delta^x(x) + u_t \end{aligned}$$

where $X_t^x = [X_t \ (x_t - x) \otimes X_t]$ and $\delta^x(x) = [\delta(x)' \ (\text{vec} D(x))']'$. The LL varying coefficient estimator of $\delta^x(x)$ can then be obtained by minimizing

$$\sum_{t=1}^n [y_t - X_t^x \delta^x(x)]^2 K_{tx}$$

with respect to $\delta^x(x)$ as

$$\dot{\delta}^x(x) = (\mathbf{X}^{x'} \mathbf{K}(x) \mathbf{X}^x)^{-1} \mathbf{X}^{x'} \mathbf{K}(x) \mathbf{y}. \quad (21)$$

From this $\dot{\delta}(x) = (\mathbf{I} \ 0) \dot{\delta}^x(x)$, and hence

$$\dot{m}(x) = (1 \ x \ 0) \dot{\delta}^x(x) = (1 \ x) \dot{\delta}(x).$$

The above idea can be extended to the situations where $\xi_t = (x_t \ z_t)$ such that

$$\mathbb{E}(y_t | \xi_t) = m(\xi_t) = m(x_t, z_t) = X_t \delta(z_t),$$

where the coefficients are varying with respect to only a subset of ξ_t ; z_t is $1 \times l$ and ξ_t is $1 \times p$, $p = k + l$. Examples of these include functional coefficient autoregressive models of Chen and Tsay (1993) and CFY (2000), random coefficient models of Raj and Ullah (1981), smooth transition autoregressive models of Granger and Teräsvirta (1993), and threshold autoregressive models of Tong (1990).

To estimate $\delta(z_t)$ we can again do a local constant approximation $\delta(z_t) \simeq \delta(z)$ and then minimize $\sum [y_t - X_t \delta(z)]^2 K_{tz}$ with respect to $\delta(z)$, where $K_{tz} = K(\frac{z_t - z}{h})$. This gives the LC varying coefficient estimator

$$\tilde{\delta}(z) = (\mathbf{X}'\mathbf{K}(z)\mathbf{X})^{-1}\mathbf{X}'\mathbf{K}(z)\mathbf{y} \quad (22)$$

where $\mathbf{K}(z)$ is a diagonal matrix of K_{tz} , $t = 1, \dots, n$. When $z = x$, (22) reduces to the LLS estimator $\tilde{\delta}(x)$ in (20).

CFY (2000) consider a local linear approximation $\delta(z_t) \simeq \delta(z) + D(z)(z_t - z)'$. The LL varying coefficient estimator of CFY is then obtained by minimizing

$$\begin{aligned} \sum_{t=1}^n [y_t - X_t \delta(z_t)]^2 K_{tz} &= \sum_{t=1}^n [y_t - X_t \delta(z) - [(z_t - z) \otimes X_t] \text{vec} D(z)]^2 K_{tz} \\ &= \sum_{t=1}^n [y_t - X_t^z \delta^z(z)]^2 K_{tz} \end{aligned}$$

with respect to $\delta^z(z) = [\delta(z)' \text{vec} D(z)]'$ where $X_t^z = [X_t \ (z_t - z) \otimes X_t]$. This gives

$$\ddot{\delta}^z(z) = (\mathbf{X}^{z'}\mathbf{K}(z)\mathbf{X}^z)^{-1}\mathbf{X}^{z'}\mathbf{K}(z)\mathbf{y}, \quad (23)$$

and $\ddot{\delta}(z) = (\mathbf{I} \ 0)\ddot{\delta}^z(z)$. Hence

$$\ddot{m}(\xi) = (1 \ x \ 0)\ddot{\delta}^z(z) = (1 \ x)\ddot{\delta}(z).$$

For the asymptotic properties of these varying coefficient estimators, see CFY (2000). When $z = x$, (23) reduces to the LL varying coefficient estimator $\dot{\delta}^x(x)$ in (21). See Lee and Ullah (2001) for more discussion of these models and issues of testing nonlinearity.

3.11 Regime switching autoregressive model between unit root and stationary root

To avoid the usual dichotomy between unit-root nonstationarity and stationarity, we may consider models that permit two regimes of unit root nonstationarity and stationarity.

One model is the Innovation Regime-Switching (IRS) model of Kuan, Huang, and Tsay (2005). Intuitively, it may be implausible to believe that all random shocks exert only one effect (permanent or transitory) on future financial asset prices in a long time span. This intuition underpins the models that allow for breaks, stochastic unit root, or regime switching. As an alternative, Kuan, Huang, and Tsay (2005) propose the IRS model that permits the random shock in each period to be permanent or transitory, depending on a switching mechanism, and hence admits distinct dynamics (unit-root nonstationarity or stationarity) in different periods. Under the IRS framework, standard unit-root models and stationarity models are just two extreme cases. By applying the IRS model to real exchange rate, they circumvent the difficulties arising from unit-root (or stationarity) testing. They allow the data to speak for themselves, rather than putting them in the straitjacket of unit-root nonstationarity or stationarity. Huang and Kuan (2007) re-examine long-run PPP based on the IRS model and their empirical study on U.S./U.K. real exchange rates shows that there are both temporary and permanent influences on the real exchange rate such that approximately 42% of the shocks in the long run are more likely to have a permanent effect. They also found that transitory shocks dominate in the fixed-rate regimes, yet permanent shocks play a more important role during the floating regimes. Thus, the long-run PPP is rejected due to the presence of a significant amount of permanent shocks, but there are still long periods of time in which the deviations from long-run PPP are only transitory.

Another model is a threshold unit root (TUR) model or threshold integrated moving average (TIMA) model of Gonzalo and Martínez (2006). Based on this model they examine whether large and small shocks have different long-run effects, as well as whether one of them is purely transitory. They develop a new nonlinear permanent–transitory decomposition, that is applied to US stock prices to analyze the quality of the stock market.

Comparison of these two models with the linear autoregressive model with a unit root or

a stationary AR model for the out-of-sample forecasting remains to be examined empirically.

3.12 Bagging nonlinear forecasts

To improve on unstable forecast, bootstrap aggregating or bagging is introduced by Breiman (1996a). Lee and Yang (2006) show how bagging works for binary and quantile predictions. Lee and Yang (2006) attributed a part of success of the bagging predictors to the small sample estimation uncertainties. Therefore, a question that may arise is that whether the good performance of bagging predictors critically depends on algorithms we employ in nonlinear estimation.

They find that bagging improves the forecasting performance of predictors on highly nonlinear regression models – e.g., artificial neural network models, especially when the sample size is limited. It is usually hard to choose the number of hidden nodes and the number of inputs (or lags), and to estimate the large number of parameters in an ANN model. Therefore, a neural network model generate poor predictions in a small sample. In such cases, the bagging can do a valuable job to improve the forecasting performance as shown in Lee and Yang (2006), confirming the result of Breiman (1996b). Bagging predictor is a combined predictor formed over a set of training sets to smooth out the “instability” caused by parameter estimation uncertainty and model uncertainty. A predictor is said to be “unstable” if a small change in the training set will lead to a significant change in the predictor (Breiman, 1996b).

As bagging would be valuable in nonlinear forecasting, in this section, we will show how bagging predictor may improve the predicting performance of its underlying predictor. Let

$$\mathcal{D}_t \equiv \{(Y_s, \mathbf{X}_{s-1})\}_{s=t-R+1}^t \quad (t = R, \dots, T)$$

be a training set at time t and let $\varphi(\mathbf{X}_t, \mathcal{D}_t)$ be a forecast of Y_{t+1} or of the binary variable $G_{t+1} \equiv \mathbf{1}(Y_{t+1} \geq 0)$ using this training set \mathcal{D}_t and the explanatory variable vector \mathbf{X}_t . The optimal forecast $\varphi(\mathbf{X}_t, \mathcal{D}_t)$ for Y_{t+1} will be the conditional mean of Y_{t+1} given \mathbf{X}_t if we have the squared error loss function, or the conditional quantile of Y_{t+1} on \mathbf{X}_t if the loss is a tick function. Below we also consider the binary forecast for $G_{t+1} \equiv \mathbf{1}(Y_{t+1} \geq 0)$.

Suppose each training set \mathcal{D}_t consists of R observations generated from the underlying probability distribution \mathbf{P} . The forecast $\{\varphi(\mathbf{X}_t, \mathcal{D}_t)\}_{t=R}^T$ can be improved if more training sets were able to be generated from \mathbf{P} and the forecast can be formed from averaging the multiple forecasts obtained from the multiple training sets. Ideally, if \mathbf{P} were known and multiple training sets $\mathcal{D}_t^{(j)}$ ($j = 1, \dots, J$) may be drawn from \mathbf{P} , an ensemble aggregating predictor $\varphi_A(\mathbf{X}_t)$ can be constructed by the weighted averaging of $\varphi(\mathbf{X}_t, \mathcal{D}_t^{(j)})$ over j , i.e.,

$$\varphi_A(\mathbf{X}_t) \equiv \mathbb{E}_{\mathcal{D}_t} \varphi(\mathbf{X}_t, \mathcal{D}_t) \equiv \sum_{j=1}^J w_{j,t} \varphi(\mathbf{X}_t, \mathcal{D}_t^{(j)}),$$

where $\mathbb{E}_{\mathcal{D}_t}(\cdot)$ denotes the expectation over \mathbf{P} , $w_{j,t}$ is the weight function with $\sum_{j=1}^J w_{j,t} = 1$, and the subscript A in φ_A denotes ‘‘aggregation’’.

Lee and Yang (2006) show that the ensemble aggregating predictor $\varphi_A(X_t)$ has no larger expected loss than the original predictor $\varphi(X_t, \mathcal{D}_t)$. For any convex loss function $c(\cdot)$ on the forecast error z_{t+1} , we will have

$$\mathbb{E}_{\mathcal{D}_t, Y_{t+1}, \mathbf{X}_t} c(z_{t+1}) \geq \mathbb{E}_{Y_{t+1}, \mathbf{X}_t} c(\mathbb{E}_{\mathcal{D}_t}(z_{t+1})),$$

where $\mathbb{E}_{\mathcal{D}_t}(z_{t+1})$ is the aggregating forecast error, and $\mathbb{E}_{\mathcal{D}_t, Y_{t+1}, \mathbf{X}_t}(\cdot) \equiv \mathbb{E}_{\mathbf{X}_t}[\mathbb{E}_{Y_{t+1}, \mathbf{X}_t}\{\mathbb{E}_{\mathcal{D}_t}(\cdot) | X_t\}]$ denotes the expectation $\mathbb{E}_{\mathcal{D}_t}(\cdot)$ taken over \mathbf{P} (i.e., averaging over the multiple training sets generated from \mathbf{P}), then taking an expectation of Y_{t+1} conditioning on X_t , and then taking an expectation of X_t . Similarly we define the notation $\mathbb{E}_{Y_{t+1}, \mathbf{X}_t}(\cdot) \equiv \mathbb{E}_{\mathbf{X}_t}[\mathbb{E}_{Y_{t+1}, \mathbf{X}_t}(\cdot) | X_t]$. Therefore, the aggregating predictor will always have no larger expected cost than the original predictor for a convex loss function $\varphi(X_t, \mathcal{D}_t)$. The examples of the convex loss function includes the squared error loss and a tick (or check) loss $\rho_\alpha(z) \equiv [\alpha - \mathbf{1}(z < 0)]z$.

How much this aggregating predictor can improve depends on the distance between $\mathbb{E}_{\mathcal{D}_t, Y_{t+1}, \mathbf{X}_t} c(z_{t+1})$ and $\mathbb{E}_{Y_{t+1}, \mathbf{X}_t} c(\mathbb{E}_{\mathcal{D}_t}(z_{t+1}))$. We can define this distance by $\Delta \equiv \mathbb{E}_{\mathcal{D}_t, Y_{t+1}, \mathbf{X}_t} c(z_{t+1}) - \mathbb{E}_{Y_{t+1}, \mathbf{X}_t} c(\mathbb{E}_{\mathcal{D}_t}(z_{t+1}))$. Therefore, the effectiveness of the aggregating predictor depends on the *convexity* of cost function. The more convex is the cost function, the more effective this aggregating predictor can be. If the loss function is the squared error loss, then it can be shown that $\Delta = \mathbb{V}_{\mathcal{D}_t}[\varphi(\mathbf{X}_t, \mathcal{D}_t)]$ is the variance of the predictor, which measures the ‘‘instability’’ of the predictor. See Lee and Yang (2006, Proposition 1) and Breiman (1996b). If the loss is the

tick function, the effectiveness of bagging is also different for different quantile predictions: bagging works better for tail-quantile predictions than for mid-quantile predictions.

In practice, however, \mathbf{P} is not known. In that case we may estimate \mathbf{P} by its empirical distribution, $\hat{\mathbf{P}}(\mathcal{D}_t)$, for a given \mathcal{D}_t . Then, from the empirical distribution $\hat{\mathbf{P}}(\mathcal{D}_t)$, multiple training sets may be drawn by the bootstrap method. Bagging predictors, $\varphi^B(\mathbf{X}_t, \mathcal{D}_t^*)$, can then be computed by taking weighted average of the predictors trained over a set of bootstrap training sets. More specifically, the bagging predictor $\varphi^B(\mathbf{X}_t, \mathcal{D}_t^*)$ can be obtained in the following steps:

1. Given a training set of data at time t , $\mathcal{D}_t \equiv \{(Y_s, \mathbf{X}_{s-1})\}_{s=t-R+1}^t$, construct the j th bootstrap sample $\mathcal{D}_t^{*(j)} \equiv \{(Y_s^{*(j)}, \mathbf{X}_{s-1}^{*(j)})\}_{s=t-R+1}^t$, $j = 1, \dots, J$, according to the empirical distribution of $\hat{\mathbf{P}}(\mathcal{D}_t)$ of \mathcal{D}_t .
2. Train the model (estimate parameters) from the j th bootstrapped sample $\mathcal{D}_t^{*(j)}$.
3. Compute the bootstrap predictor $\varphi^{*(j)}(\mathbf{X}_t, \mathcal{D}_t^{*(j)})$ from the j th bootstrapped sample $\mathcal{D}_t^{*(j)}$.
4. Finally, for mean and quantile forecast, the bagging predictor $\varphi^B(\mathbf{X}_t, \mathcal{D}_t^*)$ can be constructed by averaging over J bootstrap predictors

$$\varphi^B(\mathbf{X}_t, \mathcal{D}_t^*) \equiv \sum_{j=1}^J \hat{w}_{j,t} \varphi^{*(j)}(\mathbf{X}_t, \mathcal{D}_t^{*(j)});$$

and for binary forecast, the bagging binary predictor $\varphi^B(\mathbf{X}_t, \mathcal{D}_t^*)$ can be constructed by majority voting over J bootstrap predictors:

$$\varphi^B(\mathbf{X}_t, \mathcal{D}_t^*) \equiv \mathbf{1} \left(\sum_{j=1}^J \hat{w}_{j,t} \varphi^{*(j)}(\mathbf{X}_t, \mathcal{D}_t^{*(j)}) > 1/2 \right)$$

with $\sum_{j=1}^J \hat{w}_{j,t} = 1$ in both cases.

One concern of applying bagging to time series is whether a bootstrap can provide a sound simulation sample for dependent data, for which the bootstrap is required to be consistent. It has been shown that some bootstrap procedure (such as moving block bootstrap) can provide consistent densities for moment estimators and quantile estimators. See, e.g., Fitzenberger (1997).

4 Nonlinear forecasting models for the conditional variance of returns

4.1 Nonlinear parametric models for volatility

Volatility models are of paramount importance in financial economics. Issues such as portfolio allocation, option pricing, risk management, and generally any decision making under uncertainty rely on the understanding and forecasting of volatility. This is one of the most active areas of research in time series econometrics. Important surveys as in Bollerslev, Chou, and Kroner (1992), Bera and Higgins (1993), Bollerslev, Engle, and Nelson (1994), Poon and Granger (2002), and Bauwens, Laurent, and Rombouts (2006) attest to the variety of issues in volatility research. The motivation for the introduction of the first generation of volatility models namely the ARCH models of Engle (1982) was to account for clusters of activity and fat-tail behavior of financial data. Subsequent models accounted for more complex issues. Among others and without being exclusive, we should mention issues related to asymmetric responses of volatility to news, probability distribution of the standardized innovations, i.i.d. behavior of the standardized innovation, persistence of the volatility process, linkages with continuous time models, intraday data and unevenly spaced observations, seasonality and noise in intraday data. The consequence of this research agenda has been a vast array of specifications for the volatility process.

Suppose that the return series $\{y_t\}_{t=1}^{T+1}$ of a financial asset follows the stochastic process $y_{t+1} = \mu_{t+1} + \varepsilon_{t+1}$, where $\mathbb{E}(y_{t+1}|\mathcal{F}_t) = \mu_{t+1}(\theta)$ and $\mathbb{E}(\varepsilon_{t+1}^2|\mathcal{F}_t) = \sigma_{t+1}^2(\theta)$ given the information set \mathcal{F}_t (σ -field) at time t . Let $z_{t+1} \equiv \varepsilon_{t+1}/\sigma_{t+1}$ have the conditional normal distribution with zero conditional mean and unit conditional variance. Volatility models can be classified in three categories: MA family, ARCH family, and stochastic volatility (SV) family.

The simplest method to forecast volatility is to calculate a historical moving average variance, denoted as MA(m), or an exponential weighted moving average (EWMA):

MA(m)	$\sigma_t^2 = \frac{1}{m} \sum_{j=1}^m (y_{t-j} - \hat{\mu}_t^m)^2$, $\hat{\mu}_t^m = \frac{1}{m} \sum_{j=1}^m y_{t-j}$
EWMA	$\sigma_t^2 = (1 - \lambda) \sum_{j=1}^{t-1} \lambda^{j-1} (y_{t-j} - \hat{\mu}_t)^2$, $\hat{\mu}_t = \frac{1}{t-1} \sum_{j=1}^{t-1} y_{t-j}$

Riskmetrics (1995) uses the EWMA specification with $\lambda = 0.94$ for predicting volatility. For these two MA family models, there are not parameters to estimate.

Second, the ARCH family is very extensive with many variations on the original model ARCH(p) of Engle (1982). Some representative models are: GARCH model of Bollerslev (1986); Threshold GARCH (T-GARCH) of Glosten *et al* (1993); Exponential GARCH (E-GARCH) of Nelson (1991); quadratic GARCH models (Q-GARCH) as in Sentana (1995); Absolute GARCH (ABS-GARCH) of Taylor (1986) and Schwert (1990); and Smooth Transition GARCH (ST-GARCH) of González-Rivera (1998).

ARCH(p)	$\sigma_t^2 = \omega + \sum_{i=1}^p \alpha_i \varepsilon_{t-i}^2$
GARCH	$\sigma_t^2 = \omega + \beta \sigma_{t-1}^2 + \alpha \varepsilon_{t-1}^2$
I-GARCH	$\sigma_t^2 = \omega + \beta \sigma_{t-1}^2 + \alpha \varepsilon_{t-1}^2, \alpha + \beta = 1$
T-GARCH	$\sigma_t^2 = \omega + \beta \sigma_{t-1}^2 + \alpha \varepsilon_{t-1}^2 + \gamma \varepsilon_{t-1}^2 \mathbf{1}(\varepsilon_{t-1} \geq 0)$
ST-GARCH	$\sigma_t^2 = \omega + \beta \sigma_{t-1}^2 + \alpha \varepsilon_{t-1}^2 + \gamma \varepsilon_{t-1}^2 F(\varepsilon_{t-1}, \delta)$ with $F(\varepsilon_{t-1}, \delta) = [1 + \exp(\delta \varepsilon_{t-1})]^{-1} - 0.5$
E-GARCH	$\ln \sigma_t^2 = \omega + \beta \ln \sigma_{t-1}^2 + \alpha [z_{t-1} - c z_{t-1}]$
Q-GARCH	$\sigma_t^2 = \omega + \beta \sigma_{t-1}^2 + \alpha (\varepsilon_{t-1} + \gamma)^2$
ABS-GARCH	$\sigma_t = \omega + \beta \sigma_{t-1} + \alpha \varepsilon_{t-1} $

The EWMA specification can be viewed as an integrated GARCH model with $\omega = 0$, $\alpha = \lambda$, and $\beta = 1 - \lambda$. In the T-GARCH model, the parameter γ allows for possible asymmetric effects of positive and negative innovations. In Q-GARCH models, the parameter γ measures the extent of the asymmetry in the news impact curve. For the ST-GARCH model, the parameter γ measures the asymmetric effect of positive and negative shocks, and the parameter $\delta > 0$ measures the smoothness of the transition between regimes, with a higher value of δ making ST-GARCH closer to T-GARCH.

Third, the stationary SV model of Taylor (1986) with η_t is i.i.d. $N(0, \sigma_\eta^2)$ and ξ_t is i.i.d. $N(0, \pi^2/2)$ is a representative member of the SV family.

SV	$\sigma_t^2 = \exp(0.5h_t), \ln(y_t^2) = -1.27 + h_t + \xi_t, h_t = \gamma + \phi h_{t-1} + \eta_t.$
----	--

With so many models, the natural question becomes which one to choose. There is not a universal answer to this question. The best model depends upon the objectives of the user. Thus, given an objective function, we search for the model(s) with the best predictive ability controlling for possible biases due to “data snooping” (Lo and MacKinlay, 1999). To compare the relative performance of volatility models, it is customary to choose either a statistical loss function or an economic loss function.

The preferred statistical loss functions are based on moments of forecast errors (mean-error, mean-squared error, mean absolute error, etc.). The best model will minimize a function of the forecast errors. The volatility forecast is often compared to a measure of realized volatility. With financial data, the common practice has been to take squared returns as a measure of realized volatility. However, this practice is questionable. Andersen and Bollerslev (1998) argued that this measure is a noisy estimate, and proposed the use of the intra-day (at each five minutes interval) squared returns to calculate the daily realized volatility. This measure requires intra-day data, which is subject to the variation introduced by the bid-ask spread and the irregular spacing of the price quotes.

Some other authors have evaluated the performance of volatility models with criteria based on economic loss functions. For example, West, Edison, and Cho (1993) considered the problem of portfolio allocation based on models that maximize the utility function of the investor. Engle, Kane, and Noh (1997) and Noh, Engle, and Kane (1994) considered different volatility forecasts to maximize the trading profits in buying/selling options. Lopez (2001) considered probability scoring rules that were tailored to a forecast user's decision problem and confirmed that the choice of loss function directly affected the forecast evaluation of different models. Brooks and Persaud (2003) evaluated volatility forecasting in a financial risk management setting in terms of Value-at-Risk (VaR). The common feature to these branches of the volatility literature is that none of these has controlled for forecast dependence across models and the inherent biases due to data-snooping.

Controlling for model dependence (White, 2000), González-Rivera, Lee, and Mishra (2004) evaluate fifteen volatility models for the daily returns to the SP500 index according to their out-of-sample forecasting ability. The forecast evaluation is based on two economic loss functions: an option pricing formula and a utility function; and two statistical loss functions: a goodness-of-fit based on a Value-at-Risk (VaR) calculation, and the predictive likelihood function. For option pricing, volatility is the only component that is not observable and it needs to be estimated. The loss function assess the difference between the actual price of a call option and the estimated price, which is a function of the estimated volatility of the stock. The second economic loss function refers to the problem of wealth allocation. An investor wishes to maximize her utility allocating wealth between a risky asset and a risk-free

asset. The loss function assesses the performance of the volatility estimates according to the level of utility they generate. The statistical function based on the goodness-of-fit of a VaR calculation is important for risk management. The main objective of VaR is to calculate extreme losses within a given probability of occurrence, and the estimation of the volatility is central to the VaR measure. The preferred models depend very strongly upon the loss function chosen by the user. González-Rivera, Lee, and Mishra (2004) find that, for option pricing, simple models such as the exponential weighted moving average (EWMA) proposed by Riskmetrics (González-Rivera, Lee, and Yoldas, 2007) performed as well as any GARCH model. For an utility loss function, an asymmetric quadratic GARCH model is the most preferred. For VaR calculations, a stochastic volatility model dominates all other models. And, for a predictive likelihood function, modeling the conditional standard deviation instead of the variance results in a dominant model.

4.2 Nonparametric models for volatility

Ziegelmann (2002) considers the kernel smoothing techniques that free the traditional parametric volatility estimators from the constraints related to their specific models. Ziegelmann (2002) applies the nonparametric local ‘exponential’ estimator to estimate conditional volatility functions, ensuring its nonnegativity. Its asymptotic properties are established and compared with those for the local linear estimator for the volatility model of Fan and Yao (1998). Long, Su, and Ullah (2007) extend this idea to semiparametric multivariate GARCH and show that there may exist substantial out-of-sample forecasting gain over the parametric models. This gain accounts for the presence of nonlinearity in the conditional variance-covariance that is neglected in parametric linear models.

4.3 Forecasting volatility using high frequency data – nonlinear HAR models

Using high-frequency data, quadratic variation may be estimated using realized volatility (RV). Andersen, Bollerslev, Diebold, and Labys (2001) and Barndorff-Nielsen and Shephard (2002) establish that RV, defined as the sum of squared intraday returns of small intervals,

is an asymptotically unbiased estimator of the unobserved quadratic variation as the interval length approaches zero. Besides the use of high frequency information in volatility estimation, volatility forecasting using high frequency information has been addressed as well. In an application to volatility prediction, Ghysels, Santa-Clara, and Valkanov (2006) investigate the predictive power of various regressors (lagged realized volatility, squared return, realized power, and daily range) for future volatility forecasting. They find that the best predictor is realized power (sum of absolute intraday returns), and more interestingly, direct use of intraday squared returns in mixed data sampling (MIDAS) regressions does not necessarily lead to better volatility forecasts.

Andersen, Bollerslev, Diebold, and Labys (2003) represent another approach to forecasting volatility using RV. The model they propose is a fractional integrated AR model: ARFI(5, d) for logarithmic RV's obtained from foreign exchange rates data of 30-minute frequency and demonstrate the superior predictive power of their model.

Alternatively, Corsi (2004) proposes the heterogeneous autoregressive (HAR) model of RV, which is able to reproduce long memory. McAleer and Medeiros (2007) propose a new model that is a multiple regime smooth transition (ST) extension of the HAR model, which is specifically designed to model the behavior of the volatility inherent in financial time series. The model is able to describe simultaneously long memory as well as sign and size asymmetries. They apply the model to several Dow Jones Industrial Average index stocks using transaction level data from the Trades and Quotes database that covers ten years of data, and find strong support for long memory and both sign and size asymmetries. Furthermore, they show that the multiple regime smooth transition HAR model, when combined with the linear HAR model, is flexible for purposes of forecasting volatility.

5 Forecasting quantiles, directions, durations, density, intervals, and etc.

In statistical literature the major development in nonlinear time series models has been in the conditional mean, it is not clear yet that any of those nonlinear models may generate profits after accounting for various market frictions and transactions cost. In financial econometrics,

therefore the main objective has been forecasting higher moments, quantiles, durations, directions, and some other aspects. In the previous two sections, we have considered the conditional mean and variance in some detail. In this section, we provide a brief survey on the other aspects.

5.1 Forecasting quantiles

Optimal forecast of a time series model extensively depends on the specification of the loss function. Symmetric quadratic loss function is the most prevalent in applications due to its simplicity. The optimal forecast under quadratic loss is simply the conditional mean, but an asymmetric loss function implies a more complicated forecast that depends on the distribution of the forecast error as well as the loss function itself (Granger 1999), as the expected loss function if formulated with the expectation taken with respect to the conditional distribution. Specification of the loss function defines the model under consideration.

Consider a stochastic process $Z_t \equiv (Y_t, X_t)'$ where Y_t is the variable of interest and X_t is a vector of other variables. Suppose there are $T + 1$ ($\equiv R + P$) observations. We use the observations available at time t , $R \leq t < T + 1$, to generate P forecasts using each model. For each time t in the prediction period, we use either a rolling sample $\{Z_{t-R+1}, \dots, Z_t\}$ of size R or the whole past sample $\{Z_1, \dots, Z_t\}$ to estimate model parameters $\hat{\beta}_t$. We can then generate a sequence of one-step-ahead forecasts $\{f(Z_t, \hat{\beta}_t)\}_{t=R}^T$.

Suppose that there is a decision maker who takes an one-step point forecast $f_{t,1} \equiv f(Z_t, \hat{\beta}_t)$ of Y_{t+1} and uses it in some relevant decision. The one-step forecast error $e_{t+1} \equiv Y_{t+1} - f_{t,1}$ will result in a cost of $c(e_{t+1})$, where the function $c(e)$ will increase as e increases in size, but not necessarily symmetrically or continuously. The optimal forecast $f_{t,1}^*$ will be chosen to produce the forecast errors that minimize the expected loss

$$\min_{f_{t,1}} \int_{-\infty}^{\infty} c(y - f_{t,1}) dF_t(y),$$

where $F_t(y) \equiv \Pr(Y_{t+1} \leq y | I_t)$ is the conditional distribution function, with I_t being some proper information set at time t that includes Z_{t-j} , $j \geq 0$. The corresponding optimal forecast error will be

$$e_{t+1}^* = Y_{t+1} - f_{t,1}^*.$$

Then the optimal forecast would satisfy

$$\frac{\partial}{\partial f_{t,1}} \int_{-\infty}^{\infty} c(y - f_{t,1}^*) dF_t(y) = 0.$$

When we may interchange the operations of differentiation and integration,

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial f_{t,1}} c(y - f_{t,1}^*) dF_t(y) \equiv \mathbb{E} \left(\frac{\partial}{\partial f_{t,1}} c(Y_{t+1} - f_{t,1}^*) | I_t \right)$$

the “generalized forecast error”, $g_{t+1} \equiv \frac{\partial}{\partial f_{t,1}} c(Y_{t+1} - f_{t,1}^*)$, forms the condition of forecast optimality:

$$H_0 : \mathbb{E}(g_{t+1} | I_t) = 0 \quad a.s.,$$

that is a martingale difference (MD) property of the generalized forecast error. This forms the optimality condition of the forecasts and gives an appropriate regression function corresponding to the specified loss function $c(\cdot)$.

To see this we consider the following two examples. First, when the loss function is the squared error loss

$$c(Y_{t+1} - f_{t,1}) = (Y_{t+1} - f_{t,1})^2,$$

the generalized forecast error will be $g_{t+1} \equiv \frac{\partial}{\partial f_{t,1}} c(Y_{t+1} - f_{t,1}^*) = -2e_{t+1}^*$ and thus $\mathbb{E}(e_{t+1}^* | I_t) = 0$ *a.s.*, which implies that the optimal forecast

$$f_{t,1}^* = \mathbb{E}(Y_{t+1} | I_t)$$

is the conditional mean. Next, when the loss is the check function, $c(e) = [\alpha - \mathbf{1}(e < 0)] \cdot e \equiv \rho_\alpha(e_{t+1})$, the optimal forecast $f_{t,1}$, for given $\alpha \in (0, 1)$, minimizing

$$\min_{f_{t,1}} \mathbb{E}[c(Y_{t+1} - f_{t,1}) | I_t]$$

can be shown to satisfy

$$\mathbb{E}[\alpha - \mathbf{1}(Y_{t+1} < f_{t,1}^*) | I_t] = 0 \quad a.s.$$

Hence, $g_{t+1} \equiv \alpha - \mathbf{1}(Y_{t+1} < f_{t,1}^*)$ is the generalized forecast error. Therefore,

$$\alpha = \mathbb{E}[\mathbf{1}(Y_{t+1} < f_{t,1}^*) | I_t] = \Pr(Y_{t+1} \leq f_{t,1}^* | I_t),$$

and the optimal forecast $f_{t,1}^* = q_\alpha(Y_{t+1}|I_t)$ is the conditional α -quantile.

Increasing financial fragility in emerging markets and the extensive use of derivative products in developed countries can be characterized as two distinct features of financial world over the last decade. Consequently, effective use of risk measurement tools has been suggested as a main panacea for mitigating growing financial risk. Uniform risk measurement methodology called Value-at-Risk (VaR) has received a great attention from both regulatory and academic fronts. During a short span of time, numerous papers have studied various aspects of VaR methodology. Bao, Lee, and Saltoglu (2006) examine the relative out-of-sample predictive performance of various VaR models in the literature.

An interesting VaR model is the CaViaR (conditional autoregressive Value-at-Risk) model suggested by Engle and Manganelli (1999). They estimate the VaR from a quantile regression rather than inverting a conditional distribution. The idea is similar to the GARCH modeling and VaR is modeled autoregressively

$$q_t(\alpha) = a_0 + a_1 q_{t-1}(\alpha) + h(x_t|\theta),$$

where $x_t \in \mathcal{F}_{t-1}$, θ is a parameter vector, and $h(\cdot)$ is a function to explain the VaR model. Depending on the specification of $h(\cdot)$, the CaViaR model may be

$$q_t(\alpha) = a_0 + a_1 q_{t-1}(\alpha) + a_2 |r_{t-1}|,$$

$$q_t(\alpha) = a_0 + a_1 q_{t-1}(\alpha) + a_2 |r_{t-1}| + a_3 |r_{t-1}| \cdot \mathbf{1}(r_{t-1} < 0),$$

where the second model allow nonlinearity (asymmetry) similar to the asymmetric GARCH model of Glosten et al (1993).

Bao, Lee, and Saltoglu (2006) compare various VaR models. Their results show that the CaViaR quantile regression models of Engle and Manganelli (2004) have shown some success in predicting the VaR risk measure for various periods, generally more stable than those that invert a distribution function.

5.2 Forecasting directions

It is well known that, while financial returns $\{Y_t\}$ may not be predictable, their variance, sign, and quantiles may be predictable. Christofferson and Diebold (2006) show that binary

variable $G_{t+1} \equiv \mathbf{1}(Y_{t+1} > 0)$ is predictable when some conditional moments are time varying, where $\mathbf{1}(\cdot)$ takes the value of 1 if the statement in the parenthesis is true, and 0 otherwise. Hong and Lee (2003), Hong and Chung (2003), Linton and Whang (2004), Lee and Yang (2006) among many others find some evidence that the directions of stock returns and foreign exchange rate changes are predictable. Lee and Yang (2006) also show that forecasting quantiles and forecasting binary (directional) forecasts are related, in that the former may lead to the latter. Therefore, predictability of financial returns in conditional quantiles may imply the predictability of financial returns in conditional directions.

5.3 Probability forecasts

Diebold and Rudebush (1989) consider the probability forecasts for business cycle turning points. To measure the accuracy of predicted probabilities, that is the average distance between the predicted probabilities and observed realization (as measured by a zero-one dummy variable). Suppose there are $T + 1$ ($\equiv R + P$) observations. We use the observations available at time t ($R \leq t < T + 1$), to estimate a model. We then have time series of $P = T - R + 1$ probability forecasts $\{p_{t+1}\}_{t=R}^T$ where p_t is the predicted probability of the occurrence of an event (e.g., business cycle turning point) in the next period $t + 1$. Let $\{d_{t+1}\}_{t=R}^T$ be the corresponding realization with $d_t = 1$ if a business cycle turning point (or any defined event) occurs in period t and $d_t = 0$ otherwise. The loss function analogous to the squared error is Brier's score based on quadratic probability score (QPS):

$$QPS = P^{-1} \sum_{t=R}^T 2(p_t - d_t)^2.$$

The QPS ranges from 0 to 2, with 0 for perfect accuracy. As noted by Diebold and Rudebush (1989), the use of the symmetric loss function may not be appropriate as a forecaster may be penalized more heavily for missing a call (making a type II error) than for signaling a false alarm (making a type I error). Another loss function is given by the log probability score (LPS)

$$LPS = -P^{-1} \sum_{t=R}^T \ln (p_t^{d_t} (1 - p_t)^{(1-d_t)}),$$

which is similar to the loss for the interval forecast. A large mistakes are penalized more heavily under LPS than under QPS. More loss functions are discussed in Diebold and Rudebush (1989).

Another loss function useful in this context is the Kuipers score (KS), which is defined by

$$KS = \text{Hit Rate} - \text{False Alarm Rate},$$

where Hit Rate is the fraction of the bad events that were correctly predicted as good events (power, or $1 -$ probability of type II error), and False Alarm Rate is the fraction of good events that had been incorrectly predicted as bad events (probability of type I error).

5.4 Forecasting interval

Suppose Y_t is a stationary series. Let the one-period ahead conditional interval forecast made at time t from a model be denoted as

$$J_{t,1}(\alpha) = (L_{t,1}(\alpha), U_{t,1}(\alpha)), \quad t = R, \dots, T,$$

where $L_{t,1}(\alpha)$ and $U_{t,1}(\alpha)$ are the lower and upper limits of the ex ante interval forecast for time $t + 1$ made at time t with the coverage probability α . Define the indicator variable $X_{t+1}(\alpha) = \mathbf{1}[Y_{t+1} \in J_{t,1}(\alpha)]$. The sequence $\{X_{t+1}(\alpha)\}_{t=R}^T$ is i.i.d. Bernoulli (α) . The optimal interval forecast would satisfy $\mathbb{E}(X_{t+1}(\alpha)|I_t) = \alpha$, so that $\{X_{t+1}(\alpha) - \alpha\}$ will be an MD. A better model has a larger expected Bernoulli log-likelihood

$$\mathbb{E}\alpha^{X_{t+1}(\alpha)}(1 - \alpha)^{[1-X_{t+1}(\alpha)]}.$$

Hence, we can choose a model for interval forecasts with the largest out-of-sample mean of the predictive log-likelihood defined by

$$P^{-1} \sum_{t=R}^T \ln (\alpha^{x_{t+1}(\alpha)}(1 - \alpha)^{[1-x_{t+1}(\alpha)]}).$$

5.5 Forecasting density

Consider a financial return series $\{y_t\}_{t=1}^T$. This observed data on a univariate series is a realization of a stochastic process $\mathbf{Y}^T \equiv \{Y_\tau : \Omega \rightarrow \mathbb{R}, \tau = 1, 2, \dots, T\}$ on a complete

probability space $(\Omega, \mathcal{F}_T, P_0^T)$, where $\Omega = \mathbb{R}^T \equiv \times_{\tau=1}^T \mathbb{R}$ and $\mathcal{F}_T = \mathcal{B}(\mathbb{R}^T)$ is the Borel σ -field generated by the open sets of \mathbb{R}^T , and the *joint* probability measure $P_0^T(B) \equiv P_0[\mathbf{Y}^T \in B]$, $B \in \mathcal{B}(\mathbb{R}^T)$ completely describes the stochastic process. A sample of size T is denoted as $\mathbf{y}^T \equiv (y_1, \dots, y_T)'$.

Let σ -finite measure ν^T on $\mathcal{B}(\mathbb{R}^T)$ be given. Assume $P_0^T(B)$ is absolutely continuous with respect to ν^T for all $T = 1, 2, \dots$, so that there exists a measurable Radon-Nikodým density $g^T(\mathbf{y}^T) = dP_0^T/d\nu^T$, unique up to a set of zero measure- ν^T .

Following White (1994), we define a probability model \mathcal{P} as a collection of distinct probability measures on the measurable space (Ω, \mathcal{F}_T) . A probability model \mathcal{P} is said to be correctly specified for \mathbf{Y}^T if \mathcal{P} contains P_0^T . Our goal is to evaluate and compare a set of parametric probability models $\{P_\theta^T\}$, where $P_\theta^T(B) \equiv P_\theta[Y^T \in B]$. Suppose there exists a measurable Radon-Nikodým density $f^T(\mathbf{y}^T) = dP_\theta^T/d\nu^T$ for each $\theta \in \Theta$, where θ is a finite-dimensional vector of parameters and is assumed to be identified on Θ , a compact subset of \mathbb{R}^k . See White (1994, Theorem 2.6).

In the context of forecasting, instead of the joint density $g^T(\mathbf{y}^T)$, we consider forecasting the *conditional* density of \mathbf{Y}^t , given the information \mathcal{F}_{t-1} generated by \mathbf{Y}^{t-1} . Let $\varphi_t(y_t) \equiv \varphi_t(y_t|\mathcal{F}_{t-1}) \equiv g^t(\mathbf{y}^t)/g^{t-1}(\mathbf{y}^{t-1})$ for $t = 2, 3, \dots$ and $\varphi_1(y_1) \equiv \varphi_1(y_1|\mathcal{F}_0) \equiv g^1(\mathbf{y}^1) = g^1(y_1)$. Thus the goal is to forecast the (true, unknown) conditional density $\varphi_t(y_t)$.

For this, we use an one-step-ahead conditional density forecast model $\psi_t(y_t; \theta) \equiv \psi_t(y_t|\mathcal{F}_{t-1}; \theta) \equiv f^t(\mathbf{y}^t)/f^{t-1}(\mathbf{y}^{t-1})$ for $t = 2, 3, \dots$ and $\psi_1(y_1) \equiv \psi_1(y_1|\mathcal{F}_0) \equiv f^1(\mathbf{y}^1) = f^1(y_1)$. If $\psi_t(y_t; \theta_0) = \varphi_t(y_t)$ almost surely for some $\theta_0 \in \Theta$, then the one-step-ahead density forecast is correctly specified, and it is said to be optimal because it dominates all other density forecasts for any loss functions as discussed in the previous section (see Granger and Pesaran, 2000; Diebold *et al.*, 1998; Granger 1999).

In practice, it is rarely the case that we can find an optimal model. As it is very likely that “the true distribution is in fact too complicated to be represented by a simple mathematical function” (Sawa, 1978), all the models proposed by different researchers can be possibly misspecified and thereby we regard each model as an approximation to the truth. Our task is then to investigate which density forecast model can approximate the true conditional density most closely. We have to first define a metric to measure the distance of a given

model to the truth, and then compare different models in terms of this distance.

The adequacy of a density forecast model can be measured by the conditional Kullback-Leibler (1951) Information Criterion (KLIC) divergence measure between two conditional densities,

$$\mathbb{I}_t(\varphi : \psi, \boldsymbol{\theta}) = \mathbb{E}_{\varphi_t}[\ln \varphi_t(y_t) - \ln \psi_t(y_t; \boldsymbol{\theta})],$$

where the expectation is with respect to the true conditional density $\varphi_t(\cdot|\mathcal{F}_{t-1})$, $\mathbb{E}_{\varphi_t} \ln \varphi_t(y_t|\mathcal{F}_{t-1}) < \infty$, and $\mathbb{E}_{\varphi_t} \ln \psi_t(y_t|\mathcal{F}_{t-1}; \boldsymbol{\theta}) < \infty$. Following White (1994), we define the distance between a density model and the true density as the minimum of the KLIC

$$\mathbb{I}_t(\varphi : \psi, \boldsymbol{\theta}_{t-1}^*) = \mathbb{E}_{\varphi_t}[\ln \varphi_t(y_t) - \ln \psi_t(y_t; \boldsymbol{\theta}_{t-1}^*)],$$

where $\boldsymbol{\theta}_{t-1}^* = \arg \min \mathbb{I}_t(\varphi : \psi, \boldsymbol{\theta})$ is the pseudo-true value of $\boldsymbol{\theta}$ (Sawa, 1978). We assume that $\boldsymbol{\theta}_{t-1}^*$ is an interior point of Θ . The smaller this distance is, the closer the density forecast $\psi_t(\cdot|\mathcal{F}_{t-1}; \boldsymbol{\theta}_{t-1}^*)$ is to the true density $\varphi_t(\cdot|\mathcal{F}_{t-1})$.

However, $\mathbb{I}_t(\varphi : \psi, \boldsymbol{\theta}_{t-1}^*)$ is unknown since $\boldsymbol{\theta}_{t-1}^*$ is not observable. We need to estimate $\boldsymbol{\theta}_{t-1}^*$. If our purpose is to compare the out-of-sample predictive abilities among competing density forecast models, we split the data into two parts, one for estimation and the other for out-of-sample validation. At each period t in the out-of-sample period ($t = R + 1, \dots, T$), we estimate the unknown parameter vector $\boldsymbol{\theta}_{t-1}^*$ and denote the estimate as $\hat{\boldsymbol{\theta}}_{t-1}$. Using $\{\hat{\boldsymbol{\theta}}_{t-1}\}_{t=R+1}^T$, we can obtain the out-of-sample estimate of $\mathbb{I}_t(\varphi : \psi, \boldsymbol{\theta}_{t-1}^*)$ by

$$\mathbb{I}_P(\varphi : \psi) \equiv \frac{1}{P} \sum_{t=R+1}^T \ln[\varphi_t(y_t)/\psi_t(y_t; \hat{\boldsymbol{\theta}}_{t-1})]$$

where $P = T - R$ is the size of the out-of-sample period. Note that

$$\mathbb{I}_P(\varphi : \psi) = \frac{1}{P} \sum_{t=R+1}^T \ln[\varphi_t(y_t)/\psi_t(y_t; \boldsymbol{\theta}_{t-1}^*)] + \frac{1}{P} \sum_{t=R+1}^T \ln[\psi_t(y_t; \boldsymbol{\theta}_{t-1}^*)/\psi_t(y_t; \hat{\boldsymbol{\theta}}_{t-1})],$$

where the first term in $\mathbb{I}_P(\varphi : \psi)$ measures model uncertainty (the distance between the optimal density $\varphi_t(y_t)$ and the model $\psi_t(y_t; \boldsymbol{\theta}_{t-1}^*)$) and the second term measures parameter estimation uncertainty due to the distance between $\boldsymbol{\theta}_{t-1}^*$ and $\hat{\boldsymbol{\theta}}_{t-1}$.

Since the KLIC measure takes on a smaller value when a model is closer to the truth, we can regard it as a loss function and use $\mathbb{I}_P(\varphi : \psi)$ to formulate the loss-differential. The out-of-sample average of the loss-differential between model 1 and model 2 is

$$\mathbb{I}_P(\varphi : \psi^1) - \mathbb{I}_P(\varphi : \psi^2) = \frac{1}{P} \sum_{t=R+1}^T \ln[\psi_t^2(y_t; \hat{\boldsymbol{\theta}}_{t-1}^2) / \psi_t^1(y_t; \hat{\boldsymbol{\theta}}_{t-1}^1)],$$

which is the ratio of the two predictive log-likelihood functions. With treating model 1 as a benchmark model (for model selection) or as the model under the null hypothesis (for hypothesis testing), $\mathbb{I}_P(\varphi : \psi^1) - \mathbb{I}_P(\varphi : \psi^2)$ can be considered as a loss function to minimize. To sum up, the KLIC differential can serve as a *loss* function for density forecast evaluation as discussed in Bao, Lee, and Saltoglu (2007).

Using the KLIC divergence measure to characterize the extent of misspecification of a forecast model, Bao, Lee, and Saltoglu (2007), in an empirical study with the S&P500 and NASDAQ daily return series, find strong evidence for rejecting the Normal-GARCH benchmark model, in favor of the models that can capture skewness in the conditional distribution and asymmetry and long-memory in the conditional variance. Also, Bao and Lee (2006) investigate the nonlinear predictability of stock returns when the density forecasts are evaluated/compared instead of the conditional mean point forecasts. The conditional mean models they use for the daily closing S&P500 index returns include the martingale difference model, the linear ARMA models, the STAR and SETAR models, the ANN model, and the polynomial model. Their empirical findings suggest the out-of-sample predictive abilities of nonlinear models for stock returns are asymmetric in the sense that the right tails of the return series are predictable via many of the nonlinear models while we find no such evidence for the left tails or the entire distribution.

6 Conclusions

This article is intended to provide an introduction to selective nonlinear time series models that are useful in forecasting financial returns, volatility, duration, directional changes, and etc. Given the space, the coverage is quite limited and the readers may find the following additional references to complement the present paper. Campbell, Lo, and MacKinlay (1997,

Ch 12) provides a brief but excellent summary of the important nonlinear time series models for the conditional mean and conditional variance as well and various methods such as ANN and nonparametric methods. Similarly, the interested readers may also refer to the recent books and monographs of Granger and Teräsvirta (1993), Fan and Yao (2003), Franses and van Dijk (2000), Gao (2007), Tsay (2005), and some book chapters such as Stock (2001), Tsay (2002), Teräsvirta (2006), and White (2006).

7 Future Directions

In the article we mainly focus on forecasting the conditional mean returns, and for space reason we have been brief on other aspects of financial forecasting. In fact, it is well known that the mean return is not easy to predict, and therefore in financial forecasting, it is perhaps more important to predict the conditional variance, conditional quantile, conditional direction, conditional duration, etc.

We have not discuss the important issue of statistical inference for space reason. One issue is inference about comparing the predictive ability of competing models. There are now large literature, that are to still grow much more. Examples are Granger and Newbold (1986), Diebold and Mariano (1995), West (1998), White (2000), Hansen (2005), Romano and Wolf (2005), Giacomini and White (2006), etc. Another issue is inference about optimality of forecasts. Empirical tests of forecast optimality have traditionally been conducted under the assumption of mean squared error loss or some other known loss function as discussed in Section 5.1. Patton and Timmermann (2007) examine new testable properties that hold when the forecaster's loss function is unknown but testable. The readers can refer to these references and references therein for more details. We leave these for other work.

8 Bibliography

- Ait-Sahalia, Y. and L.P. Hansen (eds.) (2007) *Handbook of Financial Econometrics*, forthcoming, North Holland, Amsterdam.
- Andersen, T.G, and T. Bollerslev (1998), "Answering the Skeptics: Yes, Standard Volatility Models Do Provide Accurate Forecasts", *International Economic Review* 39(4), 885-905.

- Andersen, T.G., Bollerslev, T., Diebold, F.X. and Labys, P. (2001), "The Distribution of Realized Exchange Rate Volatility," *Journal of the American Statistical Association* 96, 42-55.
- Andersen, T.G., Bollerslev, T., Diebold, F.X. and Labys, P. (2003), "Modeling and Forecasting Realized Volatility," *Econometrica* 71, 579-625.
- Ang, A. and G. Bekaert (2002), "Regime Switcheds in Interest Rates", *Journal of Business and Economic Statistics* 20, 163-182.
- Bachelier, L. (1900), "Theory of Speculation", in Cootner, P. (ed.), *The Random Character of Stock Market Prices*, MIT Press, Cambridge, MA, 1964: reprint.
- Bai, J. and S. Ng (2007), "Forecasting Economic Time Series Using Targeted Predictors," NYU and Columbia.
- Bao, Y. and T.-H. Lee (2006), "Asymmetric Predictive Abilities of Nonlinear Models for Stock Returns: Evidence from Density Forecast Comparison", *Advances in Econometrics*, Volume 20, Part B, pages 41-62.
- Bao, Y., T.-H. Lee, B. Saltoglu (2006), "Evaluating Predictive Performance of Value-at-Risk Models in Emerging Markets: A Reality Check", *Journal of Forecasting* 25(2), 101-128.
- Bao, Y., T.-H. Lee, B. Saltoglu (2007), "Comparing Density Forecast Models", *Journal of Forecasting* 26(3), 203-225.
- Barndorff-Nielsen, O.E. and Shephard, N. (2002), "Econometric Analysis of Realised Volatility and Its Use in Estimating Stochastic Volatility Models," *Journal of the Royal Statistical Society B* 64, 853-223.
- Bauwens, L., S. Laurent, and J.V.K. Rombouts (2006), "Multivariate GARCH Models: A Survey", *Journal of Applied Econometrics* 21, 79-109.
- Bera, A.K. and M.L. Higgins (1993), "ARCH Models: Properties, Estimation, and Testing", *Journal of Economic Surveys* 7, 305-366.
- Bollerslev, T. (1986) "Generalized Autoregressive Conditional Heteroskedasticity", *Journal of Econometrics* 31, 307-327.
- Bollerslev, T., R.Y. Chou, and K.F. Kroner (1992), "ARCH Models in Finance", *Journal of Econometrics* 52, 5-59.
- Bollerslev, T., R.F., Engle, and D.B. Nelson (1994), "ARCH Models", *Handbook of Econometrics*, Volume 4.
- Bollerslev, T., R.F., Engle, and J. Wooldridge (1988), "A Capital Asset Pricing Model with Time Varying Covariances", *Journal of Political Economy* 96, 116-131.
- Boero, G. and E. Marrocu (2004), "The Performance of SETAR Models: A Regime Conditional Evaluation of Point, Interval, and Density Forecasts", *International Journal of Forecasting* 20, 305-320.

- Breiman, L. (1996a), “Bagging Predictors”, *Machine Learning*, 24, 123-140.
- Breiman, L. (1996b), “Heuristics of Instability and Stabilization in Model Selection”, *Annals of Statistics*, 24(6), 2350–2383.
- Brooks, C. and G. Persaud (2003), “Volatility Forecasting for Risk Management”, *Journal of Forecasting*, 22(1), 1-22.
- Campbell, J.Y, A.W. Lo, and A.C. MacKinlay (1997), *The Econometrics of Financial Markets*, Princeton University Press, New Jersey.
- Campbell, J.Y. and S.B. Thompson (2007), “Predicting Excess Stock Returns Out of Sample: Can Anything Beat the Historical Average?”, Harvard Institute of Economic Research, Discussion Paper No. 2084.
- Cai, Z., J. Fan, and Q. Yao (2000), “Functional-coefficient Regression Models for Nonlinear Time Series”, *Journal of the American Statistical Association* 95, 941-956.
- Chen, X. (2006), “Large Sample Sieve Estimation of Semi-Nonparametric Models”, *Handbook of Econometrics*, Vol. 6, Chapter 76.
- Chen, R. and R.S. Tsay (1993), “Functional-coefficient Autoregressive Models”, *Journal of American Statistical Association* 88, 298-308.
- Christofferson, P.F. and F.X. Diebold (2006), Financial Asset Returns, Direction-of-Change Forecasting, and Volatility Dynamics, *Management Science*, 52, 1273-1287.
- Cleveland, W.S. (1979), “Robust Locally Weighted Regression and Smoothing Scatter Plots”, *Journal of American Statistical Association* 74, 829-836.
- Corsi, F. (2004), “A Simple Long Memory Model of Realized Volatility,” University of Lugano.
- Dahl, C.M. and G. González-Rivera (2003a), “Testing for Neglected Nonlinearity in Regression Models based on the Theory of Random Fields”, *Journal of Econometrics* 114, 141-164.
- Dahl, C.M. and G. González-Rivera (2003b), “Identifying Nonlinear Components by Random Fields in the US GNP Growth. Implications for the Shape of the Business Cycle”, *Studies in Nonlinear Dynamics and Econometrics* 7(1), art2.
- Diebold, F.X. and R. Mariano (1995), “Comparing predictive accuracy”, *Journal of Business and Economic Statistics* 13, 253-265.
- Diebold, F.X., T.A. Gunther and A.S. Tay (1998), “Evaluating Density Forecasts with Applications to Financial Risk Management”, *International Economic Review* 39, 863-883.
- Diebold, F.X. and G.D. Rudebusch (1989), “Scoring the Leading Indicators”, *Journal of Business*, 62(3) 369-391.

- Dueker, M. and C.J. Neely (2007), “Can Markov Switching Models Predict Excess Foreign Exchange Returns?”, *Journal of Banking and Finance* 31, 279-296.
- Durland, J.M. and T.H. McCurdy (1994), “Duration-Dependent Transitions in a Markov Model of US GNP Growth”, *Journal of Business and Economic Statistics* 12, 279-288.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004), “Least Angle Regression”, *Annals of Statistics* 32(2), 407–499.
- Engle, R.F. (1982), “Autoregressive Conditional Heteroscedasticity with Estimates of the Variance of UK Inflation”, *Econometrica*, 50, 987-1008.
- Engle, R.F., V.K. Ng, and M. Rothschild (1990), “Asset Pricing with a Factor ARCH Covariance Structure: Empirical Estimates for Treasury Bills”, *Journal of Econometrics* 45, 213-238.
- Engle, R.F. and J.R. Russell (1998), “Autoregressive Conditional Duration: A New Model for Irregularly Spaced Transaction Data,” *Econometrica* 66, 1127-1162.
- Fan, J. and I. Gijbels (1996), *Local Polynomial Modelling and Its Applications*, Chapman and Hall, London.
- Fan, J. and Q. Yao (1998), “Efficient estimation of conditional variance functions in stochastic regression”, *Biometrika* 85, 645-660.
- Fan, J. and Q. Yao (2003), *Nonlinear Time Series*, Springer.
- Fan, J., Yao, Q. and Cai, Z. (2003), “Adaptive varying-coefficient linear models”, *Journal of Royal Statistical Society B* 65, 57-80.
- Fitzenberger, B. (1997), “The Moving Blocks Bootstrap and Robust Inference for Linear Least Squares and Quantile Regressions”, *Journal of Econometrics* 82, 235-287.
- Franses, P.H. and van Dijk, D. (2000), *Nonlinear Time Series Models in Empirical Finance*, Cambridge University Press.
- Gao, J. (2007), *Nonlinear Time Series: Semiparametric and Nonparametric Methods*, Chapman and Hall/CRC.
- Ghysels, E., Santa-Clara, P., and Valkanov, R. (2006), “Predicting Volatility: How to Get Most out of Returns Data Sampled at Different Frequencies”, *Journal of Econometrics* 131, 59-95.
- Giacomini, R. and H. White (2006), “Tests of Conditional Predictive Ability,” *Econometrica* 74, 1545-1578.
- Glosten, L.R., R. Jaganathan, and D. Runkle (1993), “On the Relationship between the Expected Value and the Volatility of the Nominal Excess Return on Stocks,” *Journal of Finance* 48, 1779-1801.
- González-Rivera, G. (1998), “Smooth-Transition GARCH Models”, *Studies in Nonlinear Dynamics and Econometrics* 3(2), 61-78.

- González-Rivera, G., T.-H. Lee, and S. Mishra (2004), “Forecasting Volatility: A Reality Check Based on Option Pricing, Utility Function, Value-at-Risk, and Predictive Likelihood”. *International Journal of Forecasting* 20(4): 629-645.
- González-Rivera, G., T.-H. Lee, and S. Mishra (2008), “Jumps in Cross-Sectional Rank and Expected Returns: A Mixture Model”, *Journal of Applied Econometrics*, forthcoming.
- González-Rivera, G., T.-H. Lee, and E. Yoldas (2007), “Optimality of the Riskmetrics VaR Model”, *Finance Research Letters* 4, 137-145.
- Gonzalo, J. and O. Martíneza (2006), “Large shocks vs. small shocks. (Or does size matter? May be so.)”, *Journal of Econometrics* 135, 311-347.
- Goyal, A. and I. Welch (2006), “A Comprehensive Look at The Empirical Performance of Equity Premium Prediction”, Emory and Brown.
- Granger, C.W.J. (1999), “Outline of Forecast Theory Using Generalized Cost Functions”, *Spanish Economic Review* 1, 161-173.
- Granger, C.W.J. and P. Newbold (1986), *Forecasting Economic Time Series*, Academic Press, 2ed.
- Granger, C.W.J. and M.H. Pesaran (2000), “A Decision Theoretic Approach to Forecasting Evaluation”, in *Statistics and Finance: An Interface*, W. S. Chan, W. K. Li, and Howell Tong (eds.), London: Imperial College Press.
- Granger, C.W.J. and T.-H. Lee (1999), “The Effect of Aggregation on Nonlinearity”, *Econometric Reviews* 18(3), 259-269.
- Granger, C.W.J. and T. Teräsvirta (1993), *Modelling Nonlinear Economic Relationships*, Oxford University Press, New York.
- Haggan, V. and T. Ozaki (1981), “Modeling Nonlinear Vibrations Using an Amplitude-dependent Autoregressive Time Series Model”, *Biometrika*, 68, 189-196.
- Hamilton, J.D. (1994), *Time Series Analysis*, Princeton University Press, New Jersey.
- Hamilton, J.D. (1989), “A New Approach to the Economic Analysis of Nonstationary Time Series and the Business Cycle”, *Econometrica*, 57, 357-384.
- Hamilton, J.D. (1996), “Specification Testing in Markov-Switching Time Series Models”, *Journal of Econometrics*, 70, 127-157.
- Hamilton, J.D. (2001), “A Parametric Approach to Flexible Nonlinear Inference”, *Econometrica*, 69, 537-573.
- Hamilton, J.D. and O. Jordà (2002), “A Model of the Federal Funds Target,” *Journal of Political Economy*, 110, 1135-1167.
- Hansen, B.E. (1996), “Inference when a Nuisance Parameter is not Identified under the Null Hypothesis”, *Econometrica*, 64, 413-430.

- Hansen, P.R. (2005), “A test for superior predictive ability”, *Journal of Business and Economic Statistics* 23: 365-380.
- Harding, D. and A. Pagan (2002), “Dissecting the Cycle: A Methodological Investigation”, *Journal of Monetary Economics*, 49, 365-381.
- Härdle, W. and A. Tsybakov (1997), “Local polynomial estimators of the volatility function in nonparametric autoregression”, *Journal of Econometrics* 81, 233-242.
- Hong, Y. and J. Chung (2003), “Are the Directions of Stock Price Changes Predictable? Statistical Theory and Evidence”, Cornell University.
- Hong, Y. and T.-H. Lee (2003a), “Inference on Predictability of Foreign Exchange Rates via Generalized Spectrum and Nonlinear Time Series Models”, *Review of Economics and Statistics*, 85(4), 1048-1062.
- Hong, Y. and T.-H. Lee (2003b), “Diagnostic Checking for Adequacy of Nonlinear Time Series Models”, *Econometric Theory*, 19(6), 1065-1121.
- Hornik, K., M. Stinchcombe, and H. White (1989), “Multi-Layer Feedforward Networks Are Universal Approximators,” *Neural Network*, 2: 359-366.
- Huang, Y.-L. and C.-M. Kuan (2007), “Re-examining Long-Run PPP under an Innovation Regime Switching Framework”, Academia Sinica, Taipei.
- Inoue, A. and L. Kilian (2008), “How Useful is Bagging in Forecasting Economic Time Series? A Case Study of U.S. CPI Inflation”, forthcoming, *Journal of the American Statistical Association*.
- Judd, K.L. (1998), *Numerical Methods in Economics*, MIT Press.
- Koenker, R. and G. Bassett Jr. (1978), “Regression Quantiles”, *Econometrica* 46(1): 33-50.
- Kuan, C.-M., Y.-L. Huang, and R.-S. Tsay (2005), “An unobserved component model with switching permanent and transitory innovations”, *Journal of Business and Economic Statistics*, 23, 443-454.
- Kullback, L. and R. A. Leibler (1951), “On Information and Sufficiency”, *Annals of Mathematical Statistics* 22, 79-86.
- Lee, T.-H. and A. Ullah (2001), “Nonparametric Bootstrap Tests for Neglected Nonlinearity in Time Series Regression Models”, *Journal of Nonparametric Statistics* 13, 425-451.
- Lee, T.-H., H. White and C.W.J. Granger (1993), “Testing for Neglected Nonlinearity in Time Series Models: A Comparison of Neural Network Methods and Alternative Tests”, *Journal of Econometrics* 56, 269-290.
- Lee, T.-H. and Y. Yang (2006), “Bagging Binary and Quantile Predictors for Time Series”, *Journal of Econometrics* 135, 465-497.
- Lewellen, J. (2004), “Predicting Returns with Financial Ratios”, *Journal of Financial Economics* 74, 209-235.

- Lintner, J. (1965), "Security Prices, Risk and Maximal Gains from Diversification", *Journal of Finance* 20, 587-615.
- Linton, O. and Y.-J. Whang (2004), "A Quantile Approach to Evaluating Directional Predictability", Cowles Foundation Discussion Paper No. 1454.
- Lo, A.W. and A.C. MacKinlay (1999). *A Non-Random Walk Down Wall Street*, Princeton University Press, Princeton.
- Long, X., L. Su, and A. Ullah (2007), "Estimation and Forecasting of Dynamic Conditional Covariance: A Semiparametric Multivariate Model," UC Riverside.
- Lopez, J.A. (2001), "Evaluating the Predictive Accuracy of Volatility Models", *Journal of Forecasting* 20, 87-109.
- Ludvigson, S. and S. Ng (2007), "The Empirical Risk Return Relation: A Factor Analysis Approach", *Journal of Financial Economics* 83, 171-222.
- Lundbergh, S. and T. Teräsvirta (2002), "Forecasting with smooth transition autoregressive models", *A Companion to Economic Forecasting*, edited by M.P. Clements and D.F. Hendry, Chapter 21, Blackwell.
- Luukkonen, R., P. Saikkonen, and T. Teräsvirta (1988), "Testing Linearity in Univariate Time Series Models", *Scandinavian Journal of Statistics* 15, 161-175.
- Maheu, J.M. and T.H. McCurdy (2000), "Identifying Bull and Bear Markets in Stock Returns", *Journal of Business and Economic Statistics* 18, 100-112.
- Manski, C.F. (1975), "Maximum Score Estimation of the Stochastic Utility Model of Choice", *Journal of Econometrics* 3(3), 205-228.
- Markowitz, H. (1959), *Portfolio Selection: Efficient Diversification of Investments*, John Wiley, New York.
- Marsh, I.W. (2000), "High-frequency Markov Switching Models in the Foreign Exchange Market", *Journal of Forecasting* 19, 123-134.
- McAleer, M. and M.C. Medeiros (2007), "A multiple regime smooth transition heterogeneous autoregressive model for long memory and asymmetries", *Journal of Econometrics*, forthcoming.
- Nadaraya, É.A. (1964), "On Estimating Regression", *Theory of Probability and its Applications* 9, 141-142.
- Nelson, C.R. and Siegel, A.F. (1987), "Parsimonious Modeling of Yield Curves," *Journal of Business* 60, 473-489.
- Nelson, D.B. (1991), "Conditional Heteroscedasticity in Asset Returns: A New Approach", *Econometrica*, 59(2), 347-370.
- Pagan, A.R. and A. Ullah (1999), *Nonparametric Econometrics*, Cambridge University Press.

- Patton, A.J. and A. Timmermann (2007), “Testing Forecast Optimality Under Unknown Loss”, *Journal of the American Statistical Association*, Volume 102, Number 480, 1172-1184.
- Perez-Quiros, G. and A. Timmermann (2001), “Business Cycle Asymmetries in Stock Returns: Evidence form Higher Order Moments and Conditional Densities”, *Journal of Econometrics* 103, 259-306.
- Poon, S.-H. and C.W.J. Granger (2002), “Forecasting Volatility in Financial Markets: A Review”, Strathclyde University and UCSD, Working Paper.
- Raj, B. and A. Ullah (1981), *Econometrics: A Varying Coefficients Approach*, Croom Helm, London.
- Riskmetrics (1995), *Technical Manual*, 3ed.
- Romano, J.P. and M. Wolf (2005), “Stepwise multiple testing as formalized data snooping”, *Econometrica* 73, 1237-1282.
- Ross, S. (1976), “The Arbitrage Theory of Capital Asset Pricing”, *Journal of Economic Theory* 13, 341-360.
- Ruppert, D. and M.P. Wand (1994), “Multivariate Locally Weighted Least Squares Regression”, *Annals of Statistics* 22, 1346-1370.
- Sawa, T. (1978), “Information Criteria for Discriminating among Alternative Regression Models”, *Econometrica* 46: 1273-1291
- Sentana, E. (1995), “Quadratic ARCH models”, *Review of Economic Studies* 62(4), 639-661.
- Sichel, D.E. (1994), “Inventories and the Three Phases of the Business Cycle”, *Journal of Business and Economic Statistics* 12, 269-277.
- Sharpe, W. (1964), “Capital Asset Prices: A Theory of Market Equilibrium under Conditions of Risk”, *Journal of Finance* 19, 425-442.
- Stinchcombe, M. and H. White (1998), “Consistent Specification Testing with Nuisance Parameters Present only under the Alternative,” *Econometric Theory*, 14: 295-325.
- Stock, J.H. (2001), “Forecasting Economic Time Series”, *A Companion to Theoretical Econometrics*, edited by B.P. Baltagi, Chapter 27, Blackwell.
- Stock, J.H. and M.W. Watson (2002), “Forecasting Using Principal Components from a Large Number of Predictors,” *Journal of the American Statistical Association* 97, 1167-1179.
- Stock, J.H. and M.W. Watson (2002), “Forecasting Using Many Predictors,” *Handbook of Economic Forecasting*, Volume 1, edited by G. Elliott, C.W.J. Granger, and A. Timmermann, Amsterdam: Elsevier.
- Stone, C.J. (1977), “Consistent Nonparametric Regression”, *Annals of Statistics* 5, 595-645.

- Taylor, S.J. (1986), *Modelling Financial Time Series*, Wiley, New York.
- Teräsvirta, T. (1994), “Specification, Estimation and Evaluation of Smooth Transition Autoregressive Models”, *Journal of the American Statistical Association* 89, 208-218.
- Teräsvirta, T. (2006), “Forecasting economic variables with nonlinear models”, in G. Elliott, C.W.J. Granger and A. Timmermann (eds.), *Handbook of Economic Forecasting*, Vol .1, 413-457, Amsterdam: Elsevier.
- Teräsvirta, T. and H. Anderson (1992), “Characterizing Nonlinearities in Business Cycles using Smooth Transition Autoregressive Models”, *Journal of Applied Econometrics* 7, 119-139.
- Teräsvirta, T., C.-F. Lin and C.W.J. Granger (1993), “Power of the Neural Network Linearity Test”, *Journal of Time Series Analysis* 14, 209-220.
- Tong, H. (1983), *Threshold Models in Nonlinear Time Series Analysis*, Springer-Verlag, New York.
- Tong, H. (1990), *Nonlinear Time Series: A Dynamical Systems Approach*, Oxford University Press, Oxford.
- Trippi, R. and Turban, E. (1992). *Neural Networks in Finance and Investing: Using Artificial Intelligence to Improve Real World Performance*, New York: McGraw-Hill.
- Tsay, R.S. (1998), “Testing and Modeling Multivariate Threshold Models”, *Journal of the American Statistical Association* 93, 1188-1202.
- Tsay, R.S. (2002), “Nonlinear Models and Forecasting”, *A Companion to Economic Forecasting*, edited by M.P. Clements and D.F. Hendry, Chapter 20, Blackwell.
- Tsay, R.S. (2005), *Analysis of Financial Time Series*, 2ed., Wiley.
- Varian, H.R. (1975), “A Bayesian Approach to Real Estate Assessment”, in *Studies in Bayesian Econometrics and Statistics in Honor of L.J. Savage*, eds. S.E. Fienberg and A. Zellner, Amsterdam: North Holland, pp 195-208.
- Watson, G.S. (1964), “Smooth Regression Analysis”, *Sankhya*, Series A, 26, 359-372.
- West, K.D. (1996), “Asymptotic Inference about Prediction Ability”, *Econometrica* 64, 1067-1084.
- West, K.D., H.J. Edison, and D. Cho (1993), “A Utility Based Comparison of Some Models of Exchange Rate Volatility”, *Journal of International Economics* 35, 23-45.
- White, H. (1989), “An Additional Hidden Unit Tests for Neglected Nonlinearity in Multilayer Feedforward Networks,” *Proceedings of the International Joint Conference on Neural Networks*, Washington, DC. (IEEE Press, New York, NY), II: 451-455.
- White, H. (1994), *Estimation, Inference, and Specification Analysis*, Cambridge University Press.

- White, H. (2000), "A Reality Check for Data Snooping," *Econometrica* 68(5), 1097-1126.
- White, H. (2006), "Approximate Nonlinear Forecasting Methods", *Handbook of Economic Forecasting*, Volume 1, Ch. 9., edited by G. Elliott, C.W.J. Granger, and A. Timmermann, Amsterdam: Elsevier.
- Zellner, A. (1986), "Bayesian Estimation and Prediction Using Asymmetric Loss Functions". *Journal of the American Statistical Association* 81, 446-451.
- Ziegelmann, F.A.(2002), "Nonparametric estimation of volatility functions: the local exponential estimator," *Econometric Theory* 18, 985-991.