

Nonparametric Econometrics: Theory and Applications¹

ZONGWU CAI^{a,b}

E-mail address: zcai@uncc.edu

^aDepartment of Mathematics & Statistics and Department of Economics,
University of North Carolina, Charlotte, NC 28223, U.S.A.

^bWang Yanan Institute for Studies in Economics, Xiamen University, China

May 6, 2007

©2007, ALL RIGHTS RESERVED by ZONGWU CAI

¹This manuscript may be printed and reproduced for individual or instructional use, but may not be printed for commercial purposes.

Preface

This is the advanced level of nonparametric econometrics with theory and applications. Here, the focus is on both the theory and the skills of analyzing real data using nonparametric econometric techniques and statistical softwares such as **R**. This is along the line with the spirit “STRONG THEORETICAL FOUNDATION and SKILL EXCELLENCE”. In other words, this course covers the advanced topics in analysis of economic and financial data using nonparametric techniques, particularly in nonlinear time series models and some models related to economic and financial applications. The topics covered start from classical approaches to modern modeling techniques even up to the research frontiers. The difference between this course and others is that you will learn not only the theory but also step by step how to build a model based on data (or so-called “let data speak themselves”) through real data examples using statistical softwares or how to explore the real data using what you have learned. Therefore, there is no a single book serviced as a textbook for this course so that materials from some books and articles will be provided. However, some necessary handouts, including computer codes like **R** codes, will be provided with your help (You might be asked to print out the materials by yourself).

Several projects, including the heavy computer works, are assigned throughout the term. The purpose of projects is to train student to understand the theoretical concepts and to know how to apply the methodology to real problems. The group discussion is allowed to do the projects, particularly writing the computer codes. But, writing the final report to each project must be in your own language. Copying each other will be regarded as a cheating. If you use the **R** language, similar to **SPLUS**, you can download it from the public web site at <http://www.r-project.org/> and install it into your own computer or you can use PCs at our labs. You are STRONGLY encouraged to use (but not limited to) the package **R** since it is a very convenient programming language for doing statistical analysis and Monte Carol simulations as well as various applications in quantitative economics and finance. Of course, you are welcome to use any one of other packages such as **SAS**, **GAUSS**, **STATA**, **SPSS** and **EVIEW**. But, I might not have an ability of giving you a help if doing so.

Contents

1	Package R and Simple Applications	1
1.1	Computational Toolkits	1
1.2	How to Install R ?	2
1.3	Data Analysis and Graphics Using R – An Introduction (109 pages)	3
1.4	CRAN Task View: Empirical Finance	3
1.5	CRAN Task View: Computational Econometrics	7
2	Estimation of Covariance Matrix	12
2.1	Methodology	12
2.2	An Example	14
2.3	R Commands	17
2.4	Reading Materials – the paper by Zeileis (2004)	19
2.5	Computer Codes	19
2.6	References	22
3	Density, Distribution & Quantile Estimations	23
3.1	Time Series Structure	23
3.1.1	Mixing Conditions	23
3.1.2	Martingale and Mixingale	25
3.2	Nonparametric Density Estimate	26
3.2.1	Asymptotic Properties	27
3.2.2	Optimality	30
3.2.3	Boundary Problems	32
3.2.4	Bandwidth Selection	35
3.2.5	Project for Density Estimation	38
3.2.6	Multivariate Density Estimation	39
3.2.7	Reading Materials	41
3.3	Distribution Estimation	41
3.3.1	Smoothed Distribution Estimation	41
3.3.2	Relative Efficiency and Deficiency	43
3.4	Quantile Estimation	44
3.4.1	Value at Risk	45
3.4.2	Nonparametric Quantile Estimation	46
3.5	Computer Code	47
3.6	References	50

4	Nonparametric Regression Models	54
4.1	Prediction and Regression Functions	54
4.2	Kernel Estimation	55
4.2.1	Asymptotic Properties	56
4.2.2	Boundary Behavior	58
4.3	Local Polynomial Estimate	59
4.3.1	Formulation	59
4.3.2	Implementation in R	60
4.3.3	Complexity of Local Polynomial Estimator	62
4.3.4	Properties of Local Polynomial Estimator	64
4.3.5	Bandwidth Selection	68
4.4	Project for Regression Function Estimation	70
4.5	Functional Coefficient Model	71
4.5.1	Model	71
4.5.2	Local Linear Estimation	72
4.5.3	Bandwidth Selection	73
4.5.4	Smoothing Variable Selection	74
4.5.5	Goodness-of-Fit Test	74
4.5.6	Asymptotic Results	76
4.5.7	Conditions and Proofs	78
4.5.8	Monte Carlo Simulations and Applications	85
4.6	Additive Model	87
4.6.1	Model	87
4.6.2	Backfitting Algorithm	90
4.6.3	Projection Method	91
4.6.4	Two-Stage Procedure	93
4.6.5	Monte Carlo Simulations and Applications	94
4.6.6	New Developments	95
4.6.7	Additive Model to to Boston House Price Data	95
4.7	Computer Code	95
4.7.1	Example 4.1	95
4.7.2	Codes for Additive Modeling Analysis of Boston Data	101
4.8	References	102
5	Nonparametric Quantile Models	107
5.1	Introduction	107
5.2	Modeling Procedures	112
5.2.1	Local Linear Quantile Estimate	112
5.2.2	Asymptotic Results	114
5.2.3	Bandwidth Selection	119
5.2.4	Covariance Estimate	121
5.3	Empirical Examples	122
5.3.1	A Simulated Example	122
5.3.2	Real Data Examples	125
5.4	Derivations	135

5.5	Proofs of Lemmas	138
5.6	Computer Codes	142
5.7	References	142
6	Conditional VaR and Expected Shortfall	148
6.1	Introduction	148
6.2	Setup	152
6.3	Nonparametric Estimating Procedures	153
6.3.1	Estimation of Conditional PDF and CDF	154
6.3.2	Estimation of Conditional VaR and ES	157
6.4	Distribution Theory	157
6.4.1	Assumptions	157
6.4.2	Asymptotic Properties for Conditional PDF and CDF	159
6.4.3	Asymptotic Theory for CVaR and CES	162
6.5	Empirical Examples	165
6.5.1	Bandwidth Selection	166
6.5.2	Simulated Examples	166
6.5.3	Real Examples	170
6.6	Proofs of Theorems	173
6.7	Proofs of Lemmas	179
6.8	Computer Codes	182
6.9	References	182

List of Tables

3.1	Sample sizes required for p-dimensional nonparametric regression to have comparable performance with that of 1-dimensional nonparametric regression using size 100	40
-----	--	----

List of Figures

2.1	Time plots of U.S. weekly interest rates (in percentages) from January 5, 1962 to September 10, 1999. The solid line (black) is the Treasury 1-year constant maturity rate and the dashed line the Treasury 3-year constant maturity rate (red).	15
2.2	Scatterplots of U.S. weekly interest rates from January 5, 1962 to September 10, 1999: the left panel is 3-year rate versus 1-year rate, and the right panel is changes in 3-year rate versus changes in 1-year rate.	15
2.3	Residual series of linear regression Model I for two U.S. weekly interest rates: the left panel is time plot and the right panel is ACF.	16
2.4	Time plots of the change series of U.S. weekly interest rates from January 12, 1962 to September 10, 1999: changes in the Treasury 1-year constant maturity rate are denoted by black solid line, and changes in the Treasury 3-year constant maturity rate are indicated by red dashed line.	16
2.5	Residual series of the linear regression models: Model II (top) and Model III (bottom) for two change series of U.S. weekly interest rates: time plot (left) and ACF (right).	17
3.1	Bandwidth is taken to be 0.25, 0.5, 1.0 and the optimal one (see later) with the Epanechnikov kernel.	30
3.2	The left panel is for the built-in function density() and the right panel is for own code.	31
4.1	Scatterplots of Δx_t , $ \Delta x_t $, and $(\Delta x_t)^2$ versus x_t with the smoothed curves computed using scatter.smooth() and the local constant estimation.	61
4.2	Scatterplots of Δx_t , $ \Delta x_t $, and $(\Delta x_t)^2$ versus x_t with the smoothed curves computed using scatter.smooth() and the local linear estimation.	62
4.3	The results from model (4.66).	96
4.4	(a) Residual plot for model (4.66). (b) Plot of $g_1(x_6)$ versus x_6 . (c) Residual plot for model (4.67). (d) Density estimate of Y	97

5.1	<i>Simulated Example</i> : The plots of the estimated coefficient functions for three quantiles $\tau = 0.05$ (dashed line), $\tau = 0.50$ (dotted line), and $\tau = 0.95$ (dot-dashed line) with their true functions (solid line): $\sigma(u)$ versus u in (a), $a_1(u)$ versus u in (b), and $a_2(u)$ versus u in (c), together with the 95% point-wise confidence interval (thick line) with the bias ignored for the $\tau = 0.5$ quantile estimate.	125
5.2	<i>Boston Housing Price Data</i> : Displayed in (a)-(d) are the scatter plots of the house price versus the covariates U , X_1 , X_2 and $\log(X_2)$, respectively.	127
5.3	<i>Boston Housing Price Data</i> : The plots of the estimated coefficient functions for three quantiles $\tau = 0.05$ (solid line), $\tau = 0.50$ (dashed line), and $\tau = 0.95$ (dotted line), and the mean regression (dot-dashed line): $\hat{a}_{0,\tau}(u)$ and $\hat{a}_0(u)$ versus u in (e), $\hat{a}_{1,\tau}(u)$ and $\hat{a}_1(u)$ versus u in (f), and $\hat{a}_{2,\tau}(u)$ and $\hat{a}_2(u)$ versus u in (g). The thick dashed lines indicate the 95% point-wise confidence interval for the median estimate with the bias ignored.	128
5.4	<i>Exchange Rate Series</i> : (a) Japanese-dollar exchange rate return series $\{Y_t\}$; (b) autocorrelation function of $\{Y_t\}$; (c) moving average trading technique rule.	131
5.5	<i>Exchange Rate Series</i> : The plots of the estimated coefficient functions for three quantiles $\tau = 0.05$ (solid line), $\tau = 0.50$ (dashed line), and $\tau = 0.95$ (dotted line), and the mean regression (dot-dashed line): $\hat{a}_{0,0.50}(u)$ and $\hat{a}_0(u)$ versus u in (d), $\hat{a}_{0,0.05}(u)$ and $\hat{a}_{0,0.95}(u)$ versus u in (e), $\hat{a}_{1,\tau}(u)$ and $\hat{a}_1(u)$ versus u in (f), and $\hat{a}_{2,\tau}(u)$ and $\hat{a}_2(u)$ versus u in (g). The thick dashed lines indicate the 95% point-wise confidence interval for the median estimate with the bias ignored.	134
6.1	Simulation results for Example 1 when $p = 0.05$. Displayed in (a) - (c) are the true CVaR functions (solid lines), the estimated WDKLL CVaR functions (dashed lines), and the estimated NW CVaR functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both the WDKLL and NW estimations of CVaR are plotted in (d).	167
6.2	Simulation results for Example 1 when $p = 0.05$. Displayed in (a) - (c) are the true CES functions (solid lines), the estimated WDKLL CES functions (dashed lines), and the estimated NW CES functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both the WDKLL and NW estimations of CES are plotted in (d).	168
6.3	Simulation results for Example 1 when $p = 0.01$. Displayed in (a) - (c) are the true CVaR functions (solid lines), the estimated WDKLL CVaR functions (dashed lines), and the estimated NW CVaR functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both WDKLL and NW estimation of the conditional VaR are plotted in (d).	169
6.4	Simulation results for Example 1 when $p = 0.01$. Displayed in (a) - (c) are the true CES functions (solid lines), the estimated WDKLL CES functions (dashed lines), and the estimated NW CES functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both the WDKLL and NW estimations of CVaR are plotted in (d).	170

6.5	Simulation results for Example 2 when $p = 0.05$. (a) Boxplots of MADEs for both the WDKLL and NW CVaR estimates. (b) Boxplots of MADEs for Both the WDKLL and NW CES estimates.	171
6.6	(a) 5% CVaR estimate for DJI index. (b) 5% CES estimate for DJI index. .	172
6.7	(a) 5% CVaR estimates for IBM stock returns. (b) 5% CES estimates for IBM stock returns index. (c) 5% CVaR estimates for three different values of lagged negative IBM returns ($-0.275, -0.025, 0.325$). (d) 5% CVaR estimates for three different values of lagged negative DJI returns ($-0.225, 0.025, 0.425$). (e) 5% CES estimates for three different values of lagged negative IBM returns ($-0.275, -0.025, 0.325$). (f) 5% CES estimates for three different values of lagged negative DJI returns ($-0.225, 0.025, 0.425$).	173

Chapter 1

Package R and Simple Applications

1.1 Computational Toolkits

When you work with large data sets, messy data handling, models, etc, you need to choose the computational tools that are useful for dealing with these kinds of problems. There are “**menu driven systems**” where you click some buttons and get some work done - but these are useless for anything nontrivial. To do serious economics and finance in the modern days, you have to write computer programs. And this is true of any field, for example, applied econometrics, empirical macroeconomics - and not just of “computational finance” which is a hot buzzword recently.

The **question** is how to choose the computational tools. According to Ajay Shah (December 2005), you should pay attention to three elements: **price, freedom, elegant** and **powerful computer science**, and **network effects**. Low price is better than high price. Price = 0 is obviously best of all. Freedom here is in many aspects. A good software system is one that does not tie you down in terms of hardware/OS, so that you are able to keep moving. Another aspect of freedom is in working with colleagues, collaborators and students. With commercial software, this becomes a problem, because your colleagues may not have the same software that you are using. Here free software really wins spectacularly. Good practice in research involves a great accent on reproducibility. Reproducibility is important both so as to avoid mistakes, and because the next person working in your field should be standing on your shoulders. This requires an ability to release code. This is only possible with free software. Systems like **SAS** and **Gauss** use archaic computer science. The code is inelegant. The language is not powerful. In this day and age, writing **C** or **Fortran** by hand is “too low level”. Hell, with **Gauss**, even a minimal thing like online help is tawdry.

One prefers a system to be built by people who know their computer science - it should be an elegant, powerful language. All standard CS knowledge should be nicely in play to give you a gorgeous system. Good computer science gives you more productive humans. Lots of economists use **Gauss**, and give out **Gauss** source code, so there is a network effect in favor of **Gauss**. A similar thing is right now happening with statisticians and **R**.

Here I cite comparisons among most commonly used packages (see Ajay Shah (December 2005)); see the web site at <http://www.mayin.org/ajayshah/COMPUTING/mytools.html>.

R is a very convenient programming language for doing statistical analysis and Monte Carol simulations as well as various applications in quantitative economics and finance. Indeed, we prefer to think of it of an environment within which statistical techniques are implemented. I will teach it at the introductory level, but NOTICE that you will have to learn **R** on your own. Note that about 97% of commands in **S-PLUS** and **R** are same. In particular, for analyzing time series data, **R** has a lot of bundles and packages, which can be downloaded for free, for example, at <http://www.r-project.org/>.

R, like **S**, is designed around a true computer language, and it allows users to add additional functionality by defining new functions. Much of the system is itself written in the **R** dialect of **S**, which makes it easy for users to follow the algorithmic choices made. For computationally-intensive tasks, **C**, **C++** and **Fortran** code can be linked and called at run time. Advanced users can write **C** code to manipulate **R** objects directly.

1.2 How to Install R ?

- (1) go to the web site <http://www.r-project.org/>;
- (2) click **CRAN**;
- (3) choose a site for downloading, say <http://cran.cnr.Berkeley.edu>;
- (4) click **Windows (95 and later)**;
- (5) click **base**;
- (6) click **R-2.4.1-win32.exe** (Version of 12-18-2006) to save this file first and then run it to install.

The basic R is installed into your computer. If you need to install other packages, you need

to do the followings:

- (7) After it is installed, there is an icon on the screen. Click the icon to get into **R**;
- (8) Go to the top and find **packages** and then click it;
- (9) Go down to **Install package(s)...** and click it;
- (10) There is a new window. Choose a location to download the packages, say **USA(CA1)**, move mouse to there and click OK;
- (11) There is a new window listing all packages. You can select any one of packages and click OK, or you can select all of them and then click OK.

1.3 Data Analysis and Graphics Using R – An Introduction (109 pages)

See the file **r-notes.pdf** (109 pages) which can be downloaded from

<http://www.math.uncc.edu/~zcair/r-notes.pdf>.

I encourage you to download this file and learn it by yourself.

1.4 CRAN Task View: Empirical Finance

This CRAN Task View contains a list of packages useful for empirical work in Finance, grouped by topic. Besides these packages, a very wide variety of functions suitable for empirical work in Finance is provided by both the basic **R** system (and its set of recommended core packages), and a number of other packages on the Comprehensive **R** Archive Network (CRAN). Consequently, several of the other CRAN Task Views may contain suitable packages, in particular the **Econometrics** Task View. The web site is

<http://cran.r-project.org/src/contrib/Views/Finance.html>

1. **Standard regression models:** Linear models such as ordinary least squares (**OLS**) can be estimated by **lm()** (from by the stats package contained in the basic R distribution). Maximum Likelihood (ML) estimation can be undertaken with the **optim()** function. Non-linear least squares can be estimated with the **nls()** function, as well as with **nlme()** from the **nlme** package. For the linear model, a variety of regression diagnostic tests are provided by the **car**, **lmtest**, **strucchange**, **urca**, **uroot**, and sandwich packages. The **Rcmdr** and **Zelig** packages provide user interfaces that may be of interest as well.

2. **Time series:** Classical time series functionality is provided by the arima() and KalmanLike() commands in the basic R distribution. The dse packages provides a variety of more advanced estimation methods; fracdiff can estimate fractionally integrated series; longmemo covers related material. For volatility modeling, the standard GARCH(1,1) model can be estimated with the garch() function in the tseries package. Unit root and cointegration tests are provided by tseries, urca and uroot. The Rmetrics packages fSeries and fMultivar contain a number of estimation functions for ARMA, GARCH, long memory models, unit roots and more. The ArDec implements autoregressive time series decomposition in a Bayesian framework. The dyn and dynlm are suitable for dynamic (linear) regression models. Several packages provide wavelet analysis functionality: rwt, wavelets, waveslim, wavethresh. Some methods from chaos theory are provided by the package tseriesChaos.
3. **Finance:** The Rmetrics bundle comprised of the fBasics, fCalendar, fSeries, fMultivar, fPortfolio, fOptions and fExtremes packages contains a very large number of relevant functions for different aspect of empirical and computational finance. The RQuantLib package provides several option-pricing functions as well as some fixed-income functionality from the QuantLib project to **R**. The portfolio package contains classes for equity portfolio management.
4. **Risk Management:** The VaR package estimates Value-at-Risk, and several packages provide functionality for Extreme Value Theory models: evd, evdbayes, evir, extRemes, ismec, POT. The mvtnorm package provides code for multivariate Normal and t-distributions. The Rmetrics packages fPortfolio and fExtremes also contain a number of relevant functions. The copula and fgac packages cover multivariate dependency structures using copula methods.
5. **Data and Date Management:** The its, zoo and fCalendar (part of Rmetrics) packages provide support for irregularly-spaced time series. fCalendar also addresses calendar issues such as recurring holidays for a large number of financial centers, and provides code for high-frequency data sets.

CRAN packages:

- * ArDec
- * car
- * copula
- * dse
- * dyn
- * dynlm
- * evd
- * evdbayes
- * evir
- * extRemes
- * fBasics (core)
- * fCalendar (core)
- * fExtremes (core)
- * fgac
- * fMultivar (core)
- * fOptions (core)
- * fPortfolio (core)
- * fracdiff
- * fSeries (core)
- * ismev
- * its (core)
- * lntest
- * longmemo
- * mvtnorm
- * portfolio
- * POT
- * Rcmdr

- * RQuantLib (core)
- * rwt
- * sandwich
- * strucchange
- * tseries (core)
- * tseriesChaos
- * urca (core)
- * uroot
- * VaR
- * wavelets
- * waveslim
- * wavethresh
- * Zelig
- * zoo (core)

Related links:

- * CRAN Task View: Econometrics. The web site is
<http://cran.cnr.berkeley.edu/src/contrib/Views/Econometrics.html>
or see the next section.
- * Rmetrics by Diethelm Wuertz contains a wealth of **R** code for Finance. The web site is
<http://www.itp.phys.ethz.ch/econophysics/R/>
- * Quantlib is a C++ library for quantitative finance. The web site is
<http://quantlib.org/>
- * Mailing list: R Special Interest Group Finance

1.5 CRAN Task View: Computational Econometrics

Base **R** ships with a lot of functionality useful for computational econometrics, in particular in the **stats** package. This functionality is complemented by many packages on CRAN, a brief overview is given below. There is also a considerable overlap between the tools for econometrics in this view and finance in the **Finance** view. Furthermore, the **finance SIG** is a suitable mailing list for obtaining help and discussing questions about both computational finance and econometrics. The packages in this view can be roughly structured into the following topics. The web site is

<http://cran.r-project.org/src/contrib/Views/Econometrics.html>

1. **Linear regression models:** Linear models can be fitted (via OLS) with `lm()` (from **stats**) and standard tests for model comparisons are available in various methods such as `summary()` and `anova()`. Analogous functions that also support asymptotic tests (z instead of t tests, and Chi-squared instead of F tests) and plug-in of other covariance matrices are `coeftest()` and `waldtest()` in **lmtest**. Tests of more general linear hypotheses are implemented in `linear.hypothesis()` in **car**. HC and HAC covariance matrices that can be plugged into these functions are available in **sandwich**. The packages **car** and **lmtest** also provide a large collection of further methods for diagnostic checking in linear regression models.
2. **Microeconometrics:** Many standard micro-econometric models belong to the family of generalized linear models (GLM) and can be fitted by `glm()` from package **stats**. This includes in particular logit and probit models for modelling choice data and poisson models for count data. Negative binomial GLMs are available via `glm.nb()` in package **MASS** from the **VR** bundle. Zero-inflated count models are provided in **zicounts**. Further over-dispersed and inflated models, including hurdle models, are available in package **pscl**. Bivariate poisson regression models are implemented in **bivpois**. Basic censored regression models (e.g., tobit models) can be fitted by `survreg()` in **survival**. Further more refined tools for microeconometrics are provided in **micEcon**. The package **bayesm** implements a Bayesian approach to microeconometrics and marketing. Inference for relative distributions is contained in package **reldist**.
3. **Further regression models:** Various extensions of the linear regression model and other model fitting techniques are available in base R and several CRAN packages.

Nonlinear least squares modelling is available in **nls()** in package **stats**. Relevant packages include **quantreg** (quantile regression), **sem** (linear structural equation models, including two-stage least squares), **systemfit** (simultaneous equation estimation), **betareg** (beta regression), **nlme** (nonlinear mixed-effect models), **VR** (multinomial logit models in package **nnet**) and **MNP** (Bayesian multinomial probit models). The packages **Design** and **Hmisc** provide several tools for extended handling of (generalized) linear regression models.

4. **Basic time series infrastructure:** The class **ts** in package **stats** is R's standard class for regularly spaced time series which can be coerced back and forth without loss of information to **zooreg** from package **zoo**. **zoo** provides infrastructure for both regularly and irregularly spaced time series (the latter via the class "zoo") where the time information can be of arbitrary class. Several other implementations of irregular time series building on the "POSIXt" time-date classes are available in **its**, **tseries** and **fCalendar** which are all aimed particularly at finance applications (see the **Finance** view).
5. **Time series modelling:** Classical time series modelling tools are contained in the **stats** package and include **arima()** for ARIMA modelling and Box-Jenkins-type analysis. Furthermore **stats** provides **StructTS()** for fitting structural time series and **decompose()** and **HoltWinters()** for time series filtering and decomposition. For estimating VAR models, several methods are available: simple models can be fitted by **ar()** in **stats**, more elaborate models are provided by **estVARXls()** in **dse** and a Bayesian approach is available in **MSBVAR**. A convenient interface for fitting dynamic regression models via OLS is available in **dynlm**; a different approach that also works with other regression functions is implemented in **dyn**. More advanced dynamic system equations can be fitted using **dse**. Unit root and cointegration techniques are available in **urca**, **uroot** and **tseries**. Time series factor analysis is available in **tsfa**.
6. **Matrix manipulations:** As a vector- and matrix-based language, base **R** ships with many powerful tools for doing matrix manipulations, which are complemented by the packages **Matrix** and **SparseM**.
7. **Inequality:** For measuring inequality, concentration and poverty the package **ineq** provides some basic tools such as Lorenz curves, Pen's parade, the Gini coefficient and

many more.

8. **Structural change:** **R** is particularly strong when dealing with structural changes and changepoints in parametric models, see strucchange and segmented.
9. **Data sets:** Many of the packages in this view contain collections of data sets from the econometric literature and the package Ecdat contains a complete collection of data sets from various standard econometric textbooks. micEcdat provides several data sets from the Journal of Applied Econometrics and the Journal of Business & Economic Statistics data archives. Package CDNmoney provides Canadian monetary aggregates and pwt provides the Penn world table.

CRAN packages:

- * bayesm
- * betareg
- * bivpois
- * car (core)
- * CDNmoney
- * Design
- * dse
- * dyn
- * dynlm
- * Ecdat
- * fCalendar
- * Hmisc
- * ineq
- * its
- * lmtest (core)
- * Matrix

- * micEcdat
- * micEcon
- * MNP
- * MSBVAR
- * nlme
- * pscl
- * pwt
- * quantreg
- * reldist
- * sandwich (core)
- * segmented
- * sem
- * SparseM
- * strucchange
- * systemfit
- * tseries (core)
- * tsfa
- * urca (core)
- * uroot
- * VR
- * zicounts
- * zoo (core)

Related links:

- * CRAN Task View: Finance. The web site is
[http://cran.cnr.berkeley.edu/src/contrib/ Views/Finance.html](http://cran.cnr.berkeley.edu/src/contrib/Views/Finance.html)
or see the above section.

- * Mailing list: **R** Special Interest Group Finance
- * A Brief Guide to **R** for Beginners in Econometrics. The web site is http://people.su.se/~ma/R_intro/.
- * **R** for Economists. The web site is http://www.mayin.org/ajayshah/KB/R/R_for_economists.html.

Chapter 2

Estimation of Covariance Matrix

2.1 Methodology

Consider a regression model stated in (2.1) below. There may exist situations which the error e_t has **serial correlations** and/or **conditional heteroscedasticity**, but the main objective of the analysis is to make inference concerning the regression coefficients β . When e_t has serial correlations, we assume that e_t follows an ARIMA type model but this assumption might not be always satisfied in some applications. Here, we consider a general situation without making this assumption. In situations under which the ordinary least squares estimates of the coefficients remain consistent, methods are available to provide **consistent estimate of the covariance matrix** of the coefficients. Two such methods are widely used in economics and finance. The first method is called **heteroscedasticity consistent (HC) estimator**; see Eicker (1967) and White (1980). The second method is called **heteroscedasticity and autocorrelation consistent (HAC) estimator**; see Newey and West (1987).

To ease in discussion, we write a regression model as

$$y_t = \beta^T \mathbf{x}_t + e_t, \quad (2.1)$$

where y_t is the dependent variable, $\mathbf{x}_t = (x_{1t}, \dots, x_{pt})^T$ is a p -dimensional vector of explanatory variables including constant and lagged variables, and $\beta = (\beta_1, \dots, \beta_p)^T$ is the parameter vector. The LS estimate of β is given by

$$\hat{\beta} = \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1} \sum_{t=1}^n \mathbf{x}_t y_t,$$

and the associated covariance matrix has the so-called “sandwich” form as

$$\Sigma_{\beta} = \text{Cov}(\hat{\beta}) = \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1} \mathbf{C} \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1} \stackrel{\text{if } e_t \text{ is iid}}{=} \sigma_e^2 \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1},$$

where \mathbf{C} is called the “meat” given by

$$\mathbf{C} = \text{Var} \left(\sum_{t=1}^n e_t \mathbf{x}_t \right),$$

σ_e^2 is the variance of e_t and is estimated by the variance of residuals of the regression. In the presence of serial correlations or conditional heteroscedasticity, the prior covariance matrix estimator is inconsistent, often resulting in inflating the t -ratios of $\hat{\beta}$.

The estimator of White (1980) is based on following:

$$\hat{\Sigma}_{\beta, hc} = \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1} \hat{\mathbf{C}}_{hc} \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1},$$

where with $\hat{e}_t = y_t - \hat{\beta}^T \mathbf{x}_t$ being the residual at time t ,

$$\hat{\mathbf{C}}_{hc} = \frac{n}{n-p} \sum_{t=1}^n \hat{e}_t^2 \mathbf{x}_t \mathbf{x}_t^T.$$

The estimator of Newey and West (1987) is

$$\hat{\Sigma}_{\beta, hac} = \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1} \hat{\mathbf{C}}_{hac} \left[\sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^T \right]^{-1},$$

where $\hat{\mathbf{C}}_{hac}$ is given by

$$\hat{\mathbf{C}}_{hac} = \sum_{t=1}^n \hat{e}_t^2 \mathbf{x}_t \mathbf{x}_t^T + \sum_{j=1}^l w_j \sum_{t=j+1}^n \{ \mathbf{x}_t \hat{e}_t \hat{e}_{t-j} \mathbf{x}_{t-j}^T + \mathbf{x}_{t-j} \hat{e}_{t-j} \hat{e}_t \mathbf{x}_t^T \}$$

with l is a truncation parameter and w_j is weight function such as the Barlett weight function defined by $w_j = 1 - j/(l+1)$. Other weight function can also used. Newey and West (1987) showed that if $l \rightarrow \infty$ and $l^4/T \rightarrow 0$, then $\hat{\mathbf{C}}_{hac}$ is a consistent estimator of \mathbf{C} . Newey and West (1987) suggested choosing l to be the integer part of $4(n/100)^{1/4}$ and Newey and West (1994) suggested using some adaptive (data-driven) methods to choose l ; see Newey and West (1994) for details. In general, this estimator essentially can use a nonparametric method to estimate the covariance matrix of $\sum_{t=1}^n e_t \mathbf{x}_t$ and a class of kernel-based **heteroskedasticity and autocorrelation consistent (HAC) covariance matrix**

estimators was introduced by Andrews (1991). For example, the Barlett weight w_j above can be replaced by $w_j = K(j/(l+1))$ where $K(\cdot)$ is a kernel function such as truncated kernel $K(x) = I(|x| \leq 1)$, the Tukey-Hanning kernel $K(x) = (1 + \cos(\pi x))/2$ if $|x| \leq 1$, the Parzen kernel

$$K(x) = \begin{cases} 1 - 6x^2 + 6|x|^3 & \text{for } 0 \leq |x| \leq 1/2, \\ 2(1 - |x|)^3 & \text{for } 1/2 \leq |x| \leq 1, \\ 0 & \text{otherwise,} \end{cases}$$

and the Quadratic spectral kernel

$$K(x) = \frac{25}{12\pi^2 x^2} \left(\frac{\sin(6\pi x/5)}{6\pi x/5} - \cos(6\pi x/5) \right).$$

Andrews (1991) suggested using the data-driven method to select the bandwidth l : $\hat{l} = 2.66(\hat{\alpha} T)^{1/5}$ for the Parzen kernel, $\hat{l} = 1.7462(\hat{\alpha} T)^{1/5}$ for the Tukey-Hanning kernel, and $\hat{l} = 1.3221(\hat{\alpha} T)^{1/5}$ for the quadratic spectral kernel, where

$$\hat{\alpha} = \frac{4 \sum_{i=1}^k \hat{\rho}_i^2 \hat{\sigma}_i^4 / (1 - \hat{\rho}_i)^8}{\sum_{i=1}^n \hat{\sigma}_i^4 / (1 - \hat{\rho}_i)^4}$$

with $\hat{\rho}_i$ and $\hat{\sigma}_i$ being parameters estimated from an AR(1) model for $\hat{u}_t = \hat{e}_t \mathbf{x}_t$.

2.2 An Example

Example 2.1: We consider the relationship between two U.S. weekly interest rate series: x_t : the 1-year Treasury constant maturity rate and y_t : the 3-year Treasury constant maturity rate. Both series have 1967 observations from January 5, 1962 to September 10, 1999 and are measured in percentages. The series are obtained from the Federal Reserve Bank of St Louis.

Figure 2.1 shows the time plots of the two interest rates with solid line denoting the 1-year rate and dashed line for the 3-year rate. The left panel of Figure 2.2 plots y_t versus x_t , indicating that, as expected, the two interest rates are highly correlated. A naive way to describe the relationship between the two interest rates is to use the simple model, **Model I**: $y_t = \beta_1 + \beta_2 x_t + e_t$. This results in a fitted model $y_t = 0.911 + 0.924 x_t + e_t$, with $\hat{\sigma}_e^2 = 0.538$ and $R^2 = 95.8\%$, where the standard errors of the two coefficients are 0.032 and 0.004, respectively. This simple model (Model I) confirms the high correlation between the two interest rates. However, the model is seriously inadequate as shown by Figure 2.3, which gives the time plot and ACF of its residuals. In particular, the sample ACF of the residuals

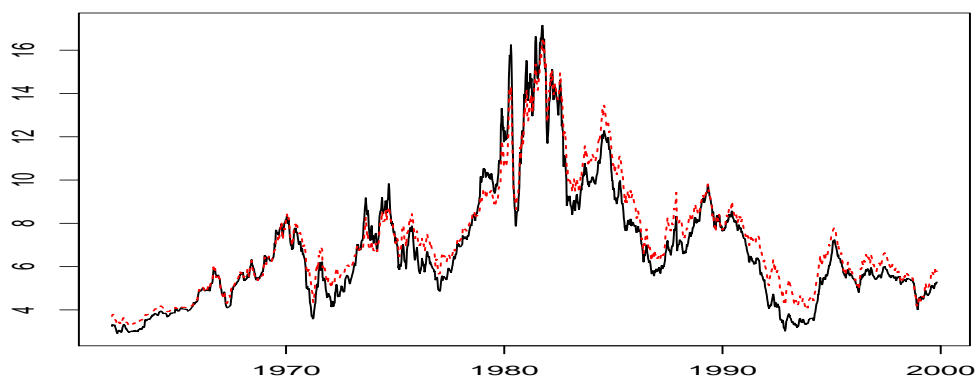


Figure 2.1: Time plots of U.S. weekly interest rates (in percentages) from January 5, 1962 to September 10, 1999. The solid line (black) is the Treasury 1-year constant maturity rate and the dashed line the Treasury 3-year constant maturity rate (red).

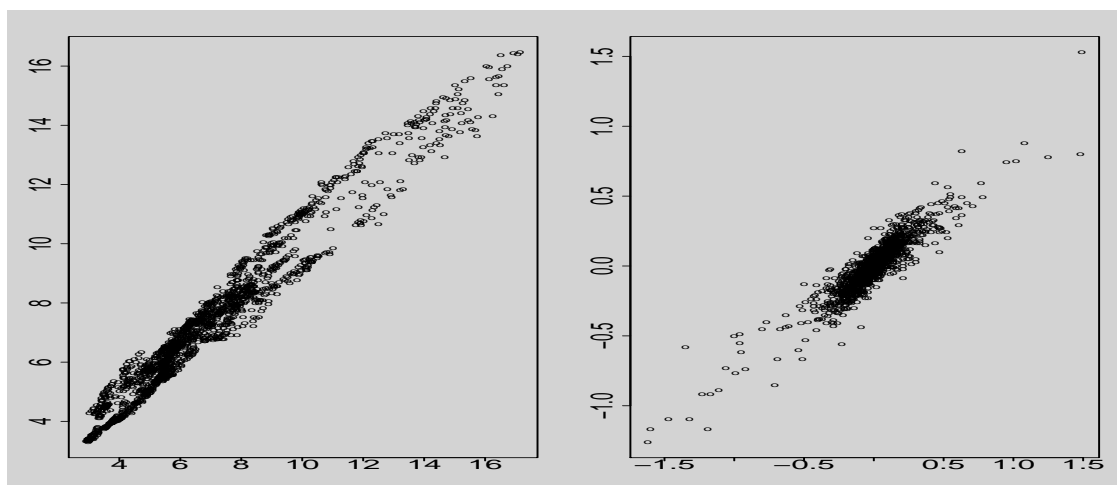


Figure 2.2: Scatterplots of U.S. weekly interest rates from January 5, 1962 to September 10, 1999: the left panel is 3-year rate versus 1-year rate, and the right panel is changes in 3-year rate versus changes in 1-year rate.

is highly significant and decays slowly, showing the pattern of a **unit root** nonstationary time series. The behavior of the residuals suggests that marked differences exist between the two interest rates. Using the modern econometric terminology, if one assumes that the two interest rate series are unit root nonstationary, then the behavior of the residuals indicates that the two interest rates are not **co-integrated**. In other words, the data fail to support the hypothesis that there exists a long-term equilibrium between the two interest rates. In some sense, this is not surprising because the pattern of “inverted yield curve” did occur during the data span. By the inverted yield curve, we mean the situation under which

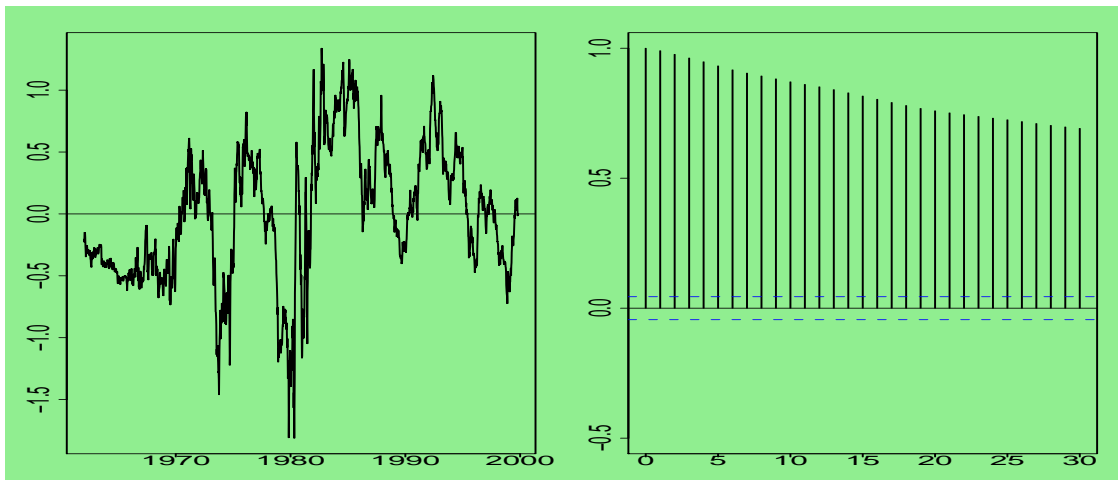


Figure 2.3: Residual series of linear regression Model I for two U.S. weekly interest rates: the left panel is time plot and the right panel is ACF.

interest rates are inversely related to their time to maturities.

The unit root behavior of both interest rates and the residuals leads to the consideration of the change series of interest rates. Let $\Delta x_t = y_t - y_{t-1} = (1 - L)x_t$ be changes in the 1-year interest rate and $\Delta y_t = y_t - y_{t-1} = (1 - L)y_t$ denote changes in the 3-year interest rate. Consider the linear regression, **Model II**: $\Delta y_t = \beta_1 + \beta_2 \Delta x_t + e_t$. Figure 2.4 shows time plots of the two change series, whereas the right panel of Figure 2.3 provides a scatterplot

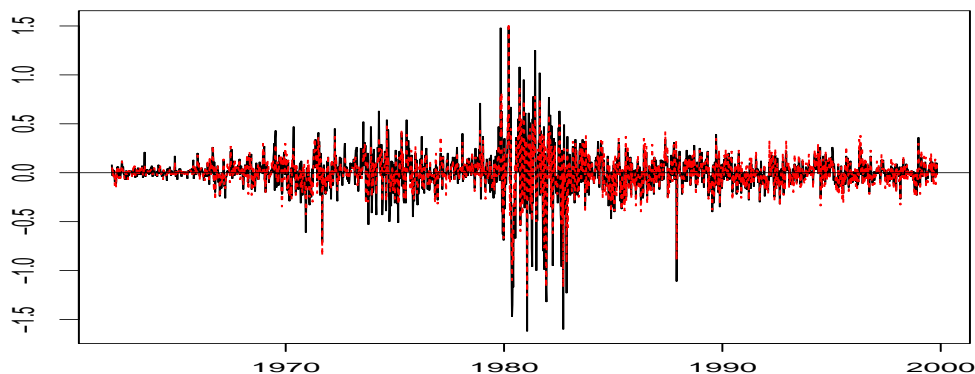


Figure 2.4: Time plots of the change series of U.S. weekly interest rates from January 12, 1962 to September 10, 1999: changes in the Treasury 1-year constant maturity rate are in denoted by black solid line, and changes in the Treasury 3-year constant maturity rate are indicated by red dashed line.

between them. The change series remain highly correlated with a fitted linear regression

model given by $\Delta y_t = 0.0002 + 0.7811 \Delta x_t + e_t$ with $\hat{\sigma}_e^2 = 0.0682$ and $R^2 = 84.8\%$. The standard errors of the two coefficients are 0.0015 and 0.0075, respectively. This model further confirms the strong linear dependence between interest rates. The two top panels of Figure 2.5 show the time plot (left) and sample ACF (right) of the residuals (Model II). Once again,

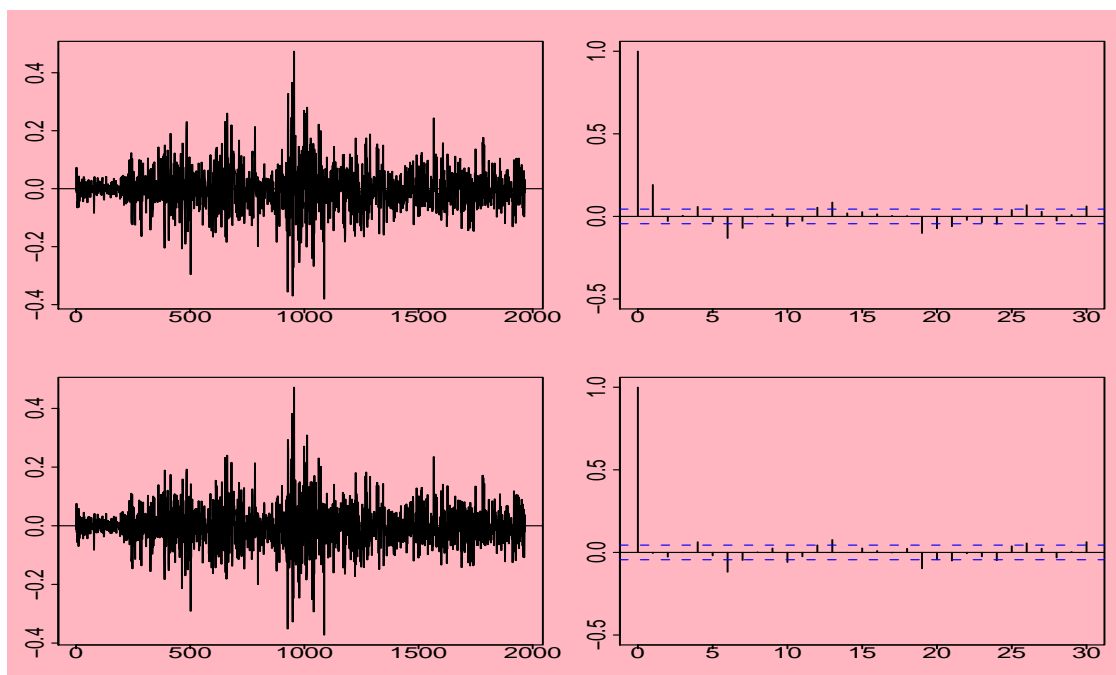


Figure 2.5: Residual series of the linear regression models: Model II (top) and Model III (bottom) for two change series of U.S. weekly interest rates: time plot (left) and ACF (right).

the ACF shows some significant serial correlation in the residuals, but the magnitude of the correlation is much smaller. This weak serial dependence in the residuals can be modeled by using the simple time series models discussed in the previous sections, and we have a linear regression with time series errors.

For illustration, we consider the first differenced interest rate series in Model II. The t -ratio of the coefficient of Δx_t is 104.63 if both serial correlation and conditional heteroscedasticity in residuals are ignored; it becomes 46.73 when the HC estimator is used, and it reduces to 40.08 when the HAC estimator is employed.

2.3 R Commands

To use HC or HAC estimator, we can use the package **sandwich** in **R** and the commands are **vcovHC()** or **vcovHAC()** or **meatHAC()**. There are a set of functions implementing

a class of kernel-based heteroskedasticity and autocorrelation consistent (HAC) covariance matrix estimators as introduced by Andrews (1991). In `vcovHC()`, these estimators differ in their choice of the ω_i in $\Omega = \text{Var}(e) = \text{diag}\{\omega_1, \dots, \omega_n\}$, an overview of the most important cases is given in the following:

$$\begin{aligned} \text{const} : \omega_i &= \hat{\sigma}^2 \\ \text{HC0} : \omega_i &= \hat{e}_i^2 \\ \text{HC1} : \omega_i &= \frac{n}{n-k} \hat{e}_i^2 \\ \text{HC2} : \omega_i &= \frac{\hat{e}_i^2}{1-h_i} \\ \text{HC3} : \omega_i &= \frac{\hat{e}_i^2}{(1-h_i)^2} \\ \text{HC4} : \omega_i &= \frac{\hat{e}_i^2}{(1-h_i)^{\delta_i}} \end{aligned}$$

where $h_i = H_{ii}$ are the diagonal elements of the hat matrix and $\delta_i = \min\{4, h_i/\bar{h}\}$.

```
vcovHC(x, type = c("HC3", "const", "HC", "HC0", "HC1", "HC2", "HC4"),
       omega = NULL, sandwich = TRUE, ...)
```

```
meatHC(x, type = , omega = NULL)
```

```
vcovHAC(x, order.by = NULL, prewhite = FALSE, weights = weightsAndrews,
        adjust = TRUE, diagnostics = FALSE, sandwich = TRUE, ar.method = "ols",
        data = list(), ...)
```

```
meatHAC(x, order.by = NULL, prewhite = FALSE, weights = weightsAndrews,
        adjust = TRUE, diagnostics = FALSE, ar.method = "ols", data = list())
```

```
kernHAC(x, order.by = NULL, prewhite = 1, bw = bwAndrews,
        kernel = c("Quadratic Spectral", "Truncated", "Bartlett", "Parzen",
        "Tukey-Hanning"), approx = c("AR(1)", "ARMA(1,1)"), adjust = TRUE,
        diagnostics = FALSE, sandwich = TRUE, ar.method = "ols", tol = 1e-7,
        data = list(), verbose = FALSE, ...)
```

```
weightsAndrews(x, order.by = NULL, bw = bwAndrews,
  kernel = c("Quadratic Spectral", "Truncated", "Bartlett", "Parzen",
    "Tukey-Hanning"), prewhite = 1, ar.method = "ols", tol = 1e-7,
  data = list(), verbose = FALSE, ...)

bwAndrews(x, order.by=NULL, kernel=c("Quadratic Spectral", "Truncated",
  "Bartlett", "Parzen", "Tukey-Hanning"), approx=c("AR(1)", "ARMA(1,1)"),
  weights = NULL, prewhite = 1, ar.method = "ols", data = list(), ...)
```

Also, there are a set of functions implementing the Newey and West (1987, 1994) heteroskedasticity and autocorrelation consistent (HAC) covariance matrix estimators.

```
NeweyWest(x, lag = NULL, order.by = NULL, prewhite = TRUE, adjust = FALSE,
  diagnostics = FALSE, sandwich = TRUE, ar.method = "ols", data = list(),
  verbose = FALSE)

bwNeweyWest(x, order.by = NULL, kernel = c("Bartlett", "Parzen",
  "Quadratic Spectral", "Truncated", "Tukey-Hanning"), weights = NULL,
  prewhite = 1, ar.method = "ols", data = list(), ...)
```

2.4 Reading Materials – the paper by Zeileis (2004)

2.5 Computer Codes

```
#####
# This is Example 2.1 for weekly interest rate series
#####

z<-read.table("c:/res-teach/xiada/teaching05-07/data/ex2-1.txt",header=F)
# first column=one year Treasury constant maturity rate;
# second column=three year Treasury constant maturity rate;
# third column=date
```

```

x=z[,1]
y=z[,2]
n=length(x)
u=seq(1962+1/52,by=1/52,length=n)
x_diff=diff(x)
y_diff=diff(y)
# Fit a simple regression model and examine the residuals
fit1=lm(y~x)           # Model 1
e1=fit1$resid

postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-2.1.eps",
horizontal=F,width=6,height=6)
matplot(u,cbind(x,y),type="l",lty=c(1,2),col=c(1,2),ylab="",xlab="")
dev.off()

postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-2.2.eps",
horizontal=F,width=6,height=6)
par(mfrow=c(1,2),mex=0.4,bg="light grey")
plot(x,y,type="p",pch="o",ylab="",xlab="",cex=0.5)
plot(x_diff,y_diff,type="p",pch="o",ylab="",xlab="",cex=0.5)
dev.off()

postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-2.3.eps",
horizontal=F,width=6,height=6)
par(mfrow=c(1,2),mex=0.4,bg="light green")
plot(u,e1,type="l",lty=1,ylab="",xlab="")
abline(0,0)
acf(e1,ylab="",xlab="",ylim=c(-0.5,1),lag=30,main="")
dev.off()

# Take different and fit a simple regression again
fit2=lm(y_diff~x_diff)           # Model 2

```

```

e2=fit2$resid

postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-2.4.eps",
horizontal=F,width=6,height=6)
matplot(u[-1],cbind(x_diff,y_diff),type="l",lty=c(1,2),col=c(1,2),
  ylab="",xlab="")
abline(0,0)
dev.off()

postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-2.5.eps",
horizontal=F,width=6,height=6)
par(mfrow=c(2,2),mex=0.4,bg="light pink")
ts.plot(e2,type="l",lty=1,ylab="",xlab="")
abline(0,0)
acf(e2, ylab="", xlab="",ylim=c(-0.5,),lag=30,main="")
# fit a model to the differenced data with an MA(1) error
fit3=arima(y_diff,xreg=x_diff, order=c(0,0,1))      # Model 3
e3=fit3$resid
ts.plot(e3,type="l",lty=1,ylab="",xlab="")
abline(0,0)
acf(e3, ylab="",xlab="",ylim=c(-0.5,1),lag=30,main="")
dev.off()
#####

library(sandwich) # HC and HAC are in the package "sandwich"
library(zoo)
z<-read.table("c:/res-teach/xiada/teaching05-07/data/ex2-1.txt",header=F)
x=z[,1]
y=z[,2]
x_diff=diff(x)
y_diff=diff(y)
# Fit a simple regression model and examine the residuals

```

```

fit1=lm(y_diff~x_diff)
print(summary(fit1))
e1=fit1$resid
# Heteroskedasticity-Consistent Covariance Matrix Estimation
#hc0=vcovHC(fit1,type="const")
#print(sqrt(diag(hc0)))
# type=c("const","HC","HC0","HC1","HC2","HC3","HC4")
# HC0 is the White estimator
hc1=vcovHC(fit1,type="HC0")
print(sqrt(diag(hc1)))
#Heteroskedasticity and autocorrelation consistent (HAC) estimation
#of the covariance matrix of the coefficient estimates in a
#(generalized) linear regression model.
hac1=vcovHAC(fit1,sandwich=T)
print(sqrt(diag(hac1)))

```

2.6 References

- Andrews, D.W.K. (1991). Heteroskedasticity and autocorrelation consistent covariance matrix estimation. *Econometrica*, **59**, 817-858.
- Eicker, F. (1967). Limit theorems for regression with unequal and dependent errors. In *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability* (L. LeCam and J. Neyman, eds.), University of California Press, Berkeley.
- Newey, W.K. and K.D. West (1987). A simple, positive-definite, heteroskedasticity and autocorrelation consistent covariance matrix. *Econometrica*, **55**, 703-708.
- Newey, W.K. and K.D. West (1994). Automatic lag selection in covariance matrix estimation. *Review of Economic Studies*, **61**, 631-653.
- White, H. (1980). A Heteroskedasticity consistent covariance matrix and a direct test for heteroskedasticity. *Econometrica*, **48**, 817-838.
- Zeileis, A. (2004). Econometric computing with HC and HAC covariance matrix estimators. *Journal of Statistical Software*, Volume **11**, Issue 10.
- Zeileis, A. (2006). Object-oriented computation of sandwich estimators. *Journal of Statistical Software*, **16**, 1-16.

Chapter 3

Density, Distribution & Quantile Estimations

3.1 Time Series Structure

Since most of economic and financial data are time series, we discuss our methodologies and theory under the framework of time series. For linear models, the time series structure can be often assumed to have some well known forms such as an **autoregressive moving average (ARMA)** model. However, under nonparametric setting, this assumption might not be valid. Therefore, we can assume a more general time series dependence, which is commonly used in the literature, described as follows.

3.1.1 Mixing Conditions

Mixing dependence is commonly used to characterize the dependent structure and it is often referred often to as **short range dependence** or **weak dependence**, which means that the distance between two observations goes farther and farther, the dependence becomes weaker and weaker very faster. It is well known that α -mixing includes many time series models as a special case. In fact, under very mild assumptions, **linear processes**, including **linear autoregressive** models and more generally **bilinear time series models** are α -mixing with mixing coefficients decaying exponentially. Many nonlinear time series models, such as **functional coefficient autoregressive processes with/without exogenous variables**, **nonlinear additive autoregressive models with/without exogenous variables**, **ARCH and GARCH type processes**, **stochastic volatility models**, and **many continuous time diffusion models (including the Black-Scholes type models)** are strong mixing under some mild conditions. See Genon-Caralot, Jeantheau and

Laredo (2000), Cai (2002), Carrasco and Chen (2002), and Chen and Tang (2005) for more details.

To simplify the notation, we only introduce mixing conditions for strictly stationary processes (in spite of the fact that a mixing process is not necessarily stationary). The idea is to define mixing coefficients to measure the strength (in different ways) of dependence for the two segments of a time series which are apart from each other in time. Let $\{X_t\}$ be a strictly stationary time series. For $n \geq 1$, define

$$\alpha(n) = \sup_{A \in \mathcal{F}_{-\infty}^0; B \in \mathcal{F}_n^\infty} |P(A)P(B) - P(AB)|,$$

where \mathcal{F}_i^j denotes the σ -algebra generated by $\{X_t; i \leq t \leq j\}$. Note that $\mathcal{F}_n^\infty \downarrow$. If $\alpha(n) \rightarrow 0$ as $n \rightarrow \infty$, $\{X_t\}$ is called α -mixing or **strong mixing**. There are several other mixing conditions such as **ρ -mixing**, **β -mixing**, **ϕ -mixing**, and **ψ -mixing**; see the books by Hall and Heyde (1980) and Fan and Yao (2003, page 68) for details. Indeed,

$$\begin{aligned} \beta(n) &= E \left\{ \sup_{A \in \mathcal{F}_n^\infty} |P(A) - P(A | X_t, t \leq 0)| \right\}, \\ \rho(n) &= \sup_{X \in \mathcal{F}_{-\infty}^0; Y \in \mathcal{F}_n^\infty} |\text{Corr}(X, Y)|, \\ \phi(n) &= \sup_{A \in \mathcal{F}_{-\infty}^0; B \in \mathcal{F}_n^\infty, P(A) > 0} |P(B) - P(B | A)|, \end{aligned}$$

and

$$\psi(n) = \sup_{A \in \mathcal{F}_{-\infty}^0; B \in \mathcal{F}_n^\infty, P(A)P(B) > 0} |1 - P(B | A)/P(B)|,$$

It is well known that the relationships among the mixing conditions are

$$\alpha(n) \leq \frac{1}{4} \rho(n) \leq \frac{1}{2} \phi(n),$$

so that ψ -mixing $\implies \phi$ -mixing $\implies \rho$ -mixing $\implies \alpha$ -mixing as well as β -mixing $\implies \alpha$ -mixing. Note that all our theoretical results are derived under mixing conditions. The following **inequalities** are very useful in applications, which can be found in the book by Hall and Heyde (1980, pp. 277-280).

Lemma 3.1: (Davydov's inequality) (i) If $E|X_i|^p + E|X_j|^q < \infty$ for some $p \geq 1$ and $q \geq 1$ and $1/p + 1/q < 1$, it holds that

$$|\text{Cov}(X_i, X_j)| \leq 8 \alpha^{1/r}(|j - i|) \|X_i\|_p \|X_j\|_q,$$

where $r = (1 - 1/p - 1/q)^{-1}$.

(ii) If $P(|X_i| \leq C_1) = 1$ and $P(|X_j| \leq C_2) = 1$ for some constants C_1 and C_2 , it holds that

$$|\text{Cov}(X_i, X_j)| \leq 4\alpha(|j - i|) C_1 C_2.$$

Note that if we allow X_i and X_j to be complex-valued random variables, (ii) still holds with the coefficient “4” on the RHS of the inequality replaced by “16”.

(iii) If $P(|X_i| \leq C_1) = 1$ and $E|X_j|^p < \infty$ for some constants C_1 and $p > 1$, then,

$$|\text{Cov}(X_i, X_j)| \leq 6 C_1 \|X_j\|_p \alpha^{1-p^{-1}}(|j - i|).$$

Lemma 3.2: If $E|X_i|^p + E|X_j|^q < \infty$ for some $p \geq 1$ and $q \geq 1$ and $1/p + 1/q = 1$, it holds that

$$|\text{Cov}(X_i, X_j)| \leq 2 \phi^{1/p}(|j - i|) \|X_i\|_p \|X_j\|_q.$$

3.1.2 Martingale and Mixingale

Martingale is very useful in applications. Here is the definition. Let $\{X_n, n \in \mathcal{N}\}$ be a sequence of random variables on a probability space (Ω, \mathcal{F}, P) , and let $\{\mathcal{F}_n, n \in \mathcal{N}\}$ be an increasing sequence of sub- σ -fields of \mathcal{F} . Suppose that the sequence $\{X_n, n \in \mathcal{N}\}$ satisfies

- (i) X_n is measurable with respect to \mathcal{F}_n ,
- (ii) $E|X_n| < \infty$,
- (iii) $E[X_n | \mathcal{F}_m] = X_m$ for all $m < n, n \in \mathcal{N}$.

Then, the sequence $\{X_n, n \in \mathcal{N}\}$ is said to be a martingale with respect to $\{\mathcal{F}_n, n \in \mathcal{N}\}$. We write that $\{X_n, \mathcal{F}_n, n \in \mathcal{N}\}$ is a martingale. If (i) and (ii) are retained and (iii) is replaced by the inequality $E[X_n | \mathcal{F}_m] \geq X_m$ ($E[X_n | \mathcal{F}_m] \leq X_m$), then $\{X_n, \mathcal{F}_n, n \in \mathcal{N}\}$ is called a **sub-martingale** (**super-martingale**). Define $Y_n = X_n - X_{n-1}$. Then $\{Y_n, \mathcal{F}_n, n \in \mathcal{N}\}$ is called a **martingale difference (MD)** if $\{X_n, \mathcal{F}_n, n \in \mathcal{N}\}$ is called a martingale. Clearly, $E[Y_n | \mathcal{F}_{n-1}] = 0$, which means that a MD is not predictable based on the past information. In a finance language, a stock market is **efficient**. Equivalently, it is a MD.

Another type of dependent structure is called **mixingale**, which is the so-called asymptotic **martingale**. The concept of mixingale, introduced by McLeish (1975), is defined as follows. Let $\{X_n, n \geq 1\}$ be a sequence of square-integrable random variables on a probability space (Ω, \mathcal{F}, P) , and let $\{\mathcal{F}_n, -\infty < n < \infty\}$ be an increasing sequence of sub- σ -fields of

\mathcal{F} . Then, $\{X_n, \mathcal{F}_n\}$ is called a **L_r -mixingale (difference) sequence** for $r \geq 1$ if, for some sequences of nonnegative constants c_n and ψ_m , where $\psi_m \rightarrow 0$ as $m \rightarrow \infty$, we have

$$(i) \quad \|E(X_n | \mathcal{F}_{n-m})\|_r \leq \psi_m c_n, \quad \text{and} \quad (ii) \quad \|X_n - E(X_n | \mathcal{F}_{n-m})\|_r \leq \psi_{m+1} c_n,$$

for all $n \geq 1$ and $m \geq 0$. The idea of mixingale is to try to build a bridge between martingale and mixing. The following examples give the idea of the scope of L_2 -mixingales.

Examples:

1. A square-integrable martingale is a mixingale with $c_n = \|X_n\|$ and $\psi_0 = 1$ and $\psi_m = 0$ for $m \geq 1$.

2. A linear process is given by $X_n = \sum_{i=-\infty}^{\infty} \alpha_{i-n} \xi_i$ with $\{\xi_i\}$ iid mean zero and variance σ^2 and $\sum_{i=-\infty}^{\infty} \alpha_i^2 < \infty$. Then, $\{X_n, \mathcal{F}_n\}$ is a mixingale with all $c_n = \sigma$ and $\psi_m^2 = \sum_{|i| \geq m} \alpha_i^2$.

3. If $\{X_n\}$ is a square-integrable sequence of ϕ -mixing, then it is a mixingale with $c_n = 2\|X_n\|_2$ and $\psi_m = \phi^{1/2}(m)$, where $\phi(m)$ is the ϕ -mixing coefficient.

4. If $\{X_n\}$ is a sequence of α -mixing with $\|X_n\|_p < \infty$ for some $p > 2$, then it is a mixingale with $c_n = 2(\sqrt{2} + 1)\|X_n\|_2$ and $\psi_m = \alpha^{1/2-1/p}(m)$, where $\alpha(m)$ is the α -mixing coefficient.

Note that Examples 3 and 4 can be derived from the following inequality, due to McLeish (1975).

Lemma 3.3: (McLeish's inequality) Suppose that X is a random variable measurable with respect to \mathcal{A} , and $\|X\|_r < \infty$ for some $1 \leq p \leq r \leq \infty$. Then

$$\|E(X | \mathcal{F}) - E(X)\|_p \leq \begin{cases} 2[\phi(\mathcal{F}, \mathcal{A})]^{1-1/r} \|X\|_r, & \text{for } \phi\text{-mixing,} \\ 2(2^{1/p} + 1)[\alpha(\mathcal{F}, \mathcal{A})]^{1/p-1/r} \|X\|_r, & \text{for } \alpha\text{-mixing.} \end{cases}$$

3.2 Nonparametric Density Estimate

Let $\{X_i\}$ be a random sample with a (unknown) marginal distribution $F(\cdot)$ (CDF) and its probability density function (PDF) $f(\cdot)$. The question is how to estimate $f(\cdot)$ and $F(\cdot)$. Since

$$F(x) = P(X_i \leq x) = E[I(X_i \leq x)] = \int_{-\infty}^x f(u) du,$$

and

$$f(x) = \lim_{h \downarrow 0} \frac{F(x+h) - F(x-h)}{2h} \approx \frac{F(x+h) - F(x-h)}{2h}$$

if h is very small, by the method of moment estimation (MME), $F(x)$ can be estimated by

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x),$$

which is called the **empirical cumulative distribution function** (ecdf), so that $f(x)$ can be estimated by

$$f_n(x) = \frac{F_n(x+h) - F_n(x-h)}{2h} = \frac{1}{n} \sum_{i=1}^n K_h(X_i - x),$$

where $K(u) = I(|u| \leq 1)/2$ and $K_h(u) = K(u/h)/h$. Indeed, the **kernel** function $K(u)$ can be taken to be any **symmetric** density function. here, h is called the **bandwidth**. $f_n(x)$ was proposed initially by Rosenblatt (1956) and Parzen (1962) explored its properties in detail. Therefore, it is called the **Rosenblatt-Parzen density estimate**.

Exercise: Please show that $F_n(x)$ is an unbiased estimate of $F(x)$ but $f_n(x)$ is a biased estimate of $f(x)$. **Think about intuitively**

- (1) why $f_n(x)$ is biased
- (2) where the bias comes from
- (3) why $K(\cdot)$ should be symmetric.

3.2.1 Asymptotic Properties

Asymptotic Properties for ECDF

If $\{X_i\}$ is stationary, then $E[F_n(x)] = 0$ and

$$\begin{aligned} n \text{Var}(F_n(x)) &= \text{Var}(I(X_1 \leq x)) + 2 \sum_{i=2}^n \left(1 - \frac{i-1}{n}\right) \text{Cov}(I(X_1 \leq x), I(X_i \leq x)) \\ &= \underbrace{F(x)[1 - F(x)] + 2 \sum_{i=2}^n \text{Cov}(I(X_1 \leq x), I(X_i \leq x))}_{\rightarrow \sigma^2(x) \text{ by assuming that } \sigma^2(x) < \infty} \\ &\quad - 2 \underbrace{\sum_{i=2}^n \frac{i-1}{n} \text{Cov}(I(X_1 \leq x), I(X_i \leq x))}_{\rightarrow 0 \text{ by Kronecker Lemma}} \\ &\rightarrow \sigma_F^2(x) \equiv F(x)[1 - F(x)] + 2 \underbrace{\sum_{i=2}^{\infty} \text{Cov}(I(X_1 \leq x), I(X_i \leq x))}_{\text{This term is called } A_d}. \end{aligned}$$

Therefore,

$$n \text{Var}(F_n(x)) \rightarrow \sigma_F^2(x). \quad (3.1)$$

One can show based on the mixing theory that

$$\sqrt{n} [F_n(x) - F(x)] \rightarrow N(0, \sigma_F^2(x)). \quad (3.2)$$

It is clear that $A_d = 0$ if $\{X_i\}$ are independent. If $A_d \neq 0$, the **question** is how to estimate it. We can use the **HC estimator** by White (1980) or the **HAC estimator** by Newey and West (1987); see Chapter 2, or the **kernel method** by Andrew (1991).

The results in (3.2) can be used to construct a test statistic to test the null hypothesis

$$H_0 : F(x) = F_0(x) \quad \text{versus} \quad H_a : F(x) \neq (>)(<)F_0(x).$$

This test statistic is the well-known **Kolmogorov-Smirnov test**, defined as

$$D_n = \sup_{-\infty < x < \infty} |F_n(x) - F_0(x)|$$

for the two-sided test. One can show (see Serfling (1980)) that under some regularity conditions,

$$P(\sqrt{n} D_n \leq d) \rightarrow 1 - 2 \sum_{j=1}^{\infty} (-1)^{j+1} \exp(-2j^2 d^2)$$

and

$$P(\sqrt{n} D_n^+ \leq d) = P(\sqrt{n} D_n^- \geq -d) \rightarrow 1 - \exp(-2d^2),$$

where $D_n^+ = \sup_{-\infty < x < \infty} [F_n(x) - F_0(x)]$ and $D_n^- = \sup_{-\infty < x < \infty} [F_0(x) - F_n(x)]$ for one-sided tests. In **R**, there is a built-in command for the Kolmogorov-Smirnov test, which is **ts.test()**.

Asymptotic Properties for Density Estimation

Next, we derive the asymptotic variance for $f_n(x)$. First, define $Z_i = K_h(X_i - x)$. Then,

$$\begin{aligned} E[Z_1 Z_i] &= \int \int K_h(u - x) K_h(v - x) f_{1,i}(u, v) du dv \\ &= \int \int K(u) K(v) f_{1,i}(x + u h, x + v h) du dv \\ &\rightarrow f_{1,i}(x, x), \end{aligned}$$

where $f_{1,i}(u, v)$ is the joint density of (X_1, X_i) , so that

$$\text{Cov}(Z_1, Z_i) \rightarrow f_{1,i}(x, x) - f^2(x).$$

It is easy to show that

$$h \text{Var}(Z_1) \rightarrow \nu_0(K) f(x),$$

where $\nu_j(K) = \int u^j K^2(u) du$. Therefore,

$$\begin{aligned} n h \text{Var}(f_n(x)) &= \text{Var}(Z_1) + 2h \underbrace{\sum_{i=2}^n \left(1 - \frac{i-1}{n}\right) \text{Cov}(Z_1, Z_i)}_{\equiv A_f \rightarrow 0 \text{ under some assumptions}} \\ &\rightarrow \nu_0(K) f(x). \end{aligned}$$

To show that $A_f \rightarrow 0$, let $d_n \rightarrow \infty$ and $d_n h \rightarrow 0$. Then,

$$|A_f| \leq h \sum_{i=2}^{d_n} |\text{Cov}(Z_1, Z_i)| + h \sum_{i=d_n+1}^n |\text{Cov}(Z_1, Z_i)|.$$

For the first term, if $f_{1,i}(u, v) \leq M_1$, then, it is bounded by $h d_n = o(1)$. For the second term, we apply the **Davydov's inequality** (see Lemma 3.1) to obtain

$$h \sum_{i=d_n+1}^n |\text{Cov}(Z_1, Z_i)| \leq M_2 \sum_{i=d_n+1}^n \alpha(i)/h = O(d_n^{-\beta+1} h^{-1})$$

if $\alpha(n) = O(n^{-\beta})$ for some $\beta > 2$. If $d_n = O(h^{-2/\beta})$, then, the second term is dominated by $O(h^{1-2/\beta})$ which goes to 0 as $n \rightarrow \infty$. Hence,

$$n h \text{Var}(f_n(x)) \rightarrow \nu_0(K) f(x). \quad (3.3)$$

By a comparison of (3.1) and (3.3), one can see clearly that there is an infinity term involved in $\sigma_F^2(x)$ due to the dependence but the asymptotic variance in (3.3) is the same as that for the iid case (without the infinity term). We can establish the following asymptotic normality for $f_n(x)$ but the proof will be discussed later.

Theorem 3.1: *Under regularity conditions, we have*

$$\sqrt{n h} \left[f_n(x) - f(x) - \frac{h^2}{2} \mu_2(K) f''(x) + o_p(h^2) \right] \rightarrow N(0, \nu_0(K) f(x)),$$

where the term $\frac{h^2}{2} \mu_2(K) f''(x)$ is called the **asymptotic bias**.

Exercise: By comparing (3.1) and (3.3), **what can you observe?**

Example 3.1: Let us examine how importance the choice of bandwidth is. The data $\{X_i\}_{i=1}^n$ are generated from $N(0, 1)$ (iid) and $n = 300$. The grid points are taken to be $[-4, 4]$ with an increment $\Delta = 0.1$. Bandwidth is taken to be 0.25, 0.5 and 1.0, respectively and the

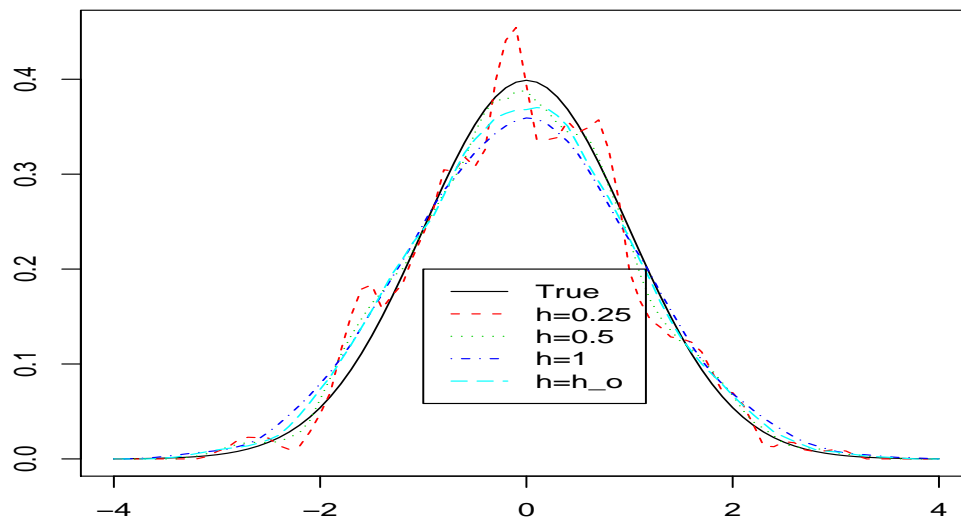


Figure 3.1: Bandwidth is taken to be 0.25, 0.5, 1.0 and the optimal one (see later) with the Epanechnikov kernel.

kernel can be the Epanechnikov kernel $K(u) = 0.75(1 - u^2)I(|u| \leq 1)$ or Gaussian kernel. Comparisons are given in Figure 3.1.

Example 3.2: Next, we apply the kernel density estimation to estimate the density of the weekly 3-month Treasury bill from January 2, 1970 to December 26, 1997. Figure 3.2 displays the estimated density together with the true standard normal density: the left panel is for the built-in function `density()` and the right panel is for own code.

Note that the computer code in **R** for the above two examples can be found in Section 3.5. **R** has a built-in function `density()` for computing the nonparametric density estimation. Also, you can use the command `plot(density())` to plot the estimated density. Further, **R** has a built-in function `ecdf()` for computing the empirical cumulative distribution function estimation and `plot(ecdf())` for plotting the step function.

3.2.2 Optimality

As we already have shown that

$$E(f_n(x)) = f(x) + \frac{h^2}{2} \mu_2(K) f''(x) + o(h^2),$$

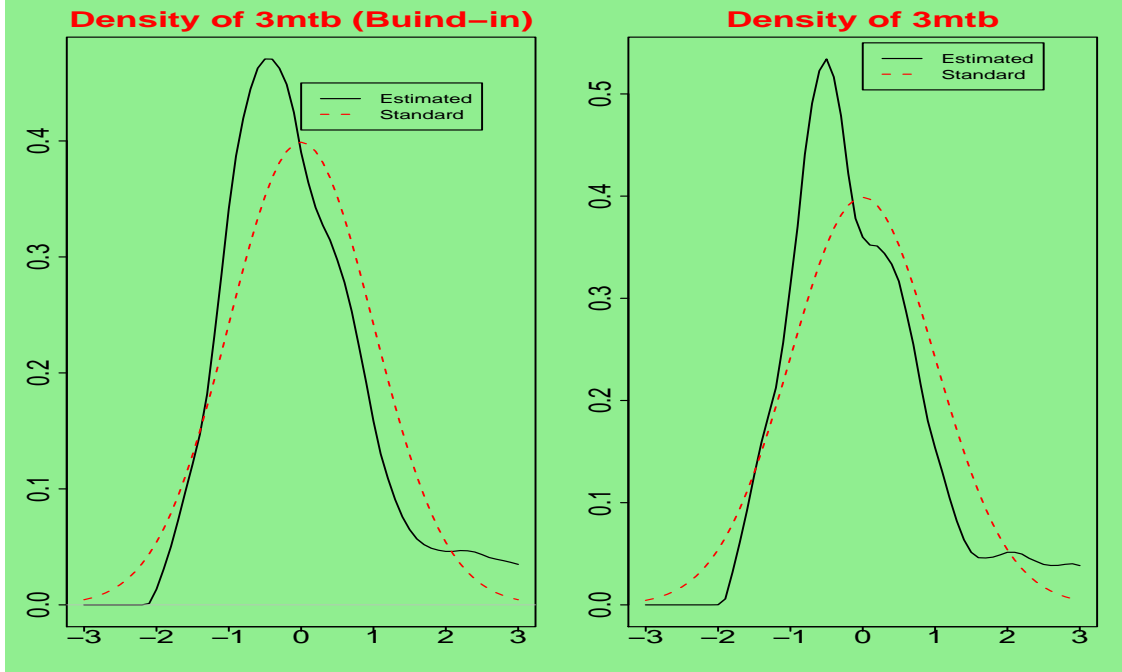


Figure 3.2: The left panel is for the built-in function **density()** and the right panel is for own code.

and

$$\text{Var}(f_n(x)) = \frac{\nu_0(K) f(x)}{n h} + o((nh)^{-1}),$$

so that the asymptotic mean integrated squares error (AMISE) is

$$\text{AMISE} = \frac{h^4}{4} \mu_2^2(K) \int [f''(x)]^2 + \frac{\nu_0(K)}{n h}.$$

Minimizing the AMISE gives the

$$h_{opt} = C_1(K) \|f''\|_2^{-2/5} n^{-1/5}, \quad (3.4)$$

where

$$C_1(K) = [\nu_0(K)/\mu_2^2(K)]^{1/5}.$$

With this asymptotically optimal bandwidth, the optimal AMISE is given by

$$\text{AMISE}_{opt} = \frac{5}{4} C_2(K) \|f''\|_2^{2/5} n^{-4/5},$$

where

$$C_2(K) = [\nu_0^2(K) \mu_2(K)]^{2/5}.$$

To choose the best kernel, it suffices to choose one to minimize $C_2(K)$.

Proposition 1: *The nonnegative probability density function K minimizing $C_2(K)$ is a re-scaling of the Epanechnikov kernel:*

$$K_{opt}(u) = \frac{3}{4a} (1 - u^2/a^2)_+$$

for any $a > 0$.

Proof: First of all, we note that $C_2(K_h) = C_2(K)$ for any $h > 0$. Let K_0 be the Epanechnikov kernel. For any other nonnegative negative K , by re-scaling if necessary, we assume that $\mu_2(K) = \mu_2(K_0)$. Thus, we need only to show that $\nu_0(K_0) \leq \nu_0(K)$. Let $G = K - K_0$. Then,

$$\int G(u)du = 0 \quad \text{and} \quad \int u^2 G(u)du = 0,$$

which implies that

$$\int (1 - u^2) G(u)du = 0.$$

Using this and the fact that K_0 has the support $[-1, 1]$, we have

$$\begin{aligned} \int G(u) K_0(u)du &= \frac{3}{4} \int_{|u| \leq 1} G(u)(1 - u^2)du \\ &= -\frac{3}{4} \int_{|u| > 1} G(u)(1 - u^2)du = \frac{3}{4} \int_{|u| > 1} K(u)(u^2 - 1)du. \end{aligned}$$

Since K is nonnegative, so is the last term. Therefore,

$$\int K^2(u)du = \int K_0^2(u)du + 2 \int K_0(u)G(u)du + \int G^2(u)du \geq \int K_0^2(u)du,$$

which proves that K_0 is the optimal kernel.

Remark: This proposition implies that the Epanechnikov kernel should be used in practice.

3.2.3 Boundary Problems

In many applications, the density $f(\cdot)$ has a bounded support. For example, the interest rate can not be less than zero and the income is always nonnegative. It is reasonable to assume that the interest rate has support $[0, 1)$. However, because a kernel density estimator spreads smoothly point masses around the observed data points, some of those near the boundary of the support are distributed outside the support of the density. Therefore, the kernel density estimator under estimates the density in the boundary regions. The problem is more severe for large bandwidth and for the left boundary where the density is high. Therefore,

some adjustments are needed. To gain some further insights, let us assume without loss of generality that the density function $f(\cdot)$ has a bounded support $[0, 1]$ and we deal with the density estimate at the left boundary. For simplicity, suppose that $K(\cdot)$ has a support $[-1, 1]$. For the left boundary point $x = ch$ ($0 \leq c < 1$), it can easily be seen that as $h \rightarrow 0$,

$$\begin{aligned} E(f_n(ch)) &= \int_{-c}^{1/h-c} f(ch + hu)K(u)du \\ &= f(0+) \mu_{0,c}(K) + h f'(0+)[c \mu_{0,c}(K) + \mu_{1,c}(K)] + o(h), \end{aligned} \quad (3.5)$$

where $f(0+) = \lim_{x \downarrow 0} f(x)$,

$$\mu_{j,c} = \int_{-c}^{\infty} u^j K(u)du, \quad \text{and} \quad \nu_{j,c}(K) = \int_{-c}^{\infty} u^j K^2(u)du.$$

Also, we can show that $\text{Var}(f_n(ch)) = O(1/nh)$. Therefore,

$$f_n(ch) = f(0+) \mu_{0,c}(K) + h f'(0+)[c \mu_{0,c}(K) + \mu_{1,c}(K)] + o_p(h).$$

Particularly, if $c = 0$ and $K(\cdot)$ is symmetric, then $E(f_n(0)) = f(0)/2 + o(1)$.

There are several methods to deal with the density estimation at boundary points. Possible approaches include the **boundary kernel** (see Gasser and Müller (1979) and Müller (1993)), **reflection** (see Schuster (1985) and Hall and Wehrly (1991)), **transformation** (see Wand, Marron and Ruppert (1991) and Marron and Ruppert (1994)) and **local polynomial fitting** (see Hjort and Jones (1996a) and Loader (1996)), and others.

Boundary Kernel

One way of choosing a boundary kernel is

$$K_{(c)}(u) = \frac{12}{(1+c)^4} (1+u) \left\{ (1-2c)u + \frac{3c^2 - 2c + 1}{2} \right\} I_{[-1,c]}.$$

Note $K_{(1)}(t) = K(t)$, the Epanechnikov kernel as defined above. Moreover, Zhang and Karunamuni (1998) have shown that this kernel is optimal in the sense of minimizing the MSE in the class of all kernels of order $(0, 2)$ with exactly one change of sign in their support. The downside to the boundary kernel is that it is not necessarily non-negative, as will be seen on densities where $f(0) = 0$.

Reflection

The reflection method is to construct the kernel density estimate based on the synthetic data $\{\pm X_t; 1 \leq t \leq n\}$ where “reflected” data are $\{-X_t; 1 \leq t \leq n\}$ and the original data are $\{X_t; 1 \leq t \leq n\}$. This results in the estimate

$$f_n(x) = \frac{1}{n} \left\{ \sum_{t=1}^n K_h(X_t - x) + \sum_{t=1}^n K_h(-X_t - x) \right\}, \quad \text{for } x \geq 0.$$

Note that when x is away from the boundary, the second term in the above is practically negligible. Hence, it only corrects the estimate in the boundary region. This estimator is twice the kernel density estimate based on the synthetic data $\{\pm X_t; 1 \leq t \leq n\}$. See Schuster (1985) and Hall and Wehrly (1991).

Transformation

The transformation method is to first transform the data by $Y_i = g(X_i)$, where $g(\cdot)$ is a given monotone increasing function, ranging from $-\infty$ to ∞ . Now apply the kernel density estimator to this transformed data set to obtain the estimate $f_n(y)$ for Y and apply the inverse transform to obtain the density of X . Therefore,

$$f_n(x) = g'(x) \frac{1}{n} \sum_{t=1}^n K_h(g(X_t) - g(x)).$$

The density at $x = 0$ corresponds to the tail density of the transformed data since $\log(0) = -\infty$, which can not usually be estimated well due to lack of the data at tails. Except at this point, the transformation method does a fairly good job. If $g(\cdot)$ is unknown in many situations, Karunamuni and Alberts (2003) suggested a parametric form and then estimated the parameter. Also, Karunamuni and Alberts (2003) considered other types of transformations.

Local Likelihood Fitting

The main idea is to consider the approximation $\log(f(X_t)) \approx P(X_t - x)$, where $P(u - x) = \sum_{j=0}^p a_j (u - x)^j$ with the localized version of log-likelihood

$$\sum_{t=1}^n \log(f(X_t)) K_h(X_t - x) - n \int K_h(u - x) f(u) du.$$

With this approximation, the local likelihood becomes

$$\mathcal{L}(a_0, \dots, d_p) = \sum_{t=1}^n P(X_t - x) K_h(X_t - x) - n \int K_h(u - x) \exp(P(u - x)) du.$$

Let $\{\hat{a}_j\}$ be the maximizer of the above local likelihood $\mathcal{L}(a_0, \dots, d_p)$. Then, the local likelihood density estimate is

$$f_n(x) = \exp(\hat{a}_0).$$

The maximizer does not exist, then $f_n(x) = 0$. See Loader (1996) and Hjort and Jones (1996a) for more details. If **R** is used for the local fit for density estimation, please use the function **density.lf()** in the package **localfit**.

Exercise: Please conduct a Monte Carol simulation to see what the boundary effects are and how the correction methods work. For example, you can consider some distribution densities with a finite support such as beta-distribution.

3.2.4 Bandwidth Selection

Simple Bandwidth Selectors

The optimal bandwidth (3.4) is not directly usable since it depends on the unknown parameter $\|f''\|_2$. When $f(x)$ is a Gaussian density with standard deviation σ , it is easy to see from (3.4) that

$$h_{opt} = (8\sqrt{\pi}/3)^{1/5} C_1(K) \sigma n^{-1/5},$$

which is called the **normal reference bandwidth selector** in literature, obtained by replacing the unknown parameter σ in the above equation by the sample standard deviation s . In particular, after calculating the constant $C_1(K)$ numerically, we have the following normal reference bandwidth selector

$$\hat{h}_{opt} = \begin{cases} 1.06 s n^{-1/5} & \text{for the Gaussian kernel} \\ 2.34 s n^{-1/5} & \text{for the Epanechnikov kernel} \end{cases}$$

Hjort and Jones (1996b) proposed an improved rule obtained by using an **Edgeworth expansion** for $f(x)$ around the Gaussian density. Such a rule is given by

$$\hat{h}_{opt}^* = h_{opt} \left(1 + \frac{35}{48} \hat{\gamma}_4 + \frac{35}{32} \hat{\gamma}_3^2 + \frac{385}{1024} \hat{\gamma}_4^2 \right)^{-1/5},$$

where $\hat{\gamma}_3$ and $\hat{\gamma}_4$ are respectively the sample skewness and kurtosis. For details about the **Edgeworth expansion**, please see the book by Hall (1992).

Note that the normal reference bandwidth selector is only a simple rule of thumb. It is a good selector when the data are nearly Gaussian distributed, and is often reasonable in many applications. However, it can lead to over-smooth when the underlying distribution is asymmetric or multi-modal. In that case, one can either subjectively tune the bandwidth, or select the bandwidth by more sophisticated bandwidth selectors. One can also transform data first to make their distribution closer to normal, then estimate the density using the normal reference bandwidth selector and apply the inverse transform to obtain an estimated density for the original data. Such a method is called the transformation method. There are quite a few important techniques for selecting the bandwidth such as **cross-validation (CV) and plug-in bandwidth selectors**. A conceptually simple technique, with theoretical justification and good empirical performance, is the plug-in technique. This technique relies on finding an estimate of the functional $||f''||_2$, which can be obtained by using a pilot bandwidth. An implementation of this approach is proposed by Sheather and Jones (1991) and an overview on the progress of bandwidth selection can be found in Jones, Marron and Sheather (1996).

Function **dpik()** in the package **KernSmooth** in **R** selects a bandwidth for estimating the kernel density estimation using the plug-in method.

Cross-Validation Method

The integrated squared error (ISE) of $f_n(x)$ is defined by

$$\text{ISE}(h) = \int [f_n(x) - f(x)]^2 dx.$$

A commonly used measure of discrepancy between $f_n(x)$ and $f(x)$ is the mean integrated squared error (MISE) $\text{MISE}(h) = E[\text{ISE}(h)]$. It can be shown easily (or see Chiu, 1991) that $\text{MISE}(h) \approx \text{AMISE}(h)$. The optimal bandwidth minimizing the AMISE is given in (3.4). The least squares cross-validation (LSCV) method proposed by Rudemo (1982) and Bowman (1984) is a popular method to estimate the optimal bandwidth h_{opt} . Cross-validation is very useful for assessing the performance of an estimator via estimating its prediction error. The basic idea is to set one of the data point aside for validation of a model and use the remaining data to build the model. The main idea is to choose h to minimize $\text{ISE}(h)$. Since

$$\text{ISE}(h) = \int f_n^2(x) dx - 2 \int f(x) f_n(x) dx + \int f^2(x) dx,$$

the question is how to estimate the second term on the right hand side. Well, let us consider the simplest case when $\{X_t\}$ are iid. Re-express $f_n(x)$ as

$$f_n(x) = \frac{n-1}{n} f_n^{(-s)}(x) + \frac{1}{n} K_h(X_s - x)$$

for any $1 \leq s \leq n$, where

$$f_n^{(-s)}(x) = \frac{1}{n-1} \sum_{t \neq s}^n K_h(X_t - x),$$

which is the kernel density estimate without the s th observation, commonly called the **jack-knife** estimate or leave-one-out estimate. It is easy to see that for any $1 \leq s \leq n$,

$$f_n(x) \approx f_n^{(-s)}(x).$$

Let $\mathcal{D}_s = \{X_1, \dots, X_{s-1}, X_{s+1}, \dots, X_n\}$. Then,

$$E[f_n^{(-s)}(X_s) | \mathcal{D}_s] = \int f_n^{(-s)}(x) f(x) dx \approx \int f_n(x) f(x) dx,$$

which, by using the method of moment, can be estimated by $\frac{1}{n} \sum_{s=1}^n f_n^{(-s)}(X_s)$. Therefore, the cross-validation is

$$\begin{aligned} \text{CV}(h) &= \int f_n^2(x) dx - \frac{2}{n} \sum_{s=1}^n f_n^{(-s)}(X_s) \\ &= \frac{1}{n^2} \sum_{s,t} K_h^*(X_s - X_t) - \frac{2}{n(n-1)} \sum_{t \neq s}^n K_h(X_s - X_t), \end{aligned}$$

where $K_h^*(\cdot)$ is the convolution of $K_h(\cdot)$ and $K_h(\cdot)$ as

$$K_h^*(u) = \int K_h(v) K_h(u - v) dv.$$

Let \hat{h}_{cv} be the minimizer of $\text{CV}(h)$. Then, it is called the optimal bandwidth based on the cross-validation. Stone (1984) showed that \hat{h}_{cv} is a consistent estimate of the optimal bandwidth h_{opt} .

Function **lscv()** in the package **locfit** in **R** selects a bandwidth for estimating the kernel density estimation using the least squares cross-validation method.

3.2.5 Project for Density Estimation

I. Do Monte Carlo simulations to compare the performances of the kernel density estimations for different settings and to make your own conclusions based on your simulations. Please do the followings:

1. Use the Rosenblatt-Parzen method by choosing different **sample sizes** (you take several different sample sizes, say 250, 400, 600 and 1000), **different kernels** (say the normal and Epanechnikov kernel), **different bandwidths**, and **different bandwidth selection methods** such as cross-validation and plug-in as well as normal reference. Any conclusions and comments?
2. Compare the Rosenblatt-Parzen method with **local density method** as in Loader (1996) or Hjort and Jones (1996). Any conclusions and comments?
3. Compare the various methods for boundary correction.

To assess the performance of finite samples, for each setting, you need to compute the mean absolute deviation errors (MADE) for $\hat{f}(\cdot)$, defined as

$$\text{MADE} = n_0^{-1} \sum_{k=1}^{n_0} \left| \hat{f}(u_k) - f(u_k) \right|,$$

where $\hat{f}(\cdot)$ is the nonparametric estimate of $f(\cdot)$ and $\{u_k\}$ are the grid points, taken to be arbitrary within the range of data. Note that you can choose any distribution to generate your samples for your simulation. Also, note that the choice of the grid points is not important so that they can be chosen arbitrarily. In general, the number of replications can be taken to be $n_{sim} = 500$ or 1000. The question is how to report the simulation results. There are two ways of doing so. You can display the n_{sim} values of MADE either in a boxplot form (**boxplot() in R**) or in a table by presenting the *median* and *standard deviation* of the n_{sim} values of MADE. Either one is okay but the boxplot is preferred by most people.

II. Consider three real data sets for the US Treasury bill (Secondary Market Rate): the **daily** 3-month Treasury bill from January 4, 1954 to May 2, 2007, in the data file *DTB3.txt* or *DTB3.csv*, the **weekly** 3-month Treasury bill from January 8, 1954 to

April 27, 2007, in the data file *WTB3MS.txt* or *WTB3MS.csv*, and the **monthly** 3-month Treasury bill from January 1, 1934 to March 1, 2007, in the data file *TB3MS.txt* or *TB3MS.csv*.

1. Apply Ljung-Box test [**Box.test()** in **R**] to see if three series are autocorrelated or not. Also, you might look at the autocorrelation function (ACF) [**acf()** in **R**] or/and partial autocorrelation function (PACF) [**pacf()** in **R**].
2. Apply the kernel density estimation to estimate three density functions.
3. Any conclusions and comments on three density functions?

♡ Note that the real data sets can be downloaded from the web site for Federal Reserve Bank of Saint Louis at <http://research.stlouisfed.org/fred2/categories/46>. You can use any statistical package to do your simulation. You try to use **R** since it is very simple. You need to hand in all necessary materials (tables or graphs) to support your conclusions. If you need any help, please come to see me.

3.2.6 Multivariate Density Estimation

As we discussed earlier, the kernel density or distribution estimation is basically one-dimensional. For multivariate case, the kernel density estimate is given by

$$f_n(x) = \frac{1}{n} \sum_{t=1}^n K_H(X_t - x), \quad (3.6)$$

where $K_H(u) = K(H^{-1}u)/\det(H)$, $K(u)$ is a multivariate kernel function, and H is the bandwidth matrix such as for all $1 \leq i, j \leq p$, $n h_{ij} \rightarrow \infty$ and $h_{ij} \rightarrow 0$ where h_{ij} is the (i, j) th element of H . The bandwidth matrix is introduced to capture the dependent structure in the independent variables. Particularly, if H is a diagonal matrix and $K(u) = \prod_{j=1}^p K_j(u_j)$ where $K_j(\cdot)$ is a univariate kernel function, then, $f_n(x)$ becomes

$$f_n(x) = \frac{1}{n} \sum_{t=1}^n \prod_{j=1}^p K_{h_j}(X_{jt} - x_j),$$

which is called the **product kernel density estimation**. This case is commonly used in practice. Similar to the univariate case, it is easy to derive the theoretical results for the multivariate case, **which is left as an exercise**. See Wand and Jones (1995) for details.

Table 3.1: Sample sizes required for p -dimensional nonparametric regression to have comparable performance with that of 1-dimensional nonparametric regression using size 100

dimension	2	3	4	5	6	7	8	9	10
sample size	252	631	1,585	3,982	10,000	25,119	63,096	158,490	398,108

Curse of Dimensionality

For the product kernel estimate with $h_j = h$, we can show easily that

$$E(f_n(x)) = f(x) + \frac{h^2}{2} \text{tr}(\mu_2(K) f''(x)) + o(h^2),$$

where $\mu_2(K) = \int u u^T K(u) du$, and

$$\text{Var}(f_n(x)) = \frac{\nu_0(K) f(x)}{n h^p} + o((nh)^{-1}),$$

so that the AMSE is given by

$$\text{AMSE} = \frac{\nu_0(K) f(x)}{n h^p} + \frac{h^4}{4} B(x),$$

where $B(x) = (\text{tr}(\mu_2(K) f''(x)))^2$. By minimizing the AMSE, we obtain the optimal bandwidth

$$h_{opt} = \left(\frac{p \nu_0(K) f(x)}{B(x)} \right)^{1/(p+4)} n^{-1/(p+4)},$$

which leads to the optimal rate of convergence for MSE which is $O(n^{-4/(4+p)})$ by trading off the rates between the bias and variance. When p is large, the so called **“curse of dimensionality”** exists. To understand this problem quantitatively, let us look at the rate of convergence. To have a comparable performance with one-dimensional nonparametric regression with n_1 data points, for p -dimensional nonparametric regression, we need the number of data points n_p ,

$$O(n_p^{-4/(4+p)}) = O(n_1^{-4/5}),$$

or $n_p = O(n_1^{(p+4)/5})$. Note that here we only emphasize on the rate of convergence for MSE by ignoring the constant part. Table 3.1 shows the result with $n_1 = 100$. The increase of required sample sizes is exponentially fast.

Exercise: Please derive the asymptotic results given in (3.6) for the general multivariate case.

In **R**, the built-in function **density()** is only for univariate case. For multivariate situations, there are two packages **ks** and **KernSmooth**. Function **kde()** in **ks** can compute the multivariate density estimate for 2- to 6- dimensional data and Function **bkde2D()** in **KernSmooth** computes the 2D kernel density estimate. Also, **ks** provides some functions for some bandwidth matrix selection such as **Hbcv()** and **Hscv** for 2D case and **Hlscv()** and **Hpi()**.

3.2.7 Reading Materials

Applications in Finance: Please read the papers by Aït-Sahalia and Lo (1998, 2000), Pritsker (1998) and Hong and Li (2005) on how to apply the kernel density estimation to the nonparametric estimation of the state-price densities (SPD) or risk neutral densities (RND) and nonparametric risk estimation based on the state-price density. Please download the data from <http://finance.yahoo.com/> (say, S&P500 index) to estimate the SPD.

3.3 Distribution Estimation

3.3.1 Smoothed Distribution Estimation

The question is how to obtain a smoothed estimate of CDF $F(x)$. Well, one way of doing so is to integrate the estimated PDF $f_n(x)$, given by

$$\hat{F}_n(x) = \int_{-\infty}^x f_n(u) du = \frac{1}{n} \sum_{i=1}^n \mathcal{K}\left(\frac{x - X_i}{h}\right),$$

where $\mathcal{K}(x) = \int_{-\infty}^x K(u) du$; the distribution of $K(\cdot)$. **Why do we need this smoothed estimate of CDF?** To answer this question, we need to consider the **mean squares error (MSE)**.

First, we derive the asymptotic bias. By the integration by parts, we have

$$\begin{aligned} E\left[\hat{F}_n(x)\right] &= E\left[\mathcal{K}\left(\frac{x - X_i}{h}\right)\right] = \int F(x - hu)K(u)du \\ &= F(x) + \frac{h^2}{2} \mu_2(K) f'(x) + o(h^2). \end{aligned}$$

Next, we derive the asymptotic variance.

$$E\left[\mathcal{K}^2\left(\frac{x - X_i}{h}\right)\right] = \int F(x - hu)b(u)du = F(x) - h f(x) \theta + o(h),$$

where $b(u) = 2 K(u) \mathcal{K}(u)$ and $\theta = \int u b(u) du$. Then,

$$\text{Var} \left[\mathcal{K} \left(\frac{x - X_i}{h} \right) \right] = F(x)[1 - F(x)] - h f(x) \theta + o(h).$$

Define $I_j(x) = \text{Cov}(I(X_1 \leq x), I(X_{j+1} \leq t)) = F_j(x, x) - F^2(x)$ and

$$I_{nj}(x) = \text{Cov} \left(\mathcal{K} \left(\frac{x - X_1}{h} \right), \mathcal{K} \left(\frac{x - X_{j+1}}{h} \right) \right).$$

By means of Lemma 2 in Lehmann (1966), the covariance $I_{nj}(x)$ may be written as follows

$$\begin{aligned} I_{nj}(t) = & \int \left\{ P \left[\mathcal{K} \left(\frac{x - X_1}{h} \right) > u, \mathcal{K} \left(\frac{x - X_{j+1}}{h} \right) > v \right] \right. \\ & \left. - P \left[\mathcal{K} \left(\frac{x - X_1}{h} \right) > u \right] P \left[\mathcal{K} \left(\frac{x - X_{j+1}}{h} \right) > v \right] \right\} dudv. \end{aligned}$$

Inverting the CDF $\mathcal{K}(\cdot)$ and making two changes of variables, the above relation becomes

$$I_{nj}(x) = \int [F_j(x - hu, x - hv) - F(x - hu)F(x - hv)] K(u)K(v) dudv.$$

Expanding the right-hand side of the above equation according to Taylor's formula, we obtain

$$|I_{nj}(x) - I_j(x)| \leq C h^2.$$

By the Davydov's inequality (see Lemma 3.1), we have

$$|I_{nj}(x) - I_j(x)| \leq C \alpha(j),$$

so that for any $1/2 < \tau < 1$,

$$|I_{nj}(x) - I_j(x)| \leq C h^{2\tau} \alpha^{1-\tau}(j).$$

Therefore,

$$\frac{1}{n} \sum_{j=1}^{n-1} (n-j) |I_{nj}(x) - I_j(x)| \leq \sum_{j=1}^{n-1} |I_{nj}(x) - I_j(x)| \leq C h^{2\tau} \sum_{j=1}^{\infty} \alpha^{1-\tau}(j) = O(h^{2\tau})$$

provided that $\sum_{j=1}^{\infty} \alpha^{1-\tau}(j) < \infty$ for some $1/2 < \tau < 1$. Indeed, this assumption is satisfied if $\alpha(n) = O(n^{-\beta})$ for some $\beta > 2$. By the stationarity, it is clear that

$$n \text{Var} \left(\widehat{F}_n(x) \right) = \text{Var} \left(\mathcal{K} \left(\frac{x - X_1}{h} \right) \right) + \frac{2}{n} \sum_{j=1}^{n-1} (n-j) I_{nj}(x).$$

Therefore,

$$\begin{aligned} n \operatorname{Var} \left(\widehat{F}_n(x) \right) &= F(x)[1 - F(x)] - h f(x) \theta + o(h) + 2 \sum_{j=1}^{\infty} I_j(x) + O(h^{2\tau}) \\ &= \sigma_F^2(x) - h f(x) \theta + o(h). \end{aligned}$$

We can establish the following asymptotic normality for $\widehat{F}_n(x)$ but the proof will be discussed later.

Theorem 3.2: *Under regularity conditions, we have*

$$\sqrt{n} \left[\widehat{F}_n(x) - F(x) - \frac{h^2}{2} \mu_2(K) f'(x) + o_p(h^2) \right] \rightarrow N(0, \sigma_F^2(x)).$$

Similarly, we have

$$n \operatorname{AMSE} \left(\widehat{F}_n(x) \right) = \frac{n h^4}{4} \mu_2^2(K) [f'(x)]^2 + \sigma_F^2(x) - h f(x) \theta.$$

If $\theta > 0$, minimizing the AMSE gives the

$$h_{opt} = \left(\frac{\theta f(x)}{\mu_2^2(K) [f'(x)]^2} \right)^{1/3} n^{-1/3},$$

and with this asymptotically optimal bandwidth, the optimal AMSE is given by

$$n \operatorname{AMSE}_{opt} \left(\widehat{F}_n(x) \right) = \sigma_F^2(x) - \frac{3}{4} \left(\frac{\theta^2 f^2(x)}{\mu_2(K) f'(x)} \right)^{2/3} n^{-1/3}.$$

Remark: From the aforementioned equation, we can see that if $\theta > 0$, the AMSE of $\widehat{F}_n(x)$ can be smaller than that for $F_n(x)$ in the second order. Also, it is easy to see that if $K(\cdot)$ is the Epanechnikov kernel, $\theta > 0$.

3.3.2 Relative Efficiency and Deficiency

To measure the **relative efficiency and deficiency** of $\widehat{F}_n(x)$ over $F_n(x)$, we define

$$i(n) = \min \left\{ k \in \{1, 2, \dots\}; \operatorname{MSE}(F_k(x)) \leq \operatorname{MSE} \left(\widehat{F}_n(x) \right) \right\}.$$

We have the following results without the detailed proofs which can be found in Cai and Roussas (1998).

Proposition 2: (i) Under regularity conditions,

$$\frac{i(n)}{n} \rightarrow 1, \quad \text{if and only if} \quad nh_n^4 \rightarrow 0.$$

(ii) Under regularity conditions,

$$\frac{i(n) - n}{nh} \rightarrow \theta(x), \quad \text{if and only if} \quad nh_n^3 \rightarrow 0,$$

where $\theta(x) = f(x)\theta/\sigma_F^2(x)$.

Remark: It is clear that the quantity $\theta(x)$ may be looked upon as a way of measuring the performance of the estimate $\hat{F}_n(x)$. Suppose that the kernel $K(\cdot)$ is chosen, so that $\theta > 0$, which is equivalent to $\theta(x) > 0$. Then, for sufficiently large n , $i(n) > n + nh_n(\theta(x) - \varepsilon)$. Thus, $i(n)$ is substantially larger than n , and, indeed, $i(n) - n$ tends to ∞ . Actually, Reiss (1981) and Falk (1983) posed the question of determining the exact value of the superiority of θ over a certain class of kernels. More specifically, let \mathcal{K}_m be the class of kernels $\mathcal{K} : [-1, 1] \rightarrow \mathfrak{R}$ which are absolutely continuous and satisfy the requirements: $\mathcal{K}(-1) = 0$, $\mathcal{K}(1) = 1$, and $\int_{-1}^1 u^\mu K(u) du = 0$, $\mu = 1, \dots, m$, for some $m = 0, 1, \dots$ (where the moment condition is vacuous for $m = 0$). Set $\Psi_m = \sup\{\theta; \mathcal{K} \in \mathcal{K}_m\}$. Then, Mammitzsch (1984) answered the question posed by showing in an elegant manner. See Cai and Roussas (1998) for more details and simulation results.

Exercise: Please conduct a Monte Carol simulation to see what the differences are for smoothed and non-smoothed distribution estimations.

3.4 Quantile Estimation

Let $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ denote the order statistics of $\{X_t\}_{t=1}^n$. Define the inverse of $F(x)$ as $F^{-1}(p) = \inf\{x \in \mathfrak{R}; F(x) \geq p\}$, where \mathfrak{R} is the real line. The traditional estimate of $F(x)$ has been the empirical distribution function $F_n(x)$ based on X_1, \dots, X_n , while the estimate of the p -th quantile $\xi_p = F^{-1}(p)$, $0 < p < 1$, is the sample quantile function $\xi_{pn} = F_n^{-1}(p) = X_{([np])}$, where $[x]$ denotes the integer part of x . It is a consistent estimator of ξ_p for α -mixing data (Yoshihara, 1995). However, as stated in Falk (1983), $F_n(x)$ does not take into account the smoothness of $F(x)$; i.e., the existence of a probability density function $f(x)$. In order to incorporate this characteristic, investigators proposed several smoothed quantile estimates, one of which is based on $\hat{F}_n(x)$ obtained as a convolution between $F_n(x)$

and a properly scaled kernel function; see the previous section. Finally, note that **R** has a command **quantile()** which can be used for computing ξ_{pn} , the nonparametric estimate of quantile.

3.4.1 Value at Risk

Value at Risk (VaR) is a popular measure of market risk associated with an asset or a portfolio of assets. It has been chosen by the Basel Committee on Banking Supervision as a benchmark risk measure and has been used by financial institutions for asset management and minimization of risk. Let $\{X_t\}_{t=1}^n$ be the market value of an asset over n periods of $t = 1$ a time unit, and let $Y_t = -\log(X_t/X_{t-1})$ be the **negative log-returns (loss)**. Suppose $\{Y_t\}_{j=1}^n$ is a strictly stationary dependent process with marginal distribution function $F(y)$. Given a positive value p close to zero, the $1 - p$ level VaR is

$$\nu_p = \inf\{u : F(u) \geq 1 - p\} = F^{-1}(1 - p),$$

which specifies the smallest amount of loss such that the probability of the loss in market value being larger than ν_p is less than p . Comprehensive discussions on VaR are available in Duffie and Pan (1997) and Jorion (2001), and references therein. Therefore, VaR can be regarded as a special case of quantile. **R** has a built-in package called **VaR** for a set of methods for calculation of VaR, particularly, for some parametric models such as the **General Pareto Distribution (GPD)**. But the restrict parametric specifications might be misspecified.

A more general form for the **generalized Pareto distribution** with shape parameter $k \neq 0$, scale parameter σ , and threshold parameter θ , is

$$f(x) = \frac{1}{\sigma} \left(1 + k \frac{x - \theta}{\sigma}\right)^{-1/k-1}, \quad \text{and} \quad F(x) = 1 - \left(1 + k \frac{x - \theta}{\sigma}\right)^{-1/k}$$

for $\theta < x$, when $k > 0$. In the limit for $k = 0$, the density is $f(x) = \frac{1}{\sigma} \exp(-(x - \theta)/\sigma)$ for $\theta < x$. If $k = 0$ and $\theta = 0$, the generalized Pareto distribution is equivalent to the exponential distribution. If $k > 0$ and $\theta = \sigma$, the generalized Pareto distribution is equivalent to the Pareto distribution.

Another popular risk measure is the **expected shortfall** (ES) which is the expected loss, given that the loss is at least as large as some given quantile of the loss distribution (e.g.,

VaR), defined as

$$\mu_p = E(Y_t | Y_t > \nu_p) = \int_{\nu_p}^{\infty} y f(y) dy / p.$$

It is well known from Artzner, Delbaen, Eber and Heath (1999) that ES is a **coherent risk** measure such as it satisfies the **four axioms**: **homogeneity** (increasing the size of a portfolio by a factor should scale its risk measure by the same factor), **monotonicity** (a portfolio must have greater risk if it has systematically lower values than another), **risk-free condition or translation invariance** (adding some amount of cash to a portfolio should reduce its risk by the same amount), and **subadditivity** (the risk of a portfolio must be less than the sum of separate risks or merging portfolios cannot increase risk). VaR satisfies homogeneity, monotonicity, and risk-free condition but is not sub-additive. See Artzner, *et al.* (1999) for details.

3.4.2 Nonparametric Quantile Estimation

The smoothed sample quantile estimate of ξ_p , $\hat{\xi}_p$, based on $\hat{F}_n(x)$, is defined by:

$$\hat{\xi}_p = \hat{F}_n^{-1}(1 - p) = \inf \left\{ x \in \mathbb{R}; \hat{F}_n(x) \geq 1 - p \right\}.$$

$\hat{\xi}_p$ is referred to in literature as the perturbed (smoothed) sample quantile. Asymptotic properties of $\hat{\xi}_p$, both under independence as well as under certain modes of dependence, have been investigated extensively in literature; see Cai and Roussas (1997) and Chen and Tang (2005).

By the differentiability of $\hat{F}_n(x)$, we use the Taylor expansion and ignore the higher terms to obtain

$$\hat{F}_n(\hat{\xi}_p) = 1 - p \approx \hat{F}_n(\xi_p) - f_n(\xi_p) (\hat{\xi}_p - \xi_p), \quad (3.7)$$

then,

$$\hat{\xi}_p - \xi_p \approx [\hat{F}_n(\xi_p) - (1 - p)] / f_n(\xi_p) \approx [\hat{F}_n(\xi_p) - (1 - p)] / f(\xi_p)$$

since $f_n(x)$ is a consistent estimator of $f(x)$. As an application of Theorem 3.2, we can establish the following theorem for the asymptotic normality of $\hat{\xi}_p$ but the proof is omitted since it is similar to that for Theorem 3.2.

Theorem 3.3: *Under regularity conditions, we have*

$$\sqrt{n} \left[\hat{\xi}_p - \xi_p - \frac{h^2}{2} \mu_2(K) f'(\xi_p) / f(\xi_p) + o_p(h^2) \right] \rightarrow N(0, \sigma_F^2(\xi_p) / f^2(\xi_p)).$$

Next, let us examine the AMSE. To this effect, we can derive the asymptotic bias and variance. From the previous section, we have

$$E \left[\widehat{\xi}_p \right] = \xi_p + \frac{h^2}{2} \mu_2(K) f'(\xi_p)/f(\xi_p) + o_p(h^2),$$

and

$$n \operatorname{Var} \left[\widehat{\xi}_p \right] = \sigma_F^2(\xi_p)/f^2(\xi_p) - h \theta / f(\xi_p) + o(h).$$

Therefore, the AMSE is

$$n \operatorname{AMSE}(\widehat{\xi}_p) = \frac{n h^4}{4} \mu_2^2(K) [f'(\xi_p)/f(\xi_p)]^2 + \sigma_F^2(\xi_p)/f^2(\xi_p) - h \theta / f(\xi_p).$$

If $\theta > 0$, minimizing the AMSE gives the

$$h_{opt} = \left(\frac{\theta f(\xi_p)}{\mu_2^2(K) [f'(\xi_p)]^2} \right)^{1/3} n^{-1/3},$$

and with this asymptotically optimal bandwidth, the optimal AMSE is given by

$$n \operatorname{AMSE}_{opt}(\widehat{\xi}_p) = \sigma_F^2(\xi_p)/f^2(\xi_p) - \frac{3}{4} \left(\frac{\theta^2}{\mu_2(K) f'(\xi_p) f(\xi_p)} \right)^{2/3} n^{-1/3},$$

which indicates a reduction to the AMSE of the second order. Chen and Tang (2005) conducted an intensive study on simulations to demonstrate the advantages of nonparametric estimation $\widehat{\xi}_p$ over the **sample quantile** ξ_{pn} under the VaR setting. We refer to the paper by Chen and Tang (2005) for simulation results and empirical examples.

Exercise: Please use the above procedures to estimate nonparametrically the ES and discuss its properties as well as conduct simulation studies and empirical applications.

3.5 Computer Code

```
# April 10, 2007
graphics.off() # clean the previous graphs on the screen

#####

# Example 3.1

#####
```



```
#####
# Define the Epanechnikov kernel function
kernel<-function(x){0.75*(1-x^2)*(abs(x)<=1)}
#####
# Define the kernel density estimator
kernden=function(x,z,h,ker){
  # parameters: x=variable; h=bandwidth; z=grid point; ker=kernel
  nz<-length(z)
  nx<-length(x)
  x0=rep(1,nx*nz)
  dim(x0)=c(nx,nz)
  x1=t(x0)
  x0=x*x0
  x1=z*x1
  x0=x0-t(x1)
  if(ker==1){x1=kernel(x0/h)}          # Epanechnikov kernel
  if(ker==0){x1=dnorm(x0/h)}          # normal kernel
  f1=apply(x1,2,mean)/h
  return(f1)
}
#####
#####
# Simulation for different bandwidths and different kernels
n=300                                # n=300
ker=1                                # ker=1 => Epan; ker=0 => Gaussian
h0=c(0.25,0.5,1)                    # set initial bandwidths
z=seq(-4,4,by=0.1)                   # grid points
nz=length(z)                         # number of grid points
x=rnorm(n)                           # simulate  $x \sim N(0, 1)$ 
if(ker==1){h_o=2.34*n^(-0.2)}        # bandwidth for Epanechnikov kernel
if(ker==0){h_o=1.06*n^(-0.2)}        # bandwidth for normal kernel
```

```

f1=kernden(x,z,h0[1],ker)
f2=kernden(x,z,h0[2],ker)
f3=kernden(x,z,h0[3],ker)
f4=kernden(x,z,h_o,ker)
text1=c("True","h=0.25","h=0.5","h=1","h=h_o")
data=cbind(dnorm(z),f1,f2,f3,f4)      # combine them as a matrix
win.graph()
matplot(z,data,type="l",lty=1:5,col=1:5,xlab="",ylab="")
legend(-1,0.2,text1,lty=1:5,col=1:5)
#####

#####
# Example 3.2
#####
z1=read.table("c:/res-teach/xiada/teaching05-07/data/ex3-2.txt")
# dada: weekly 3-month Treasury bill from 1970 to 1997
x=z1[,4]/100                        # decimal
n=length(x)
y=diff(x)                          # Delta  $x_t = x_t - x_{t-1}$  = change rate
x=x[1:(n-1)]
n=n-1
x_star=(x-mean(x))/sqrt(var(x))    # standardized
den_3mtb=density(x_star,bw=0.30,kernel=c("epanechnikov"),
  from=-3,to=3,n=61)
den_est=den_3mtb$y                 # estimated density values
z_star=seq(-3,3,by=0.1)
text1=c("Estimated Density","Standard Norm")

win.graph()
par(bg="light green")
plot(den_3mtb,main="Density of 3mtb (Buind-in)",ylab="",xlab="",
  col.main="red")

```

```

points(z_star,dnorm(z_star),type="l",lty=2,col=2,ylab="",xlab="")
legend(0,0.45,text1,lty=c(1,2),col=c(1,2),cex=0.7)

h_den=0.5
f_hat=kernden(x_star,z_star,h_den,1)
ff=cbind(f_hat,dnorm(z_star))

win.graph()
par(bg="light blue")
matplot(z_star,ff,type="l",lty=c(1,2),col=c(1,2),ylab="",xlab="")
title(main="Density of 3mtb",col.main="red")
legend(0,0.55,text1,lty=c(1,2),col=c(1,2),cex=0.7)
#####

```

3.6 References

- Aït-Sahalia, Y. and A.W. Lo (1998). Nonparametric estimation of state-price densities implicit in financial asset prices. *Journal of Finance*, **53**, 499-547.
- Aït-Sahalia, Y. and A.W. Lo (2000). Nonparametric risk management and implied risk aversion. *Journal of Econometric*, **94**, 9-51.
- Andrews, D.W.K. (1991). Heteroskedasticity and autocorrelation consistent covariance matrix estimation. *Econometrica*, **59**, 817-858.
- Artzner, P., F. Delbaen, J.M. Eber, and D. Heath (1999). Coherent measures of risk. *Mathematical Finance*, **9**, 203-228.
- Bowman, A. (1984). An alternative method of cross-validation for the smoothing of density estimate. *Biometrika*, **71**, 353-360.
- Cai, Z. (2002). Regression quantile for time series. *Econometric Theory*, **18**, 169-192.
- Cai, Z. and G.G. Roussas (1997). Smooth estimate of quantiles under association. *Statistics and Probability Letters*, **36**, 275-287.
- Cai, Z. and G.G. Roussas (1998). Efficient estimation of a distribution function under quadrant dependence. *Scandinavian Journal of Statistics*, **25**, 211-224.
- Carrasco, M. and X. Chen (2002). Mixing and moments properties of various GARCH and stochastic volatility models. *Econometric Theory*, **18**, 17-39.

- Chen, S.X. and C.Y. Tang (2005). Nonparametric inference of value at risk for dependent financial returns. *Journal of Financial Econometrics*, **3**, 227-255.
- Chiu, S.T. (1991). Bandwidth selection for kernel density estimation. *The Annals of Statistics*, **19**, 1883-1905.
- Duffie, D. and J. Pan (1997). An overview of value at risk. *Journal of Derivatives*, **4**, 7-49.
- Fan, J. and Q. Yao (2003). *Nonlinear Time Series: Nonparametric and Parametric Methods*. Springer-Verlag, New York.
- Gasser, T. and H.-G. Müller (1979). Kernel estimation of regression functions. In *Smoothing Techniques for Curve Estimation*, Lecture Notes in Mathematics, **757**, 23-68. Springer-Verlag, New York.
- Falk, M.(1983). Relative efficiency and deficiency of kernel type estimators of smooth distribution functions. *Statistica Neerlandica*, **37**, 73-83.
- Genon-Caralot, V., T. Jeantheau and C. Laredo (2000). Stochastic volatility models as hidden Markov models and statistical applications. *Bernoulli*, **6**, 1051-1079.
- Hall, P. (1992). *The Bootstrap and Edgeworth Expansion*. Springer-Verlag, New York.
- Hall, P. and C.C. Heyde (1980). *Martingale Limit Theory and its Applications*. Academic Press, New York.
- Hall, P. and T.E. Wehrly (1991). A geometrical method for removing edge effects from kernel-type nonparametric regression estimators. *Journal of American Statistical Association*, **86**, 665-672.
- Hjort, N.L. and M.C. Jones (1996a). Locally parametric nonparametric density estimation. *The Annals of Statistics*, **24**, 1619-1647.
- Hjort, N.L. and M.C. Jones (1996b). Better rules of thumb for choosing bandwidth in density estimation. *Working paper*, Department of Mathematics, University of Oslo, Norway.
- Hong, Y. and H. Li (2005). Nonparametric specification testing for continuous-time models with applications to interest rate term structures. *Review of Financial Studies*, **18**, 37-84.
- Jones, M.C., J.S. Marron and S.J. Sheather (1996). A brief survey of bandwidth selection for density estimation. *Journal of American Statistical Association*, **91**, 401-407.
- Jorion, P. (2001). *Value at Risk*, 2nd Edition. New York: McGraw-Hill.
- Karunamuni, R.J. and T. Alberts (2003). On boundary correction in kernel density estimation. *Working paper*, Department of Mathematical and Statistical Sciences, University of Alberta, Canada.

- Lehmann, E. (1966). Some concepts of dependence. *Annals of Mathematical Statistics*, **37**, 1137-1153.
- Loader, C. R. (1996). Local likelihood density estimation. *The Annals of Statistics*, **24**, 1602-1618.
- Mammitzsch, V. (1984). On the asymptotically optimal solution within a certain class of kernel type estimators. *Statistics Decisions*, **2**, 247-255.
- Marron, J.S. and D. Ruppert (1994). Transformations to reduce boundary bias in kernel density estimation. *Journal of the Royal Statistical Society Series B*, **56**, 653-671.
- McLeish, D.L. (1975). A maximal inequality and dependent strong laws. *The Annals of Probability*, **3**, 829-839.
- Müller, H.-G. (1993). On the boundary kernel method for nonparametric curve estimation near endpoints. *Scandinavian Journal of Statistics*, **20**, 313-328.
- Newey, W.K. and K.D. West (1987). A simple, positive-definite, heteroskedasticity and autocorrelation consistent covariance matrix. *Econometrica*, **55**, 703-708.
- Parzen, E. (1962). On estimation of a probability of density function and mode. *Annals of Mathematical Statistics*, **33**, 1065-1076.
- Pritsker, M. (1998). Nonparametric density estimation and tests of continuous time interest rate models. *Review of Financial Studies*, **11**, 449-487.
- Reiss, R.D. (1981). Nonparametric estimation of smooth distribution functions. *Scandinavia Journal of Statistics*, **8**, 116-119.
- Rosenblatt, M. (1956). Remarks on some nonparametric estimates of a density function. *Annals of Mathematical Statistics*, **27**, 832-837.
- Rudemo, M. (1982). Empirical choice of histograms and kernel density estimators. *Scandinavia Journal of Statistics*, **9**, 65-78.
- Schuster, E.F. (1985). Incorporating support constraints into nonparametric estimates of densities. *Communications in Statistics Theory and Methods*, **14**, 1123-1126.
- Serfling, R.J. (1980). *Approximation Theorems of Mathematical Statistics*. Wiley, New York.
- Sheather, S.J. and M.C. Jones (1991). A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society, Series B*, **53**, 683-690.
- Stone, C. J. (1984). An asymptotically optimal window selection rule for kernel density estimates. *The Annals of Statistics*, **12**, 1285-1297.
- Wand, M.P. and M.C. Jones (1995). *Kernel Smoothing*. London: Chapman and Hall.

- Wand, M.P., J.S. Marron and D. Ruppert (1991). Transformations in density estimation (with discussion). *Journal of the American Statistical Association*, **86**, 343-361.
- White, H. (1980). A Heteroskedasticity consistent covariance matrix and a direct test for heteroskedasticity. *Econometrica*, **48**, 817-838.
- Yoshihara, K. (1995). The Bahadur representation of sample quantiles for sequences of strongly mixing random variables. *Statistics and Probability Letters*, **24**, 299-304.
- Zhang, S. and R.J. Karunamuni (1998). On Kernel density estimation near endpoints. *Journal of Statistical Planning and Inference*, **70**, 301-316.

Chapter 4

Nonparametric Regression Models

4.1 Prediction and Regression Functions

Suppose that we have the information set I_t at time t and we want to forecast the future value, say Y_{t+1} (one step-ahead forecast, or Y_{t+s} , s -step ahead). There are several forecasting criteria available in the literature. The general form is

$$m(I_t) = \min_a E[\rho(Y_{t+1} - a) | I_t],$$

where $\rho(\cdot)$ is an **objective (loss) function**. Here are three major directions.

(1) If $\rho(z) = z^2$ is the quadratic function, then, $m(I_t) = E(Y_{t+1} | I_t)$, called the mean regression function. Implicitly, it requires that the distribution of Y_t should be symmetric. If the distribution of Y_t is skewed, then this is not a good criterion.

(2) If $\rho_\tau(y) = y(\tau - I_{\{y < 0\}})$ called the **“check” function**, where $\tau \in (0, 1)$ and I_A is the indicator function of any set A , then, $m(I_t)$ satisfies

$$\int_{-\infty}^{m(I_t)} f(y | I_t) du = F(m(I_t) | I_t) = \tau,$$

where $f(y | I_t)$ and $F(m(I_t) | I_t)$ are the conditional PDF and CDF of Y_{t+1} given I_t , respectively. This $m(I_t)$ becomes the **conditional quantile** or **quantile regression**, denoted by $q_\tau(I_t)$, proposed by Koenker and Bassett (1978, 1982). Particularly, if $\tau = 1/2$, then, $m(I_t)$ is the well known **least absolute deviation (LAD)** regression which is robust. If $q_\tau(I_t)$ is a linear function of regressors like $\beta_\tau^T \mathbf{X}_t$ as in Koenker and Bassett (1978, 1982), Koenker (2005) developed the **R** module **quantreg** to make statistical inferences on the linear quantile regression model.

To fit a linear quantile regression using **R**, one can use the command **rq()** in the package **quantreg**. For a nonlinear parametric model, the command is **nlrq()**. For a nonparametric quantile model for univariate case, one can use the command **lprq()** for implementing the local polynomial estimation. For an additive quantile regression, one can use the commands **rqss()** and **qss()**.

(3) If $\rho(x) = \frac{1}{2} x^2 I_{|x| \leq M} + M(|x| - M/2) I_{|x| > M}$, the so called Huber function in literature, then it is the **Huber robust regression**. We will not discuss this topic. If you have an interest, please read the book by Rousseeuw and Leroy (1987). In **R**, the library **MASS** has the function **rlm** for robust linear model. Also, the library **lqs** contains functions for bounded-influence regression.

Note that for the second and third cases, the regression functions usually do not have a close form of expression. Since the information set I_t contains too many variables (high dimension), it is often to approximate I_t by some finite numbers of variables, say $X_t = (X_{t1}, \dots, X_{tp})^T$ ($p \geq 1$), including the lagged variables and exogenous variables. First, our focus is on the mean regression $m(X_t)$. Of course, by the same token, we can consider the nonparametric estimation of the conditional variance $\sigma^2(x) = \text{Var}(Y_t|X_t = x)$. **Why do we need to consider nonlinear (nonparametric) models in economic practice?** To find the answer, please read the book by Granger and Teräsvirta (1993).

4.2 Kernel Estimation

How to estimate $m(x)$ nonparametrically? Let us look at the **Nadaraya-Watson estimate** of the mean regression $m(x)$. The main idea is as follows:

$$m(x) = \int y f(y|x) dy = \frac{\int y f(x, y) dy}{\int f(x, y) dy},$$

where $f(x, y)$ is the joint PDF of X_t and Y_t . To estimate $m(x)$, we can apply the plug-in method. That is, plug the nonparametric kernel density estimate $f_n(x, y)$ (product kernel method) into the right hand side of the above equation to obtain

$$\hat{m}_{nw}(x) = \frac{\int y f_n(x, y) dy}{\int f_n(x, y) dy} = \dots = \frac{1}{n} \sum_{t=1}^n Y_t K_h(X_t - x) / f_n(x) = \sum_{t=1}^n W_t Y_t,$$

where $f_n(x)$ is the kernel density estimation of $f(x)$, defined in Chapter 3, and

$$W_t = K_h(X_t - x) / \sum_{t=1}^n K_h(X_t - x).$$

$\hat{m}_{nw}(x)$ is the well known **Nadaraya-Watson (NW) estimator**, proposed by Nadaraya (1964) and Watson (1984). Note that the weights $\{W_t\}$ do not depend on $\{Y_t\}$. Therefore, $\hat{m}_{nw}(x)$ is called a **linear estimator**, similar to the least squares estimate (LSE).

Let us look at the NW estimator from a different angle. $\hat{m}_{nw}(x)$ can be re-expressed as the minimizer of the **weighted locally least squares**; that is,

$$\hat{m}_{nw}(x) = \min_a \sum_{t=1}^n (Y_t - a)^2 K_h(X_t - x).$$

This means that when X_t is in a neighborhood of x , $m(X_t)$ is approximated by a constant a (**local approximation**). Indeed, we consider the following working model

$$Y_t = m(X_t) + \varepsilon_t \approx a + \varepsilon_t$$

with the weights $\{K_h(X_t - x)\}$, where $\varepsilon_t = Y_t - E(Y_t | X_t)$. Therefore, the Nadaraya-Watson estimator is also called the **local constant estimator**.

In the implementation, for each x , we can fit the following transformed linear model

$$Y_t^* = \beta_1 X_t^* + \varepsilon_t,$$

where $Y_t^* = \sqrt{K_h(X_t - x)} Y_t$ and $X_t^* = \sqrt{K_h(X_t - x)}$. In **R**, we can use functions **lm()** or **glm()** with weights $\{K_h(X_t - x)\}$ to fit a weighted least squares or generalized linear model. Or, you can use the weighted least squares theory (matrix multiplication); see Section 4.7.

4.2.1 Asymptotic Properties

We derive the asymptotic properties of the nonparametric estimator for the time series situations. Note that the mathematical derivations are different for the iid case and time series situations since $E[Y_t | X_1, \dots, X_n] \neq E[Y_t | X_t] = m(X_t)$ which is true for the iid case. To easy notation, we consider only the simple case when $p = 1$.

$$\hat{m}_{nw}(x) = \underbrace{\frac{1}{n} \sum_{t=1}^n m(X_t) K_h(X_t - x) / f_n(x)}_{I_1} + \underbrace{\sum_{t=1}^n W_t \varepsilon_t}_{I_2}.$$

We will show that I_1 contributes only the **asymptotic bias** and I_2 gives the **asymptotic normality**. First, we derive the asymptotic bias for the interior boundary points. By the Taylor's expansion, when X_t is in $(x - h, x + h)$, we have

$$m(X_t) = m(x) + m'(x)(X_t - x) + \frac{1}{2} m''(x)(X_t - x)^2,$$

where $x_t = x + \theta(X_t - x)$ with $-1 < \theta < 1$. Then,

$$\begin{aligned} I_{11} &\equiv \frac{1}{n} \sum_{t=1}^n m(X_t) K_h(X_t - x) = m(x) f_n(x) + m'(x) \underbrace{\frac{1}{n} \sum_{t=1}^n (X_t - x) K_h(X_t - x)}_{J_1(x)} \\ &\quad + \frac{1}{2} \underbrace{\frac{1}{n} \sum_{t=1}^n m''(x_t) (X_t - x)^2 K_h(X_t - x)}_{J_2(x)}. \end{aligned}$$

Then,

$$\begin{aligned} E[J_1(x)] &= E[(X_t - x) K_h(X_t - x)] = \int (u - x) K_h(u - x) f(u) du \\ &= h \int u K(u) f(x + hu) du = h^2 f'(x) \mu_2(K) + o(h^2). \end{aligned}$$

Similar to the derivation of the variance of $f_n(x)$ in (3.3), we can show that

$$nh \text{Var}(J_1(x)) = O(1).$$

Therefore, $J_1(x) = h^2 f'(x) \mu_2(K) + o_p(h^2)$. By the same token, we have

$$\begin{aligned} E[J_2(x)] &= E[m''(x_t) (X_t - x)^2 K_h(X_t - x)] \\ &= h^2 \int m''(x + \theta hu) u^2 K(u) f(x + hu) du = h^2 m''(x) \mu_2(K) f(x) + o(h^2) \end{aligned}$$

and $\text{Var}(J_2(x)) = O(1/nh)$. Therefore, $J_2(x) = h^2 m''(x) \mu_2(K) f(x) + o_p(h^2)$. Hence,

$$\begin{aligned} I_1 &= m(x) + m'(x) J_1(x)/f_n(x) + \frac{1}{2} J_2(x)/f_n(x) \\ &= m(x) + \frac{h^2}{2} \mu_2(K) [m''(x) + 2m'(x)f'(x)/f(x)] + o_p(h^2) \end{aligned}$$

by the fact that $f_n(x) = f(x) + o_p(1)$. The term

$$B_{nw}(x) = \frac{h^2}{2} \mu_2(K) [m''(x) + 2m'(x)f'(x)/f(x)] \quad (4.1)$$

is regarded as the **asymptotic bias**. The bias term involves not only **curvatures** of $m(x)$ ($m''(x)$) but also the unknown density function $f(x)$ and its derivative $f'(x)$ so that the design can not be adaptive.

Under some regularity conditions, similar to (3.3), we can show that for x being an interior grid point,

$$nh \text{Var}(I_2) \rightarrow \nu_0(K) \sigma_\varepsilon^2(x)/f(x) = \sigma_m^2(x),$$

where $\sigma_\varepsilon^2(x) = \text{Var}(\varepsilon_t | X_t = x)$. Further, we can establish the asymptotic normality (the proof is provided later)

$$\sqrt{n h} [\widehat{m}_{nw}(x) - m(x) - B_{nw}(x) + o_p(h^2)] \rightarrow N \{0, \sigma_m^2(x)\},$$

where $B_{nw}(x)$ is given in (4.1).

4.2.2 Boundary Behavior

For expositional purpose, in what follows, we only consider the case when $p = 1$. As for the boundary behavior of the NW estimator, we can follow Fan and Gijbels (1996). Without loss of generality, we consider the left boundary point $x = ch$, $0 < c < 1$. From Fan and Gijbels (1996), we take $K(\cdot)$ to have support $[-1, 1]$ and $m(\cdot)$ to have support $[0, 1]$. Similar to (3.5), it is easy to see that if $x = ch$,

$$\begin{aligned} E[J_1(ch)] &= E[(X_t - ch) K_h(X_t - ch)] = \int_0^1 (u - ch) K_h(u - ch) f(u) du \\ &= h \int_{-c}^{1/h-c} u K(u) f(h(u + c)) du \\ &= h f(0+) \mu_{1,c}(K) + h^2 f'(0+) [\mu_{2,c}(K) + c \mu_{1,c}(K)] + o(h^2), \end{aligned}$$

and

$$\begin{aligned} E[J_2(ch)] &= E[m''(x_t)(X_t - ch)^2 K_h(X_t - ch)] \\ &= h^2 \int_{-c}^{1/h-c} m''(h(c + \theta u)) u^2 K(u) f(h(u + c)) du \\ &= h^2 m''(0+) \mu_{2,c}(K) f(0+) + o(h^2). \end{aligned}$$

Also, we can see that

$$\text{Var}(J_1(ch)) = O(1/nh) \quad \text{and} \quad \text{Var}(J_2(ch)) = O(1/nh),$$

which imply that

$$J_1(ch) = h f(0+) \mu_{1,c}(K) + o_p(h) \quad \text{and} \quad J_2(ch) = h^2 m''(0+) \mu_{2,c}(K) f(0+) + o(h^2).$$

This, in conjunction with (3.5), gives

$$I_1 - m(ch) = m'(ch) J_1(ch)/f_n(ch) + \frac{1}{2} J_2(ch)/f_n(ch) = a(c, K) h + b(c, K) h^2 + o_p(h^2),$$

where

$$a(c, K) = \frac{m'(0+)\mu_{1,c}(K)}{\mu_{0,c}(K)},$$

and

$$b(c, K) = \frac{\mu_{2,c}(K) m''(0+)}{2 \mu_{0,c}(K)} + \frac{f'(0+)m'(0+)[\mu_{2,c}(K) \mu_{0,c}(K) - \mu_{1,c}^2(K)]}{f(0+) \mu_{0,c}^2(K)}.$$

Here, $a(c, K) h + b(c, K) h^2$ serves as the **asymptotic bias term**, which is of the order $O(h)$.

We can show that at the boundary point, the asymptotic variance has the following form

$$n h \text{Var}(\hat{m}_{nw}(x)) \rightarrow \nu_{0,c}(K) \sigma_m^2(0+)/[\mu_{0,c}(K) f(0+)],$$

which the **same order** as that for the interior point although the scaling constant is different.

4.3 Local Polynomial Estimate

To overcome the above shortcomings of local constant estimate, we can use the local polynomial fitting scheme; see Fan and Gijbels (1996). The main idea is described as follows.

4.3.1 Formulation

Assume that the regression function $m(x)$ has $(q + 1)$ th order continuous derivative. For ease notation, assume that $p = 1$. When $X_t \in (x - h, x + h)$, then

$$m(X_t) \approx \sum_{j=0}^q \frac{m^{(j)}(x)}{j!} (X_t - x)^j = \sum_{j=0}^q \beta_j (X_t - x)^j,$$

where $\beta_j = m^{(j)}(x)/j!$. Therefore, when $X_t \in (x - h, x + h)$, the model becomes

$$Y_t \approx \sum_{j=0}^q \beta_j (X_t - x)^j + \varepsilon_t.$$

Hence, we can apply the weighted least squares method. The **weighted locally least squares** becomes

$$\sum_{t=1}^n \left(Y_t - \sum_{j=0}^q \beta_j (X_t - x)^j \right)^2 K_h(X_t - x). \quad (4.2)$$

Minimizing the above with respect to $\beta = (\beta_0, \dots, \beta_q)^T$ to obtain the **local polynomial estimate** $\hat{\beta}$;

$$\hat{\beta} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} Y, \quad (4.3)$$

where $\mathbf{W} = \text{diag}\{K_h(X_1 - x), \dots, K_h(X_n - x)\}$,

$$\mathbf{X} = \begin{pmatrix} 1 & (X_1 - x) & \cdots & (X_1 - x)^q \\ 1 & (X_2 - x) & \cdots & (X_2 - x)^q \\ \vdots & \vdots & \ddots & \vdots \\ 1 & (X_n - x) & \cdots & (X_n - x)^q \end{pmatrix}, \quad \text{and} \quad Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}.$$

Therefore, for $1 \leq j \leq q$,

$$\hat{m}^{(j)}(x) = j! \hat{\beta}_j.$$

This means that the local polynomial method estimates not only the regression function itself but also derivatives of regression.

4.3.2 Implementation in R

There are several ways of implementing the local polynomial estimator. One way you can do so is to write your own code by using matrix multiplication as in (4.3) or employing function **lm()** or **glm()** with weights $\{K_h(X_t - x)\}$. Recently, in **R**, there are some build-in packages for implementing the local polynomial estimate. For example, the package **KernSmooth** contains several functions. Function **bkde()** computes the kernel density estimate and Function **bkde2D()** computes the 2D kernel density estimate as well as Function **bkfe()** computes the kernel functional (derivative) density estimate. Function **dpik()** selects a bandwidth for estimating the kernel density estimation using the plug-in method and Function **dpill()** chooses a bandwidth for the local linear ($q = 1$) regression estimation using the plug-in approach. Finally, Function **locpoly()** is for the local polynomial fitting including a local polynomial estimate of the density of X (or its derivative) if the dependent variable is omitted.

Example 4.1: We apply the kernel regression estimation and local polynomial fitting methods to estimate the drift and diffusion of the weekly 3-month Treasury bill from January 2, 1970 to December 26, 1997. Let x_t denote the weekly 3-month Treasury bill. It is often to model x_t by assuming that it satisfies the continuous-time stochastic differential equation (Black-Scholes model)

$$dx_t = \mu(x_t) dt + \sigma(x_t) dW_t,$$

where W_t is a Wiener process, $\mu(x_t)$ is called the drift function and $\sigma(x_t)$ is called the diffusion function. Our interest is to identify $\mu(x_t)$ and $\sigma(x_t)$. Assume a time series sequence $\{X_{t\Delta}, 1 \leq t \leq n\}$ is observed at **equally spaced** time points. Using the **infinitesimal generator** (Øksendal, 1985), the first-order approximations of moments of x_t , a discretized version of the Ito's process, are given by Stanton (1997) (see Fan and Zhang (2003) for the higher orders)

$$\Delta x_t = \mu(x_t) \Delta + \sigma(x_t) \varepsilon \sqrt{\Delta},$$

where $\Delta x_t = x_{t+\Delta} - x_t$, $\varepsilon \sim N(0, 1)$, and x_t and ε_t are independent. Therefore,

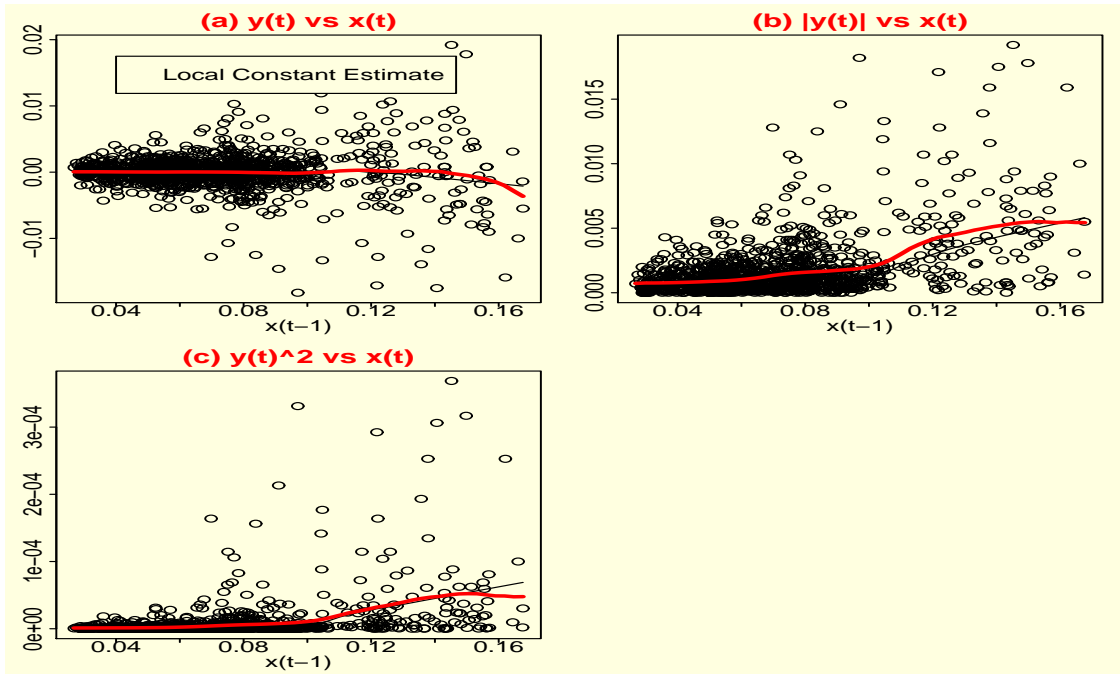


Figure 4.1: Scatterplots of Δx_t , $|\Delta x_t|$, and $(\Delta x_t)^2$ versus x_t with the smoothed curves computed using `scatter.smooth()` and the local constant estimation.

$$\mu(x_t) = \lim_{\Delta \rightarrow 0} E[\Delta x_t | x_t] / \Delta \quad \text{and} \quad \sigma^2(x_t) = \lim_{\Delta \rightarrow 0} E[(\Delta x_t)^2 | x_t] / \Delta.$$

Hence, estimating $\mu(x)$ and $\sigma^2(x)$ becomes a nonparametric regression problem. We can use both local constant and local polynomial method to estimate $\mu(x)$ and $\sigma^2(x)$. As a result, the local constant estimators (red line) together with the **lowess** smoothers (black line) and the scatterplots of Δx_t [in (a)], $|\Delta x_t|$ [in (b)], and $(\Delta x_t)^2$ [in (c)] versus x_t are presented in Figure 4.1 and the local linear estimators (red line) together with the **lowess** smoothers (black line) and the scatterplots of Δx_t [in (a)], $|\Delta x_t|$ [in (b)], and $(\Delta x_t)^2$ [in (c)] versus x_t are displaced in Figure 4.2. An alternative approach can be found in Aït-Sahalia (1996).

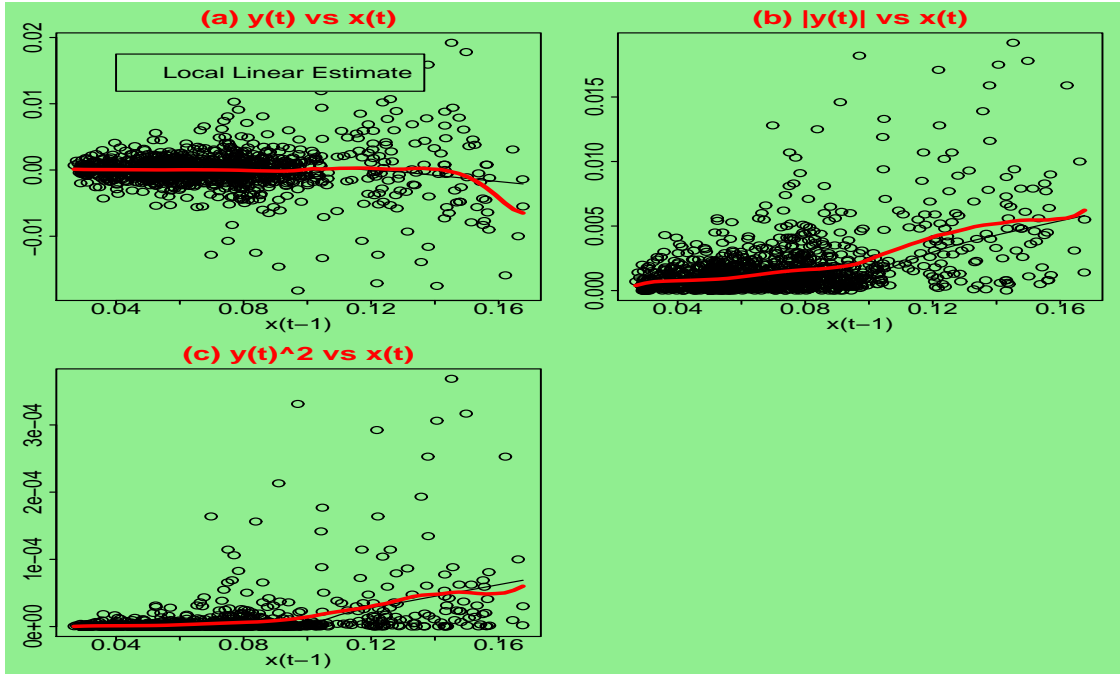


Figure 4.2: Scatterplots of Δx_t , $|\Delta x_t|$, and $(\Delta x_t)^2$ versus x_t with the smoothed curves computed using `scatter.smooth()` and the local linear estimation.

4.3.3 Complexity of Local Polynomial Estimator

To implement the local polynomial estimator, we have to choose the order of the polynomial q , the bandwidth h and the kernel function $K(\cdot)$. These parameters are of course confounded each other. Clearly, when $h = \infty$, the local polynomial fitting becomes a global polynomial fitting and the order q determines the model complexity. Unlike in the parametric models, the complexity of local polynomial fitting is primarily controlled by the **bandwidth**, as shown in Fan and Gijbels (1996) and Fan and Yao (2003). Hence q is usually small and the issue of choosing q becomes less critical. We discuss those issues in detail as follows.

(1) If the objective is to estimate $m^{(j)}(\cdot)$ ($j \geq 0$), the local polynomial fitting corrects automatically the boundary bias when $q - j$ is odd. Further, when $q - j$ is odd, comparing with the order $q - 1$ fit (so that $q - j - 1$ is even), the order q fit contains one extra parameter without increasing the variance for estimating $m^{(j)}(\cdot)$. But this extra parameter creates opportunities for bias reduction, particularly in the boundary regions; see the next section and the books by Fan and Gijbels (1996) and Ruppert and Wand (1994). For these reasons, the odd order fits (the order q is chosen so that $q - j$ is odd) outperforms the even order fits [the order $(q - 1)$ fit so that q is even]. Based on theoretical and practical considerations,

the order $q = j + 1$ is recommended in Fan and Gijbels (1996). **If the primary objective is to estimate the regression function, one uses local linear fit and if the target function is the first order derivative, one uses the local quadratic fit and so on.**

(2) It is well known that the choice of the bandwidth h plays an important role in any kernel smoothing, including the local polynomial fitting. A too large bandwidth causes over-smoothing (reducing variance), creating excessive modeling bias, while a too small bandwidth results in under-smoothing (reducing bias but increasing variance), obtaining wiggly estimates. **The bandwidth can be subjectively chosen by users via visually inspecting resulting estimates, or automatically chosen by data via minimizing an estimated theoretical risk** (discussed later). Since the choice of bandwidth is not easy task, it is often attached by people who do not know well nonparametric techniques.

(3) Since the estimate is based on the local regression (4.2), it is reasonable to require a non-negative weight function $K(\cdot)$. It can be shown (see Fan and Gijbels (1996)) that for all choices of q and j , the optimal weight function is $K(z) = 3/4(1 - z^2)_+$, **the Epanechnikov kernel**, based on minimizing the asymptotic variance of the local polynomial estimator. Thus, it is a universal weighting scheme and provides a useful benchmark for other kernels to compare with. As shown in Fan and Gijbels (1996) and Fan and Yao (2003), other kernels have nearly the same efficiency for practical use of q and j . **Hence the choice of the kernel function is not critical.**

The local polynomial estimator compares favorably with other estimators, including the **Nadaraya-Watson (local constant) estimator** and other linear estimators such as the **Gasser and Müller estimator** of Gasser and Müller (1979) and the **Priestley and Chao estimator** of Priestley and Chao (1972). Indeed, it was shown by Fan (1993) that the local linear fitting is **asymptotically minimax** based on the quadratic loss function among all linear estimators and is nearly minimax among all possible linear estimators. This minimax property is extended by Fan, Gasser, Gijbels, Brockmann and Engel (1995) to more general local polynomial fitting. For the detailed comparisons of the above four estimators, see Fan and Gijbels (1996).

Note that the Gasser and Müller estimator and the Priestley and Chao estimator are particularly for the fixed design. That is, $X_t = t$. Let $s_t = (2t + 1)/2$ ($t = 1, \dots, n - 1$) with

$s_0 = -\infty$ and $s_n = \infty$. The Gasser and Müller estimator is

$$\widehat{m}_{gm}(t_0) = \sum_{t=1}^n \int_{s_{t-1}}^{s_t} K_h(u - t_0) du Y_t.$$

Unlike the local constant estimator, no denominator is needed since the total weight

$$\sum_{t=1}^n \int_{s_{t-1}}^{s_t} K_h(u - t_0) du = 1.$$

Indeed, the Gasser and Müller estimator is an improved version of the Priestley and Chao estimator, which is defined as

$$\widehat{m}_{pc}(t_0) = \sum_{t=1}^n K_h(t - t_0) Y_t.$$

Note that the Priestley and Chao estimator is only applicable for the equi-space setting.

4.3.4 Properties of Local Polynomial Estimator

Define, for $0 \leq j \leq q$,

$$s_{n,j}(x) = \sum_{t=1}^n (X_t - x)^j K_h(X_t - x)$$

and $S_n(x) = \mathbf{X}^T \mathbf{W} \mathbf{X}$. Then, the $(i+1, j+1)$ th element of $S_n(x)$ is $s_{n,i+j}(x)$. Similar to the evaluation of I_{11} , we can show easily that

$$s_{n,j}(x) = n h^j \mu_j(K) f(x) \{1 + o_p(1)\}.$$

Define, $H = \text{diag}\{1, h, \dots, h^q\}$ and $S = (\mu_{i+j}(K))_{0 \leq i, j \leq q}$. Then, it is not difficult to show that $S_n(x) = n f(x) H S H \{1 + o_p(1)\}$.

First of all, for $0 \leq j \leq q$, let e_j be a $(q+1) \times 1$ vector with $(j+1)$ th element being one and zero otherwise. Then, $\widehat{\beta}_j$ can be re-expressed as

$$\widehat{\beta}_j = e_j^T \widehat{\beta} = \sum_{t=1}^n W_{j,n,h}(X_t - x) Y_t,$$

where $W_{j,n,h}(X_t - x)$ is called the effective kernel in Fan and Gijbels (1996) and Fan and Yao (2003), given by

$$W_{j,n,h}(X_t - x) = e_j^T S_n(x)^{-1} (1, (X_t - x), \dots, (X_t - x)^q)^T K_h(X_t - x).$$

It is not difficult to show (based on the least square theory) that $W_{j,n,h}(X_t - x)$ satisfies the following the so-called **discrete moment conditions**

$$\sum_{t=1}^n (X_t - x)^l W_{j,n,h}(X_t - x) = \begin{cases} 1 & \text{if } l = j, \\ 0 & \text{otherwise.} \end{cases} \quad (4.4)$$

Note that the local constant estimator does not have this property; see $J_1(x)$ in Section 4.2.1. This property implies that the local polynomial estimator is unbiased for estimating β_j , when the true regression function $m(x)$ is a polynomial of order q .

To gain more insights about the local polynomial estimator, define the **equivalent kernel** (see Fan and Gijbels (1996))

$$W_j(u) = e_j^T S^{-1} (1, u, \dots, u^q)^T K(u).$$

Then, it can be shown (see Fan and Gijbels (1996)) that

$$W_{j,n,h}(X_t - x) = \frac{1}{n h^{j+1} f(x)} W_j((X_t - x)/h) \{1 + o_p(1)\}$$

and

$$\int u^l W_j(u) du = \begin{cases} 1 & \text{if } l = j, \\ 0 & \text{otherwise.} \end{cases}$$

The implications of these results are as follows.

As pointed out by Fan and Yao (2003), the local polynomial estimator works like a kernel regression estimation with a known design density $f(x)$. This explains why the local polynomial fit **adapts** to various design densities. In contrast, the kernel regression estimator has large bias at the region where the derivative of $f(x)$ is large, namely it can not adapt to highly-skewed designs. To see that, imagine the true regression function has large slope in this region. Since the derivative of design density is large, for a given x , there are more points on one side of x than the other. When the local average is taken, the Nadaraya-Watson estimate is biased towards the side with more local data points because the local data are asymmetrically distributed. This issue is more pronounced at the boundary regions, since the local data are even more asymmetric. On the other hand, the local polynomial fit creates asymmetric weights, if needed, to compensate for this kind of design bias. Hence, it is adaptive to various design densities and to the boundary regions.

We next derive the asymptotic bias and variance expression for local polynomial estimators. For independent data, we can obtain the bias and variance expression via conditioning on the design matrix \mathbf{X} . However, for time series data, conditioning on \mathbf{X} would mean conditioning on nearly the entire series. Hence, we derive the asymptotic bias and variance using the asymptotic normality rather than conditional expectation. As explained in Chapter 3, localizing in the state domain weakens the dependent structure for the local data. Hence, one would expect that the result for the independent data continues to hold for the stationary process with certain mixing conditions. The mixing condition and the bandwidth should be related, which can be seen later.

Set $B_n(x) = (b_1(x), \dots, b_n(x))^T$, where, for $0 \leq j \leq q$,

$$b_{j+1}(x) = \sum_{t=1}^n \left[m(X_t) - \sum_{j=0}^q \frac{m^{(j)}(x)}{j!} (X_t - x)^j \right] (X_t - x)^j K_h(X_t - x).$$

Then,

$$\hat{\beta} - \beta = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} B_n(x) + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \varepsilon,$$

where $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$. It is easy to show that if q is odd,

$$B_n(x) = n h^{q+1} H f(x) \frac{m^{(q+1)}(x)}{(q+1)!} c_{1,q} \{1 + o_p(1)\},$$

where, for $1 \leq k \leq 3$, $c_{k,q} = (\mu_{q+k}(K), \dots, \mu_{2q+k}(K))^T$. If q is even,

$$B_n(x) = n h^{q+2} H f(x) \left[c_{2,q} \frac{m^{(q+1)}(x) f'(x)}{f(x)(q+1)!} + c_{3,q} \frac{m^{(q+2)}(x)}{(q+2)!} \right] \{1 + o_p(1)\}.$$

Note that $f'(x)/f(x)$ does not appear in the right hand side of $B_n(x)$ when q is odd. In either case, we can show that

$$n h \text{Var} \left[H(\hat{\beta} - \beta) \right] \rightarrow \sigma^2(x) S^{-1} S^* S^{-1} / f(x) = \Sigma(x),$$

where S^* is a $(q+1) \times (q+1)$ matrix with the (i, j) th element being $\nu_{i+j-2}(K)$.

This shows that the leading conditional bias term depends on whether q is odd or even. By a Taylor series expansion argument, we know that when considering $|X_t - x| < h$, the remainder term from a q th order polynomial expansion should be of order $O(h^{q+1})$, so the result for odd q is quite easy to understand. When q is even, $(q+1)$ is odd hence the term h^{q+1} is associated with $\int u^l K(u) du$ for l odd, and this term is zero because $K(u)$ is a even

function. Therefore, the h^{q+1} term disappears, while the remainder term becomes $O(h^{q+2})$. Since q is either odd or even, then we see that the bias term is an even power of h . This is similar to the case where one uses higher order kernel functions based upon a symmetric kernel function (an even function), where the bias is always an even power of h .

Finally, we can show that when q is odd,

$$\sqrt{nh} \left[H(\hat{\beta} - \beta) - B(x) \right] \rightarrow N(0, \Sigma(x)),$$

the asymptotic bias term for the local polynomial estimator is

$$B(x) = \frac{h^{q+1}}{(q+1)!} m^{(q+1)}(x) S^{-1} c_{1,q} \{1 + o_p(1)\}.$$

Or,

$$\sqrt{nh^{2j+1}} \left[\hat{m}^{(j)}(x) - m^{(j)}(x) - B_j(x) \right] \rightarrow N(0, \sigma_{jj}(x)),$$

where the asymptotic bias and variance for the local polynomial estimator of $m^{(j)}(x)$ are

$$B_j(x) = \frac{j! h^{q+1-j}}{(q+1)!} m^{(q+1)}(x) \int u^{q+1} W_j(u) du \{1 + o_p(1)\}$$

and

$$\sigma_{jj}(x) = \frac{(j!)^2 \sigma^2(x)}{f(x)} \int W_j^2(u) du.$$

Similarly, we can derive the asymptotic bias and variance at boundary points if the regression function has a finite support. For details, see Fan and Gijbels (1996), Fan and Yao (2003), and Ruppert and Wand (1994). Indeed, define S_c , S_c^* , and $c_{k,q,c}$ similarly to S , S^* and $c_{k,q}$ with $\mu_j(K)$ and $\nu_j(K)$ replaced by $\mu_{j,c}(K)$ and $\nu_{j,c}(K)$ respectively. We can show that

$$\sqrt{nh} \left[H(\hat{\beta}(ch) - \beta(ch)) - B_c(0) \right] \rightarrow N(0, \Sigma_c(0)), \quad (4.5)$$

where the asymptotic bias term for the local polynomial estimator at the left boundary point is

$$B_c(0) = \frac{h^{q+1}}{(q+1)!} m^{(q+1)}(0) S_c^{-1} c_{1,q,c} \{1 + o_p(1)\},$$

and the asymptotic variance is $\Sigma_c(0) = \sigma^2(0) S_c^{-1} S_c^* S_c^{-1} / f(0)$. Or,

$$\sqrt{nh^{2j+1}} \left[\hat{m}^{(j)}(ch) - m^{(j)}(ch) - B_{j,c}(0) \right] \rightarrow N(0, \sigma_{jj,c}(0)),$$

where with $W_{j,c}(u) = e_j^T S_c^{-1} (1, u, \dots, u^q)^T K(u)$,

$$B_{j,c}(0) = \frac{j! h^{q+1-j}}{(q+1)!} m^{(q+1)}(0) \int_{-c}^{\infty} u^{q+1} W_{j,c}(u) du \{1 + o_p(1)\}$$

and

$$\sigma_{jj,c}(0) = \frac{(j!)^2 \sigma^2(0)}{f(0)} \int_{-c}^{\infty} W_{j,c}^2(u) du.$$

Exercise: Please derive the asymptotic properties for the local polynomial estimator. That is to prove (4.5).

The above conclusions show that when $q - j$ is odd, the bias at the boundary is of the same order as that for points on the interior. Hence, the local polynomial fit does not create excessive boundary bias when $q - j$ is odd. Thus, the appealing boundary behavior of local polynomial mean estimation extends to derivative estimation. However, when $q - j$ is even, the bias at the boundary is larger than in the interior, and the bias can also be large at points where $f(x)$ is discontinuous. This is referred to as boundary effect. For these reasons (and the minimax efficiency arguments), it is recommended that one strictly set $q - j$ to be odd when estimating $m^{(j)}(x)$. It is indeed an odd world!

4.3.5 Bandwidth Selection

As seen in previous sections, for stationary sequences of data under certain mixing conditions, the local polynomial estimator performs very much like that for independent data, because windowing reduces dependency among local data. Partially because of this, there are not many studies on bandwidth selection for these problems. However, it is reasonable to expect the bandwidth selectors for independent data continue to work for dependent data with certain mixing conditions. Below, we summarize a few of useful approaches. When data do not have strong enough mixing, the general strategy is to increase bandwidth in order to reduce the variance.

As what we had already seen for the nonparametric density estimation, the **cross-validation** method is very useful for assessing the performance of an estimator via estimating its prediction error. The basic idea is to set one of the data point aside for validation of a model and use the remaining data to build the model. It is defined as

$$CV(h) = \sum_{s=1}^n [Y_s - \hat{m}_{-s}(X_s)]^2$$

where $\hat{m}_{-s}(X_s)$ is the local polynomial estimator with $j = 0$ and bandwidth h , but without using the s th observation. The above summand is indeed a squared-prediction error

of the s th data point using the training set $\{(X_t, Y_t) : t \neq s\}$. This idea of the cross-validation method is simple but is computationally intensive. An improved version, in terms of computation, is the generalized cross-validation (GCV), proposed by Wahba (1977) and Craven and Wahba (1979). This criterion can be described as follows. The fitted values $\hat{Y} = (\hat{m}(X_1), \dots, \hat{m}(X_n))^T$ can be expressed as $\hat{Y} = H(h)Y$, where $H(h)$ is an $n \times n$ matrix, depending on the \mathbf{X} -variate and bandwidth h , and it is also called a smoothing matrix. Then the GCV approach selects the bandwidth h that minimizes

$$\text{GCV}(h) = [n^{-1} \text{tr}(I - H(h))]^{-2} \text{MASE}(h)$$

where $\text{MASE}(h) = \sum_{t=1}^n (Y_t - \hat{m}(X_t))^2 / n$ is the average of squared residuals.

A drawback of the cross-validation type method is its inherited variability (see Hall and Johnstone, 1992). Further, it can not be directly applied to select bandwidths for estimating derivative curves. As pointed out by Fan, Heckman, and Wand (1995), the cross-validation type method performs poorly due to its large sample variation, even worse for dependent data. Plug-in methods avoid these problems. The basic idea is to find a bandwidth h minimizing estimated mean integrated square error (MISE). See Ruppert, Sheather and Wand (1995) and Fan and Gijbels (1995) for details.

Nonparametric AIC Selector

Inspired by the nonparametric version of the Akaike final prediction error criterion proposed by Tjøstheim and Auestad (1994b) for the lag selection in nonparametric setting, Cai (2002) proposed a simple and quick method to select bandwidth for the foregoing estimation procedures, which can be regarded as a nonparametric version of the Akaike information criterion (AIC) to be attentive to the structure of time series data and the over-fitting or under-fitting tendency. Note that the idea is also motivated by its analogue of Cai and Tiwari (2000). The basic idea is described as follows.

By recalling the classical AIC for linear models under the likelihood setting

$$-2 (\text{maximized log likelihood}) + 2 (\text{number of estimated parameters}),$$

Cai (2002) proposed the following nonparametric AIC to select h minimizing

$$\text{AIC}(h) = \log \{\text{MASE}\} + \psi(\text{tr}(H(h)), n), \quad (4.6)$$

where $\psi(\text{tr}(H(h)), n)$ is chosen particularly to be the form of the bias-corrected version of the AIC, due to Hurvich and Tsai (1989),

$$\psi(\text{tr}(H(h)), n) = 2 \{ \text{tr}(H(h)) + 1 \} / [n - \{ \text{tr}(H(h)) + 2 \}], \quad (4.7)$$

and $\text{tr}(H(h))$ is the trace of the smoothing matrix $H(h)$, regarded as the nonparametric version of degrees of freedom, called the effective number of parameters. See the book by Hastie and Tibshirani (1990, Section 3.5) for the detailed discussion on this aspect for nonparametric models. Note that actually, (4.6) is a generalization of the AIC for the parametric regression and autoregressive time series contexts, in which $\text{tr}(H(h))$ is the number of regression (autoregressive) parameters in the fitting model. In view of (4.7), when $\psi(\text{tr}(H(h)), n) = -2 \log(1 - \text{tr}(H(h))/n)$, then (4.6) becomes the generalized cross-validation (GCV) criterion, commonly used to select the bandwidth in the time series literature even in the iid setting, when $\psi(\text{tr}(H(h)), n) = 2 \text{tr}(H(h))/n$, then (4.6) is the classical AIC discussed in Engle, Granger, Rice, and Weiss (1986) for time series data, and when $\psi(\text{tr}(H(h)), n) = -\log(1 - 2 \text{tr}(H(h))/n)$, (4.6) is the T-criterion, proposed and studied by Rice (1984) for iid samples. It is clear that when $\text{tr}(H(h))/n \rightarrow 0$, then the nonparametric AIC, the GCV and the T-criterion are asymptotically equivalent. However, the T-criterion requires $\text{tr}(H(h))/n < 1/2$, and, when $\text{tr}(H(h))/n$ is large, the GCV has relatively weak penalty. This is especially true for the nonparametric setting. Therefore, the criterion proposed here counteracts the over-fitting tendency of the GCV. Note that Hurvich, Simonoff, and Tsai (1998) gave the detailed derivation of the nonparametric AIC for the nonparametric regression problems under the iid Gaussian error setting and they argued that the nonparametric AIC performs reasonably well and better than some existing methods in the literature.

4.4 Project for Regression Function Estimation

Do Monte Carlo simulations to compare the performances of the local linear and local constant estimations for the nonparametric regression function for different settings and to make your own conclusions based on your simulations. Please do the followings:

1. Choosing different sample sizes, different kernels, different bandwidths, and different bandwidth selection methods. Any conclusions and comments?

2. Compare the local linear method with local constant method. Any conclusions and comments?

Note that you can choose any distribution to generate your samples for your simulation. You can use any statistical package to do your simulation. You try to use **R** since it is very simple. You have to write **report** to present what you do in details and to explain what you observe as well as to make your comments. Please hand in all necessary materials (tables or graphs) to support your conclusions. If you need any help, please come to see me.

4.5 Functional Coefficient Model

4.5.1 Model

As mentioned earlier, when p is large, there exists the so called curse of dimensionality. To overcome this shortcoming, one way to do so is to consider the functional coefficient model as studied in Cai, Fan and Yao (2000) and the additive model discussed in Section 4.6. First, we study the functional coefficient model. To use the notation from Cai, Fan and Yao (2000), we change the notation from the previous sections.

Let $\{\mathbf{U}_i, \mathbf{X}_i, Y_i\}_{i=-\infty}^{\infty}$ be jointly strictly stationary processes with \mathbf{U}_i taking values in \mathbb{R}^k and \mathbf{X}_i taking values in \mathbb{R}^p . Typically, k is small. Let $E(Y_1^2) < \infty$. We define the multivariate regression function

$$m(\mathbf{u}, \mathbf{x}) = E(Y | \mathbf{U} = \mathbf{u}, \mathbf{X} = \mathbf{x}), \quad (4.8)$$

where $(\mathbf{U}, \mathbf{X}, Y)$ has the same distribution as $(\mathbf{U}_i, \mathbf{X}_i, Y_i)$. In a pure time series context, both \mathbf{U}_i and \mathbf{X}_i consist of some lagged values of Y_i . The functional-coefficient regression model has the form

$$m(\mathbf{u}, \mathbf{x}) = \sum_{j=1}^p a_j(\mathbf{u}) x_j, \quad (4.9)$$

where the functions $\{a_j(\cdot)\}$ are measurable from \mathbb{R}^k to \mathbb{R}^1 and $\mathbf{x} = (x_1, \dots, x_p)^T$. This model has been studied extensively in the literature; see Cai, Fan and Yao (2000) for the detailed discussions.

For simplicity, in what follows, we consider only the case $k = 1$ in (4.9). Extension to the case $k > 1$ involves no fundamentally new ideas. Note that models with large k are often

not practically useful due to the “curse of dimensionality”. If k is large, to overcome the problem, one way to do so is to consider an index functional coefficient model proposed by Fan, Yao and Cai (2003)

$$m(\mathbf{u}, \mathbf{x}) = \sum_{j=1}^p a_j(\boldsymbol{\beta}^T \mathbf{u}) x_j, \quad (4.10)$$

where $\beta_1 = 1$. Fan, Yao and Cai (2003) studied the estimation procedures, bandwidth selection and applications. Hong and Lee (2003) considered the applications of model (4.10) to the exchange rates, Juhl (2005) studied the unit root behavior of nonlinear time series models, Li, Huang, Li and Fu (2002) modelled the production frontier using China’s manufacturing industry data, and Cai, Das, Xiong and Wu (2006) considered the nonparametric two-stage instrumental variable estimators for returns to education.

4.5.2 Local Linear Estimation

As recommended by Fan and Gijbels (1996), we estimate the coefficient functions $\{a_j(\cdot)\}$ using the local linear regression method from observations $\{U_i, \mathbf{X}_i, Y_i\}_{i=1}^n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^T$. We assume throughout that $a_j(\cdot)$ has a continuous second derivative. Note that we may approximate $a_j(\cdot)$ locally at u_0 by a linear function $a_j(u) \approx a_j + b_j(u - u_0)$. The local linear estimator is defined as $\hat{a}_j(u_0) = \hat{a}_j$, where $\{(\hat{a}_j, \hat{b}_j)\}$ minimize the sum of weighted squares

$$\sum_{i=1}^n \left[Y_i - \sum_{j=1}^p \{a_j + b_j(U_i - u_0)\} X_{ij} \right]^2 K_h(U_i - u_0), \quad (4.11)$$

where $K_h(\cdot) = h^{-1}K(\cdot/h)$, $K(\cdot)$ is a kernel function on \mathbb{R}^1 and $h > 0$ is a bandwidth. It follows from the least squares theory that

$$\hat{a}_j(u_0) = \sum_{k=1}^n K_{n,j}(U_k - u_0, \mathbf{X}_k) Y_k, \quad (4.12)$$

where

$$K_{n,j}(u, \mathbf{x}) = \mathbf{e}_{j,2p}^T \left(\tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}} \right)^{-1} \begin{pmatrix} \mathbf{x} \\ u \end{pmatrix} K_h(u) \quad (4.13)$$

$\mathbf{e}_{j,2p}$ is the $2p \times 1$ unit vector with 1 at the j th position, $\tilde{\mathbf{X}}$ denotes an $n \times 2p$ matrix with $(\mathbf{X}_i^T, \mathbf{X}_i^T(U_i - u_0))$ as its i th row, and $\mathbf{W} = \text{diag}\{K_h(U_1 - u_0), \dots, K_h(U_n - u_0)\}$.

4.5.3 Bandwidth Selection

Various existing bandwidth selection techniques for nonparametric regression can be adapted for the foregoing estimation; see, *e.g.*, Fan, Yao, and Cai (2003) and the nonparametric AIC as discussed in Section 4.3.5. Also, Fan and Gijbels (1996) and Ruppert, Sheather, and Wand (1995) developed data-driven bandwidth selection schemes based on asymptotic formulas for the optimal bandwidths, which are less variable and more effective than the conventional data-driven bandwidth selectors such as the cross-validation bandwidth rule. Similar algorithms can be developed for the estimation of functional-coefficient models based on (4.23); however, this will be a future research topic.

Cai, Fan and Yao (2000) proposed a simple and quick method for selecting bandwidth h . It can be regarded as a modified multi-fold cross-validation criterion that is attentive to the structure of stationary time series data. Let m and Q be two given positive integers and $n > mQ$. The basic idea is first to use Q subseries of lengths $n - qm$ ($q = 1, \dots, Q$) to estimate the unknown coefficient functions and then compute the one-step forecasting errors of the next section of the time series of length m based on the estimated models. More precisely, we choose h that minimizes the average mean squared (AMS) error

$$\text{AMS}(h) = \sum_{q=1}^Q \text{AMS}_q(h), \quad (4.14)$$

where for $q = 1, \dots, Q$,

$$\text{AMS}_q(h) = \frac{1}{m} \sum_{i=n-qm+1}^{n-qm+m} \left\{ Y_i - \sum_{j=1}^p \hat{a}_{j,q}(U_i) X_{i,j} \right\}^2,$$

and $\{\hat{a}_{j,q}(\cdot)\}$ are computed from the sample $\{(U_i, \mathbf{X}_i, Y_i), 1 \leq i \leq n - qm\}$ with bandwidth equal $h[n/(n - qm)]^{1/5}$. Note that we re-scale bandwidth h for different sample sizes according to its optimal rate, i.e. $h \propto n^{-1/5}$. In practical implementations, we may use $m = [0.1n]$ and $Q = 4$. The selected bandwidth does not depend critically on the choice of m and Q , as long as mQ is reasonably large so that the evaluation of prediction errors is stable. A weighted version of $\text{AMS}(h)$ can be used, if one wishes to down-weight the prediction errors at an earlier time. We believe that this bandwidth should be good for modeling and forecasting for time series.

4.5.4 Smoothing Variable Selection

Of importance is to choose an appropriate smoothing variable U in applying functional-coefficient regression models if U is a lagged variable. Knowledge on physical background of the data may be very helpful, as Cai, Fan and Yao (2000) discussed in modeling the lynx data. Without any prior information, it is pertinent to choose U in terms of some data-driven methods such as the Akaike information criterion (AIC) and its variants, cross-validation, and other criteria. Ideally, we would choose U as a linear function of given explanatory variables according to some optimal criterion, which can be fully explored in the work by Fan, Yao and Cai (2003). Nevertheless, we propose here a simple and practical approach: let U be one of the given explanatory variables such that AMS defined in (4.14) obtains its minimum value. Obviously, this idea can be also extended to select p (number of lags) as well.

4.5.5 Goodness-of-Fit Test

To test whether model (4.9) holds with a specified parametric form which is popular in economic and financial applications, such as the threshold autoregressive (TAR) models

$$a_j(u) = \begin{cases} a_{j1}, & \text{if } u \leq \eta \\ a_{j2}, & \text{if } u > \eta, \end{cases}$$

or generalized exponential autoregressive (EXPAR) models

$$a_j(u) = \alpha_j + (\beta_j + \gamma_j u) \exp(-\theta_j u^2),$$

or smooth transition autoregressive (STAR) models

$$a_j(u) = [1 - \exp(-\theta_j u)]^{-1} \quad (\text{logistic}),$$

or

$$a_j(u) = 1 - \exp(-\theta_j u^2) \quad (\text{exponential}),$$

or

$$a_j(u) = [1 - \exp(-\theta_j |u|)]^{-1} \quad (\text{absolute}),$$

[for more discussions on those models, please see the survey paper by van Dijk, Teräsvirta and Franses (2002)], we propose a goodness-of-fit test based on the comparison of the residual sum of squares (RSS) from both parametric and nonparametric fittings. This method is closely

related to the sieve likelihood method proposed by Fan, Zhang and Zhang (2001). Those authors demonstrated the optimality of this kind of procedures for independent samples.

Consider the null hypothesis

$$H_0 : a_j(u) = \alpha_j(u, \boldsymbol{\theta}), \quad 1 \leq j \leq p, \quad (4.15)$$

where $\alpha_j(\cdot, \boldsymbol{\theta})$ is a given family of functions indexed by unknown parameter vector $\boldsymbol{\theta}$. Let $\hat{\boldsymbol{\theta}}$ be an estimator of $\boldsymbol{\theta}$. The RSS under the null hypothesis is

$$\text{RSS}_0 = n^{-1} \sum_{i=1}^n \left\{ Y_i - \alpha_1(U_i, \hat{\boldsymbol{\theta}})X_{i1} - \cdots - \alpha_p(U_i, \hat{\boldsymbol{\theta}})X_{ip} \right\}^2.$$

Analogously, the RSS corresponding to model (4.9) is

$$\text{RSS}_1 = n^{-1} \sum_{i=1}^n \{ Y_i - \hat{a}_1(U_i)X_{i1} - \cdots - \hat{a}_p(U_i)X_{ip} \}^2.$$

The test statistic is defined as

$$T_n = (\text{RSS}_0 - \text{RSS}_1) / \text{RSS}_1 = \text{RSS}_0 / \text{RSS}_1 - 1,$$

and we reject the null hypothesis (4.15) for large value of T_n . We use the following nonparametric bootstrap approach to evaluate the p value of the test:

1. Generate the bootstrap residuals $\{\varepsilon_i^*\}_{i=1}^n$ from the empirical distribution of the centered residuals $\{\hat{\varepsilon}_i - \bar{\hat{\varepsilon}}\}_{i=1}^n$, where

$$\hat{\varepsilon}_i = Y_i - \hat{a}_1(U_i)X_{i1} - \cdots - \hat{a}_p(U_i)X_{ip}, \quad \bar{\hat{\varepsilon}} = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i,$$

and define

$$Y_i^* = \alpha_1(U_i, \hat{\boldsymbol{\theta}})X_{i1} + \cdots + \alpha_p(U_i, \hat{\boldsymbol{\theta}})X_{ip} + \varepsilon_i^*.$$

2. Calculate the bootstrap test statistic T_n^* based on the sample $\{U_i, \mathbf{X}_i, Y_i^*\}_{i=1}^n$.
3. Reject the null hypothesis H_0 when T_n is greater than the upper- α point of the conditional distribution of T_n^* given $\{U_i, \mathbf{X}_i, Y_i\}_{i=1}^n$.

The p -value of the test is simply the relative frequency of the event $\{T_n^* \geq T_n\}$ in the replications of the bootstrap sampling. For the sake of simplicity, we use the same bandwidth

in calculating T_n^* as that in T_n . Note that we bootstrap the centralized residuals from the nonparametric fit instead of the parametric fit, because the nonparametric estimate of residuals is always consistent, no matter whether the null or the alternative hypothesis is correct. The method should provide a consistent estimator of the null distribution even when the null hypothesis does not hold. Kreiss, Neumann, and Yao (1998) considered nonparametric bootstrap tests in a general nonparametric regression setting. They proved that, asymptotically, the conditional distribution of the bootstrap test statistic is indeed the distribution of the test statistic under the null hypothesis. It may be proven that the similar result holds here as long as $\hat{\boldsymbol{\theta}}$ converges to $\boldsymbol{\theta}$ at the rate $n^{-1/2}$.

It is a great challenge to derive the asymptotic property of the testing statistics T_n under time series context and general assumptions. That is to show that

$$b_n [T_n - \lambda_n] \rightarrow N(0, \sigma^2)$$

for some b_n and λ_n , which is a great project for future research. Note that Fan, Zhang and Zhang (2001) derived the above result for the iid sample.

4.5.6 Asymptotic Results

We first present a result on mean squared convergence that serves as a building block for our main result and is also of independent interest. We now introduce some notation. Let

$$\mathbf{S}_n = \mathbf{S}_n(u_0) = \begin{pmatrix} \mathbf{S}_{n,0} & \mathbf{S}_{n,1} \\ \mathbf{S}_{n,1} & \mathbf{S}_{n,2} \end{pmatrix}$$

and

$$\mathbf{T}_n = \mathbf{T}_n(u_0) = \begin{pmatrix} \mathbf{T}_{n,0}(u_0) \\ \mathbf{T}_{n,1}(u_0) \end{pmatrix}$$

with

$$\mathbf{S}_{n,j} = \mathbf{S}_{n,j}(u_0) = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \left(\frac{U_i - u_0}{h} \right)^j K_h(U_i - u_0)$$

and

$$\mathbf{T}_{n,j}(u_0) = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \left(\frac{U_i - u_0}{h} \right)^j K_h(U_i - u_0) Y_i. \quad (4.16)$$

Then, the solution to (4.11) can be expressed as

$$\hat{\boldsymbol{\beta}} = \mathbf{H}^{-1} \mathbf{S}_n^{-1} \mathbf{T}_n, \quad (4.17)$$

where $\mathbf{H} = \text{diag}(1, \dots, 1, h, \dots, h)$ with p -diagonal elements 1's and p diagonal elements h 's. To facilitate the notation, we denote

$$\mathbf{\Omega} = (\omega_{l,m})_{p \times p} = E(\mathbf{X} \mathbf{X}^T | U = u_0). \quad (4.18)$$

Also, let $f(u, \mathbf{x})$ denote the joint density of (U, \mathbf{X}) and $f_u(u)$ be the marginal density of U . We use the following convention: if $U = X_{j_0}$ for some $1 \leq j_0 \leq p$, then $f(u, \mathbf{x})$ becomes $f(\mathbf{x})$ the joint density of \mathbf{X} .

Theorem 4.1. Let condition A.1 in hold, and let $f(u, \mathbf{x})$ be continuous at the point u_0 . Let $h_n \rightarrow 0$ and $n h_n \rightarrow \infty$, as $n \rightarrow \infty$. Then it holds that

$$E(\mathbf{S}_{n,j}(u_0)) \rightarrow f_u(u_0) \mathbf{\Omega}(u_0) \mu_j,$$

and

$$n h_n \text{Var}(\mathbf{S}_{n,j}(u_0)_{l,m}) \rightarrow f_u(u_0) \nu_{2j} \omega_{l,m}$$

for each $0 \leq j \leq 3$ and $1 \leq l, m \leq p$.

As a consequence of Theorem 4.1, we have

$$\mathbf{S}_n \xrightarrow{\mathcal{P}} f_u(u_0) \mathbf{S}, \quad \text{and} \quad \mathbf{S}_{n,3} \xrightarrow{\mathcal{P}} \mu_3 f_u(u_0) \mathbf{\Omega}$$

in the sense that each element converges in probability, where

$$\mathbf{S} = \begin{pmatrix} \mathbf{\Omega} & \mu_1 \mathbf{\Omega} \\ \mu_1 \mathbf{\Omega} & \mu_2 \mathbf{\Omega} \end{pmatrix}.$$

Put

$$\sigma^2(u, \mathbf{x}) = \text{Var}(Y | U = u, \mathbf{X} = \mathbf{x}) \quad (4.19)$$

and

$$\mathbf{\Omega}^*(u_0) = E[\mathbf{X} \mathbf{X}^T \sigma^2(U, \mathbf{X}) | U = u_0]. \quad (4.20)$$

Let $c_0 = \mu_2 / (\mu_2 - \mu_1^2)$ and $c_1 = -\mu_1 / (\mu_2 - \mu_1^2)$.

Theorem 4.2. Let $\sigma^2(u, \mathbf{x})$ and $f(u, \mathbf{x})$ be continuous at the point u_0 . Then under conditions A.1 and A.2,

$$\sqrt{n h_n} \left[\hat{\mathbf{a}}(u_0) - \mathbf{a}(u_0) - \frac{h^2}{2} \frac{\mu_2^2 - \mu_1 \mu_3}{\mu_2 - \mu_1^2} \mathbf{a}''(u_0) \right] \xrightarrow{\mathcal{D}} N(0, \Theta^2(u_0)), \quad (4.21)$$

provided that $f_u(u_0) \neq 0$, where

$$\Theta^2(u_0) = \frac{c_0^2 \nu_0 + 2 c_0 c_1 \nu_1 + c_1^2 \nu_2}{f_u(u_0)} \mathbf{\Omega}^{-1}(u_0) \mathbf{\Omega}^*(u_0) \mathbf{\Omega}^{-1}(u_0). \quad (4.22)$$

Theorem 4.2 indicates that the asymptotic bias of $\hat{a}_j(u_0)$ is

$$\frac{h^2}{2} \frac{\mu_2^2 - \mu_1 \mu_3}{\mu_2 - \mu_1^2} a_j''(u_0)$$

and the asymptotic variance is $(n h_n)^{-1} \theta_j^2(u_0)$, where

$$\theta_j^2(u_0) = \frac{c_0^2 \nu_0 + 2 c_0 c_1 \nu_1 + c_1^2 \nu_2}{f_u(u_0)} \mathbf{e}_{j,p}^T \mathbf{\Omega}^{-1}(u_0) \mathbf{\Omega}^*(u_0) \mathbf{\Omega}^{-1}(u_0) \mathbf{e}_{j,p}.$$

When $\mu_1 = 0$, the bias and variance expressions can be simplified as $h^2 \mu_2 a_j''(u_0)/2$ and

$$\theta_j^2(u_0) = \frac{\nu_0}{f_u(u_0)} \mathbf{e}_{j,p}^T \mathbf{\Omega}^{-1}(u_0) \mathbf{\Omega}^*(u_0) \mathbf{\Omega}^{-1}(u_0) \mathbf{e}_{j,p}.$$

The optimal bandwidth for estimating $a_j(\cdot)$ can be defined to be the one that minimizes the squared bias plus variance. The optimal bandwidth is given by

$$h_{j,\text{opt}} = \left[\frac{\mu_2^2 \nu_0 - 2 \mu_1 \mu_2 \nu_1 + \mu_1^2 \nu_2}{f_u(u_0) (\mu_2^2 - \mu_1 \mu_3)^2} \frac{\mathbf{e}_{j,p}^T \mathbf{\Omega}^{-1}(u_0) \mathbf{\Omega}^*(u_0) \mathbf{\Omega}^{-1}(u_0) \mathbf{e}_{j,p}}{\{a_j''(u_0)\}^2} \right]^{1/5} n^{-1/5}. \quad (4.23)$$

4.5.7 Conditions and Proofs

We first impose some conditions on the regression model but they might not be the weakest possible.

Condition A.1

- a. The kernel function $K(\cdot)$ is a bounded density with a bounded support $[-1, 1]$.
- b. $|f(u, v | \mathbf{x}_0, \mathbf{x}_1; l)| \leq M < \infty$, for all $l \geq 1$, where $f(u, v, | \mathbf{x}_0, \mathbf{x}_1; l)$ is the conditional density of (U_0, U_l) given $(\mathbf{X}_0, \mathbf{X}_l)$, and $f(u | \mathbf{x}) \leq M < \infty$, where $f(u | \mathbf{x})$ is the conditional density of U given $\mathbf{X} = \mathbf{x}$.
- c. The process $\{U_i, \mathbf{X}_i, Y_i\}$ is α -mixing with $\sum k^c [\alpha(k)]^{1-2/\delta} < \infty$ for some $\delta > 2$ and $c > 1 - 2/\delta$.
- d. $E|\mathbf{X}|^{2\delta} < \infty$, where δ is given in condition A.1c.

Condition A.2

a. Assume that

$$E \{Y_0^2 + Y_l^2 \mid U_0 = u, \mathbf{X}_0 = \mathbf{x}_0; U_l = v, \mathbf{X}_l = \mathbf{x}_l\} \leq M < \infty, \quad (4.24)$$

for all $l \geq 1$, $\mathbf{x}_0, \mathbf{x}_1 \in \mathfrak{R}^p$, u , and v in a neighborhood of u_0 .

b. Assume that $h_n \rightarrow 0$ and $n h_n \rightarrow \infty$. Further, assume that there exists a sequence of positive integers s_n such that $s_n \rightarrow \infty$, $s_n = o((n h_n)^{1/2})$, and $(n/h_n)^{1/2} \alpha(s_n) \rightarrow 0$, as $n \rightarrow \infty$.

c. There exists $\delta^* > \delta$, where δ is given in Condition A.1c, such that

$$E \{|Y|^{\delta^*} \mid U = u, \mathbf{X} = \mathbf{x}\} \leq M_4 < \infty \quad (4.25)$$

for all $\mathbf{x} \in \mathfrak{R}^p$ and u in a neighborhood of u_0 , and

$$\alpha(n) = O(n^{-\theta^*}), \quad (4.26)$$

where $\theta^* \geq \delta \delta^* / \{2(\delta^* - \delta)\}$.

d. $E|\mathbf{X}|^{2\delta^*} < \infty$, and $n^{1/2-\delta/4} h^{\delta/\delta^*-1/2-\delta/4} = O(1)$.

Remark A.1. We provide a sufficient condition for the mixing coefficient $\alpha(n)$ to satisfy conditions A.1c and A.2b. Suppose that $h_n = A n^{-\rho}$ ($0 < \rho < 1$, $A > 0$), $s_n = (n h_n / \log n)^{1/2}$ and $\alpha(n) = O(n^{-d})$ for some $d > 0$. Then condition A.1c is satisfied for $d > 2(1 - 1/\delta)/(1 - 2/\delta)$ and condition A.2b is satisfied if $d > (1 + \rho)/(1 - \rho)$. Hence both conditions are satisfied if

$$\alpha(n) = O(n^{-d}), \quad d > \max \left\{ \frac{1 + \rho}{1 - \rho}, \frac{2(1 - 1/\delta)}{1 - 2/\delta} \right\}.$$

Note that this is a trade-off between the order δ of the moment of Y and the rate of decay of the mixing coefficient; the larger the order δ , the weaker the decay rate of $\alpha(n)$.

To study the joint asymptotic normality of $\hat{\mathbf{a}}(u_0)$, we need to center the vector $\mathbf{T}_n(u_0)$ by replacing Y_i with $Y_i - m(U_i, \mathbf{X}_i)$ in the expression (4.16) of $\mathbf{T}_{n,j}(u_0)$. Let

$$\mathbf{T}_{n,j}^*(u_0) = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \left(\frac{U_i - u_0}{h} \right)^j K_h(U_i - u_0) [Y_i - m(U_i, \mathbf{X}_i)],$$

and

$$\mathbf{T}_n^* = \begin{pmatrix} \mathbf{T}_{n,0}^* \\ \mathbf{T}_{n,1}^* \end{pmatrix}.$$

Because the coefficient functions $a_j(u)$ are conducted in the neighborhood of $|U_i - u_0| < h$, by Taylor's expansion,

$$m(U_i, \mathbf{X}_i) = \mathbf{X}_i^T \mathbf{a}(u_0) + (U_i - u_0) \mathbf{X}_i^T \mathbf{a}'(u_0) + \frac{h^2}{2} \left(\frac{U_i - u_0}{h} \right)^2 \mathbf{X}_i^T \mathbf{a}''(u_0) + o_p(h^2),$$

where $\mathbf{a}'(u_0)$ and $\mathbf{a}''(u_0)$ are the vectors consisting of the first and second derivatives of the functions $a_j(\cdot)$. Then,

$$\mathbf{T}_{n,0} - \mathbf{T}_{n,0}^* = \mathbf{S}_{n,0} \mathbf{a}(u_0) + h \mathbf{S}_{n,1} \mathbf{a}'(u_0) + \frac{h^2}{2} \mathbf{S}_{n,2} \mathbf{a}''(u_0) + o_p(h^2)$$

and

$$\mathbf{T}_{n,1} - \mathbf{T}_{n,1}^* = \mathbf{S}_{n,1} \mathbf{a}(u_0) + h \mathbf{S}_{n,2} \mathbf{a}'(u_0) + \frac{h^2}{2} \mathbf{S}_{n,3} \mathbf{a}''(u_0) + o_p(h^2),$$

so that

$$\mathbf{T}_n - \mathbf{T}_n^* = \mathbf{S}_n \mathbf{H} \boldsymbol{\beta} + \frac{h^2}{2} \begin{pmatrix} \mathbf{S}_{n,2} \\ \mathbf{S}_{n,3} \end{pmatrix} \mathbf{a}''(u_0) + o_p(h^2), \quad (4.27)$$

where $\boldsymbol{\beta} = (\mathbf{a}(u_0)^T, \mathbf{a}'(u_0)^T)^T$. Thus it follows from (4.17), (4.27), and Theorem .1 that

$$\mathbf{H} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = f_u^{-1}(u_0) \mathbf{S}^{-1} \mathbf{T}_n^* + \frac{h^2}{2} \mathbf{S}^{-1} \begin{pmatrix} \mu_2 \boldsymbol{\Omega} \\ \mu_3 \boldsymbol{\Omega} \end{pmatrix} \mathbf{a}''(u_0) + o_p(h^2), \quad (4.28)$$

from which the bias term of $\hat{\boldsymbol{\beta}}(u_0)$ is evident. Clearly,

$$\hat{\mathbf{a}}(u_0) - \mathbf{a}(u_0) = \frac{\boldsymbol{\Omega}^{-1}}{f_u(u_0) (\mu_2 - \mu_1^2)} [\mu_2 \mathbf{T}_{n,0}^* - \mu_1 \mathbf{T}_{n,1}^*] + \frac{h^2}{2} \frac{\mu_2^2 - \mu_1 \mu_3}{\mu_2 - \mu_1^2} \mathbf{a}''(u_0) + o_p(h^2). \quad (4.29)$$

Thus (4.29) indicates that the asymptotic bias of $\hat{\mathbf{a}}(u_0)$ is

$$\frac{h^2}{2} \frac{\mu_2^2 - \mu_1 \mu_3}{\mu_2 - \mu_1^2} \mathbf{a}''(u_0).$$

Let

$$\mathbf{Q}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{Z}_i, \quad (4.30)$$

where

$$\mathbf{Z}_i = \mathbf{X}_i \left[c_0 + c_1 \left(\frac{U_i - u_0}{h} \right) \right] K_h(U_i - u_0) [Y_i - m(U_i, \mathbf{X}_i)] \quad (4.31)$$

with $c_0 = \mu_2 / (\mu_2 - \mu_1^2)$ and $c_1 = -\mu_1 / (\mu_2 - \mu_1^2)$. It follows from (4.29) and (4.30) that

$$\sqrt{n} h_n \left[\hat{\mathbf{a}}(u_0) - \mathbf{a}(u_0) - \frac{h^2}{2} \frac{\mu_2^2 - \mu_1 \mu_3}{\mu_2 - \mu_1^2} \mathbf{a}''(u_0) \right] = \frac{\boldsymbol{\Omega}^{-1}}{f_u(u_0)} \sqrt{n} h_n \mathbf{Q}_n + o_p(1). \quad (4.32)$$

We need the following lemma, whose proof is more involved than that for Theorem 4.1. Therefore, we prove only this lemma. Throughout this section, we let C denote a generic constant, which may take different values at different places.

Lemma 4.1. Under conditions A.1 and A.2 and the assumption that $h_n \rightarrow 0$ and $n h_n \rightarrow \infty$, as $n \rightarrow \infty$, if $\sigma^2(u, \mathbf{x})$ and $f(u, \mathbf{x})$ are continuous at the point u_0 , then we have

- (a) $h_n \text{Var}(\mathbf{Z}_1) \rightarrow f_u(u_0) \boldsymbol{\Omega}^*(u_0) [c_0^2 \nu_0 + 2 c_0 c_1 \nu_1 + c_1^2 \nu_2]$;
- (b) $h_n \sum_{l=1}^{n-1} |\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{l+1})| = o(1)$; and
- (c) $n h_n \text{Var}(\mathbf{Q}_n) \rightarrow f_u(u_0) \boldsymbol{\Omega}^*(u_0) [c_0^2 \nu_0 + 2 c_0 c_1 \nu_1 + c_1^2 \nu_2]$.

Proof: First, by conditioning on (U_1, \mathbf{X}_1) and using Theorem 1 of Sun (1984), we have

$$\begin{aligned} \text{Var}(\mathbf{Z}_1) &= E \left[\mathbf{X}_1 \mathbf{X}_1^T \sigma^2(U_1, \mathbf{X}_1) \left\{ c_0 + c_1 \left(\frac{U_1 - u_0}{h} \right) \right\}^2 K_h^2(U_1 - u_0) \right] \\ &= \frac{1}{h} [f_u(u_0) \boldsymbol{\Omega}^*(u_0) \{c_0^2 \nu_0 + 2 c_0 c_1 \nu_1 + c_1^2 \nu_2\} + o(1)]. \end{aligned} \quad (4.33)$$

The result (c) follows in an obvious manner from (a) and (b) along with

$$\text{Var}(\mathbf{Q}_n) = \frac{1}{n} \text{Var}(\mathbf{Z}_1) + \frac{2}{n} \sum_{l=1}^{n-1} \left(1 - \frac{l}{n}\right) \text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{l+1}). \quad (4.34)$$

It thus remains to prove part (b). To this end, let $d_n \rightarrow \infty$ be a sequence of positive integers such that $d_n h_n \rightarrow 0$. Define

$$J_1 = \sum_{l=1}^{d_n-1} |\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{l+1})| \quad \text{and} \quad J_2 = \sum_{l=d_n}^{n-1} |\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{l+1})|.$$

It remains to show that $J_1 = o(h^{-1})$ and $J_2 = o(h^{-1})$.

We remark that because $K(\cdot)$ has a bounded support $[-1, 1]$, $a_j(u)$ is bounded in the neighborhood of $u \in [u_0 - h, u_0 + h]$. Let $B = \max_{1 \leq j \leq p} \sup_{|u-u_0| < h} |a_j(u)|$ and $g(\mathbf{x}) = \sum_{j=1}^p |x_j|$. Then $\sup_{|u-u_0| < h} |m(u, \mathbf{x})| \leq B g(\mathbf{x})$. By conditioning on (U_1, \mathbf{X}_1) and

$(U_{l+1}, \mathbf{X}_{l+1})$, and using (4.24) and condition A.1b, we have, for all $l \geq 1$,

$$\begin{aligned}
& |\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{l+1})| \\
& \leq C E [|\mathbf{X}_1 \mathbf{X}_{l+1}^T| \{ |Y_1| + B g(\mathbf{X}_1) \} \{ |Y_{l+1}| + B g(\mathbf{X}_{l+1}) \} K_h(U_1 - u_0) K_h(U_{l+1} - u_0)] \\
& \leq C E \left[|\mathbf{X}_1 \mathbf{X}_{l+1}^T| \{ M_2 + B^2 g^2(\mathbf{X}_1) \}^{1/2} \{ M_2 + B^2 g^2(\mathbf{X}_{l+1}) \}^{1/2} K_h(U_1 - u_0) K_h(U_{l+1} - u_0) \right] \\
& \leq C E [|\mathbf{X}_1 \mathbf{X}_{l+1}^T| \{ 1 + g(\mathbf{X}_1) \} \{ 1 + g(\mathbf{X}_{l+1}) \}] \leq C.
\end{aligned} \tag{4.35}$$

It follows that

$$J_1 \leq C d_n = o(h^{-1})$$

by the choice of d_n . We next consider the upper bound of J_2 . To this end, using Davydov's inequality (see Hall and Heyde 1980, Corollary A.2), we obtain, for all $1 \leq j, m \leq p$ and $l \geq 1$,

$$|\text{Cov}(Z_{1j}, Z_{l+1,m})| \leq C [\alpha(l)]^{1-2/\delta} [E|Z_j|^\delta]^{1/\delta} [E|Z_m|^\delta]^{1/\delta}. \tag{4.36}$$

By conditioning on (U, \mathbf{X}) and using conditions A.1b and A.2c, one has

$$\begin{aligned}
E[|Z_j|^\delta] & \leq C E[|X_j|^\delta K_h^\delta(U - u_0) \{ |Y|^\delta + B^\delta g^\delta(\mathbf{X}) \}] \\
& \leq C E[|X_j|^\delta K_h^\delta(U - u_0) \{ M_3 + B^\delta g^\delta(\mathbf{X}) \}] \\
& \leq C h^{1-\delta} E[|X_j|^\delta \{ M_3 + B^\delta g^\delta(\mathbf{X}) \}] \leq C h^{1-\delta}.
\end{aligned} \tag{4.37}$$

A combination of (4.36) and (4.37) leads to

$$J_2 \leq C h^{2/\delta-2} \sum_{l=d_n}^{\infty} [\alpha(l)]^{1-2/\delta} \leq C h^{2/\delta-2} d_n^{-c} \sum_{l=d_n}^{\infty} l^c [\alpha(l)]^{1-2/\delta} = o(h^{-1}) \tag{4.38}$$

by choosing d_n such that $h^{1-2/\delta} d_n^c = C$, so the requirement that $d_n h_n \rightarrow 0$ is satisfied.

Proof of Theorem 4.2

We use the small-block and large-block technique – namely, partition $\{1, \dots, n\}$ into $2q_n + 1$ subsets with large block of size $r = r_n$ and small block of size $s = s_n$. Set

$$q = q_n = \left\lfloor \frac{n}{r_n + s_n} \right\rfloor. \tag{4.39}$$

We now use the Cramér-Wold device to derive the asymptotic normality of \mathbf{Q}_n . For any unit vector $\mathbf{d} \in \mathbb{R}^p$, let $Z_{n,i} = \sqrt{h} \mathbf{d}^T \mathbf{Z}_{i+1}$, $i = 0, \dots, n-1$. Then

$$\sqrt{n} h \mathbf{d}^T \mathbf{Q}_n = \frac{1}{\sqrt{n}} \sum_{i=0}^{n-1} Z_{n,i},$$

and, by Lemma 4.1,

$$\begin{aligned} \text{Var}(Z_{n,0}) &\approx f_u(u_0) \mathbf{d}^T \boldsymbol{\Omega}^*(u_0) \mathbf{d} [c_0^2 \nu_0 + 2 c_0 c_1 \nu_1 + c_1^2 \nu_2] \\ &\equiv \theta^2(u_0) \end{aligned} \quad (4.40)$$

and

$$\sum_{l=0}^{n-1} |\text{Cov}(Z_{n,0}, Z_{n,l})| = o(1). \quad (4.41)$$

Define the random variables, for $0 \leq j \leq q-1$,

$$\eta_j = \sum_{i=j(r+s)}^{j(r+s)+r-1} Z_{n,i}, \quad \xi_j = \sum_{i=j(r+s)+r}^{(j+1)(r+s)} Z_{n,i}, \quad \text{and} \quad \zeta_q = \sum_{i=q(r+s)}^{n-1} Z_{n,i}.$$

Then,

$$\sqrt{n} h \mathbf{d}^T \mathbf{Q}_n = \frac{1}{\sqrt{n}} \left\{ \sum_{j=0}^{q-1} \eta_j + \sum_{j=0}^{q-1} \xi_j + \zeta_q \right\} \equiv \frac{1}{\sqrt{n}} \{Q_{n,1} + Q_{n,2} + Q_{n,3}\}. \quad (4.42)$$

We show that as $n \rightarrow \infty$,

$$\frac{1}{n} E [Q_{n,2}]^2 \rightarrow 0, \quad \frac{1}{n} E [Q_{n,3}]^2 \rightarrow 0, \quad (4.43)$$

$$\left| E [\exp(it Q_{n,1})] - \prod_{j=0}^{q-1} E [\exp(it \eta_j)] \right| \rightarrow 0, \quad (4.44)$$

$$\frac{1}{n} \sum_{j=0}^{q-1} E (\eta_j^2) \rightarrow \theta^2(u_0), \quad (4.45)$$

and

$$\frac{1}{n} \sum_{j=0}^{q-1} E [\eta_j^2 I \{|\eta_j| \geq \varepsilon \theta(u_0) \sqrt{n}\}] \rightarrow 0 \quad (4.46)$$

for every $\varepsilon > 0$. (4.43) implies that $Q_{n,2}$ and $Q_{n,3}$ are asymptotically negligible in probability, (4.44) shows that the summands η_j in $Q_{n,1}$ are asymptotically independent and (4.45) and (4.46) are the standard Lindeberg-Feller conditions for asymptotic normality of $Q_{n,1}$ for the independent setup.

We first establish (4.43). For this purpose, we choose the large block size. Condition A.2b implies that there is a sequence of positive constants $\gamma_n \rightarrow \infty$ such that $\gamma_n s_n = o(\sqrt{n} h_n)$ and

$$\gamma_n (n/h_n)^{1/2} \alpha(s_n) \rightarrow 0. \quad (4.47)$$

Define the large block size r_n by $r_n = \lfloor (n h_n)^{1/2} / \gamma_n \rfloor$ and the small block size s_n . Then it can easily be shown from (4.47) that as $n \rightarrow \infty$,

$$s_n/r_n \rightarrow 0, \quad r_n/n \rightarrow 0, \quad r_n (n h_n)^{-1/2} \rightarrow 0, \quad (4.48)$$

and

$$(n/r_n) \alpha(s_n) \rightarrow 0. \quad (4.49)$$

Observe that

$$E[Q_{n,2}]^2 = \sum_{j=0}^{q-1} \text{Var}(\xi_j) + 2 \sum_{0 \leq i < j \leq q-1} \text{Cov}(\xi_i, \xi_j) \equiv I_1 + I_2. \quad (4.50)$$

It follows from stationarity and Lemma 4.1 that

$$I_1 = q_n \text{Var}(\xi_1) = q_n \text{Var} \left(\sum_{j=1}^{s_n} Z_{n,j} \right) = q_n s_n [\theta^2(u_0) + o(1)]. \quad (4.51)$$

Next consider the second term I_2 in the right side of (4.50). Let $r_j^* = j(r_n + s_n)$, then $r_j^* - r_i^* \geq r_n$ for all $j > i$, we thus have

$$\begin{aligned} |I_2| &\leq 2 \sum_{0 \leq i < j \leq q-1} \sum_{j_1=1}^{s_n} \sum_{j_2=1}^{s_n} |\text{Cov}(Z_{n,r_i^*+r_n+j_1}, Z_{n,r_j^*+r_n+j_2})| \\ &\leq 2 \sum_{j_1=1}^{n-r_n} \sum_{j_2=j_1+r_n}^n |\text{Cov}(Z_{n,j_1}, Z_{n,j_2})|. \end{aligned}$$

By stationarity and Lemma 4.1, one obtains

$$|I_2| \leq 2n \sum_{j=r_n+1}^n |\text{Cov}(Z_{n,1}, Z_{n,j})| = o(n). \quad (4.52)$$

Hence, by (4.48)-(4.52), we have

$$\frac{1}{n} E[Q_{n,2}]^2 = O(q_n s_n n^{-1}) + o(1) = o(1). \quad (4.53)$$

It follows from stationarity, (4.48), and Lemma 4.1 that

$$\text{Var}[Q_{n,3}] = \text{Var} \left(\sum_{j=1}^{n-q_n(r_n+s_n)} Z_{n,j} \right) = O(n - q_n(r_n + s_n)) = o(n). \quad (4.54)$$

Combining (4.48), (4.53), and (4.54), we establish (4.43). As for (4.45), by stationarity, (4.48), (4.49), and Lemma 4.1, it is easily seen that

$$\frac{1}{n} \sum_{j=0}^{q_n-1} E(\eta_j^2) = \frac{q_n}{n} E(\eta_1^2) = \frac{q_n r_n}{n} \cdot \frac{1}{r_n} \text{Var} \left(\sum_{j=1}^{r_n} Z_{n,j} \right) \rightarrow \theta^2(u_0).$$

To establish (4.44), we use Lemma 1.1 of Volkonskii and Rozanov (1959) (see also Ibragimov and Linnik 1971, p. 338) to obtain

$$\left| E [\exp(it Q_{n,1})] - \prod_{j=0}^{q_n-1} E [\exp(it \eta_j)] \right| \leq 16 (n/r_n) \alpha(s_n)$$

tending to 0 by (4.49).

It remains to establish (4.46). For this purpose, we use theorem 4.1 of Shao and Yu (1996) and condition A.2 to obtain

$$E [\eta_1^2 I \{|\eta_1| \geq \varepsilon \theta(u_0) \sqrt{n}\}] \leq C n^{1-\delta/2} E (|\eta_1|^\delta) \leq C n^{1-\delta/2} r_n^{\delta/2} \{E (|Z_{n,0}|^{\delta^*})\}^{\delta/\delta^*}. \quad (4.55)$$

As in (4.37),

$$E (|Z_{n,0}|^{\delta^*}) \leq C h^{1-\delta^*/2}. \quad (4.56)$$

Therefore, by (4.55) and (4.56),

$$E [\eta_1^2 I \{|\eta_1| \geq \varepsilon \theta(u_0) \sqrt{n}\}] \leq C n^{1-\delta/2} r_n^{\delta/2} h^{(2-\delta^*)\delta/(2\delta^*)}. \quad (4.57)$$

Thus, by (4.39) and the definition of r_n , and using conditions A.2c and A.2d, we obtain

$$\frac{1}{n} \sum_{j=0}^{q-1} E [\eta_j^2 I \{|\eta_j| \geq \varepsilon \theta(u_0) \sqrt{n}\}] \leq C \gamma_n^{1-\delta/2} n^{1/2-\delta/4} h_n^{\delta/\delta^*-1/2-\delta/4} \rightarrow 0 \quad (4.58)$$

because $\gamma_n \rightarrow \infty$. This completes the proof of the theorem.

4.5.8 Monte Carlo Simulations and Applications

1. Applications to Time Series

See Cai, Fan and Yao (2000) for the detailed Monte Carlo simulation results and applications.

2. Boston Housing Data

1. Description of Data

The well known **Boston house price data** set¹ consists of 14 variables, collected on each of 506 different houses from a variety of locations. The Boston house-price data set was used originally by Harrison and Rubinfeld (1978) and it was re-analyzed in Belsley, Kuh and Welsch (1980) by various transformations in the table on pages 244-261. Variables are, denoted by X_1, \dots, X_{13} and Y , in order:

¹This dataset can be downloaded from the web site at <http://lib.stat.cmu.edu/datasets/boston>.

CRIM	per capita crime rate by town
ZN	proportion of residential land zoned for lots over 25,000 sq.ft.
INDUS	proportion of non-retail business acres per town
CHAS	Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
NOX	nitric oxides concentration (parts per 10 million)
RM	average number of rooms per dwelling
AGE	proportion of owner-occupied units built prior to 1940
DIS	weighted distances to five Boston employment centers
RAD	index of accessibility to radial highways
TAX	full-value property-tax rate per 10,000USD
PTRATIO	pupil-teacher ratio by town
B	$1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town
LSTAT	lower status of the population
MEDV	Median value of owner-occupied homes in \$1000's

The dependent variable is Y , the median value of owner-occupied homes in \$1,000's (house price). The major factors possibly affecting the house prices used in the literature are: X_{13} =proportion of population of lower educational status X_6 =the average number of rooms per house, X_1 =the per capita crime rate, X_{10} =the full property tax rate, and X_{11} =the pupil/teacher ratio. For the complete description of all 14 variables, see Harrison and Rubinfeld (1978) and Gilley and Pace (1996) for corrections.

2. Linear Models

Harrison and Rubinfeld (1978) was the first to analyze this data set using a standard regression model Y versus all 13 variables including some higher order terms or transformations on Y and X_j 's. The purpose of this study is to see whether there are the **effects of pollution** on housing prices via **hedonic pricing methodology**. Belsley, Kuh and Welsch (1980) used this data set to illustrate the effects of using **robust regression** and **outlier detection** strategies. From these results, we might conclude that the model might not be linear and there might exist outliers. Also, Pace and Gilley (1997) added a **georeferencing idea (spatial statistics)** and used a **spatial estimation** method to consider this data set.

Exercise: Please use all possible methods to explore this dataset to see what is the best linear model you can obtain.

3. Fit a Varying-Coefficient Model

Recently, Şentürk and Müller (2003) studied the correlation between the house price Y and the crime rate X_1 adjusted by the confounding variable X_{13} through a varying coefficient model and they concluded that **the expected effect of increasing crime rate on declining house prices seems to be only observed for lower educational status neighborhoods in Boston**. Finally, it is surprising that all the existing nonparametric models aforementioned above did not include the crime rate X_1 , which may be an important factor affecting the housing price, and did not consider the interaction terms such as X_{13} and X_1 .

See the paper by Fan and Huang (2005) for fitting a varying coefficient model to the Boston housing data.

Exercise: Please fit a a varying coefficient model to the Boston housing data.

4.6 Additive Model

4.6.1 Model

In this section, we use the notation from Cai (2002). Let $\{\mathbf{X}_t, \mathbf{Y}_t, Z_t\}_{t=-\infty}^{\infty}$ be jointly stationary processes, where \mathbf{X}_t and \mathbf{Y}_t take values in \mathfrak{R}^p and \mathfrak{R}^q with $p, q \geq 0$, respectively. The regression surface is defined by

$$m(\mathbf{x}, \mathbf{y}) = E \{Z_t | \mathbf{X}_t = \mathbf{x}, \mathbf{Y}_t = \mathbf{y}\}. \quad (4.59)$$

Here, it is assumed that $E|Z_t| < \infty$. Note that the regression function $m(\cdot, \cdot)$ defined in (4.59) can identify only the sum

$$m(\mathbf{x}, \mathbf{y}) = \mu + g_1(\mathbf{x}) + g_2(\mathbf{y}). \quad (4.60)$$

Such a decomposition holds, for example, for the following nonlinear additive autoregressive model with exogenous variables (ARX)

$$Y_t = \mu + g_1(X_{t-j_1}, \dots, X_{t-j_p}) + g_2(Y_{t-i_1}, \dots, Y_{t-i_q}) + \eta_t,$$

$$X_{t-j_1} = g_3(X_{t-j_2}, \dots, X_{t-j_p}) + \varepsilon_t.$$

For detailed discussions on the ARX model, the reader is referred to the papers by Masry and Tjøstheim (1997) and Cai and Masry (2000). For identifiability, it is assumed that $E\{g_1(\mathbf{X}_t)\} = 0$ and $E\{g_2(\mathbf{Y}_t)\} = 0$. Then, the projection of $m(\mathbf{x}, \mathbf{y})$ on the $g_1(\mathbf{x})$ -direction is defined by

$$E\{m(\mathbf{x}, \mathbf{Y}_t)\} = \mu + g_1(\mathbf{x}) + E\{g_2(\mathbf{Y}_t)\} = \mu + g_1(\mathbf{x}). \quad (4.61)$$

Clearly, $g_1(\cdot)$ can be identified up to an additive constant and $g_2(\cdot)$ can be retrieved likewise.

A thorough discussion of additive time series models defined in (4.60) can be found in Chen and Tsay (1993). Additive components can be estimated with a one-dimensional nonparametric rate. In most papers, to estimate additive components, several methods have been proposed. For example, Chen and Tsay (1993) used the iterative backfitting procedures, such as the ACE algorithm and the BRUTO approach; see Hastie and Tibshirani (1990) for details. But, their asymptotic properties are not well understood due to the implicit definition of the resulting estimators. To attenuate the drawbacks of iterative procedures, Auestad and Tjøstheim (1991) and Tjøstheim and Auestad (1994a) proposed a direct method based on an average regression surface idea, referred to as projection method in Tjøstheim and Auestad (1994a) for time series data. As pointed out by Cai and Fan (2000), a direct method has some advantages, such as it does not rely on iterations, it can make computation fast, and more importantly, it allows an asymptotic analysis. Finally, the projection method was extended to nonlinear ARX models by Masry and Tjøstheim (1997) using the kernel method and Cai and Masry (2000) coupled with the local polynomial approach. It should be remarked that the projection method, under the name of marginal integration, was proposed independently by Newey (1994) and Linton and Nielsen (1995) for iid samples, and since then, some important progresses have been made by some authors. For example, by combining the marginal integration with one-step backfitting, Linton (1997, 2000) presents an efficient estimator, Mammen, Linton, and Nielsen (1999) established rigorously the asymptotic theory of the backfitting, Cai and Fan (2000) considered estimating each component using the weighted projection method coupled with the local linear fitting in an efficient way, and Sperlich, Tjøstheim, and Yang (2002) extended the efficient method to models with simple

interactions.

The projection method has some disadvantages although it has the aforementioned merits. The projection method may not be efficient if covariates (endogenous or exogenous variables) are strongly correlated, which is particularly relevant for autoregressive models. The intuitive interpretation is that additive components are not orthogonal. To overcome this shortcoming, two efficient estimation methods have been proposed in the literature. The first one is called weight function procedure, proposed by Fan, Härdle, and Mammen (1998) for iid samples and extended to time series situations by Cai and Fan (2000). With an appropriate choice of the weight function, additive components can be efficiently estimated in the sense that an additive component can be estimated with the same asymptotic bias and variance as if the rest of components were known. The second one is to combine the marginal integration with one-step backfitting, introduced by Linton (1997, 2000) for iid samples and extended by Sperlich, Tjøstheim, and Yang (2002) to additive models with single interactions, but this method has not been advocated for time series situations. However, there has not been any attempt to discuss the bandwidth selection for the projection method and its variations in the literature due to their complexity. In practice, one bandwidth is usually used for all components although Cai and Fan (2000) argued that different bandwidths might be used theoretically to deal with the situation that additive components possess the different smoothness. Therefore, the projection method may not be optimal in practice in the sense that one bandwidth is used.

To estimate unknown additive components in (4.60) efficiently, following the spirit of the marginal integration with one-step backfitting proposed by Linton (1997) for iid samples, I use a two-stage method, due to Linton (2000), coupled with the local linear (polynomial) method, which has some attractive properties, such as mathematical efficiency, bias reduction and adaptation of edge effect (see Fan and Gijbels, 1996). The basic idea of the two-stage approach is described as follows. At the first stage, one obtains the initial estimated values for all components. More precisely, the idea for estimating any additive component is first to estimate directly high-dimensional regression surface by the local linear method and then to average the regression surface over the rest of variables to stabilize variance. Such an initial estimate, in general, is under-smoothed so that the bias should be asymptotically negligible. At the second stage, the local linear (polynomial) technique is used again to estimate any additive component by using the initial estimated values of the rest of components. In such

a way, it is shown that the estimate at the second stage is not only efficient in the sense of being equivalent to a procedure based on knowing other components, but also making the bandwidth selection much easier. Note that this technique is not novel to this chapter since the two-stage method is first used by Linton (1997, 2000) for iid samples, but many details and insights are.

4.6.2 Backfitting Algorithm

The building block of the generalized additive model algorithm is the scatterplot smoother. We will first describe scatterplot smoothing in a simple setting, and then indicate how it is used in generalized additive modelling. Here y is a response or outcome variable, and x is a prognostic factor. We wish to fit a smooth curve $f(x)$ that summarizes the dependence of y on x . If we were to find the curve that simply minimizes $\sum_{i=1}^n [y_i - f(x_i)]^2$, the result would be an interpolating curve that would not be smooth at all. The cubic spline smoother imposes smoothness on $f(x)$. We seek the function $f(x)$ that minimizes

$$\sum_{i=1}^n [y_i - f(x_i)]^2 + \lambda \int [f''(x)]^2 dx \quad (4.62)$$

Notice that $\int [f''(x)]^2 dx$ measures the “wiggleness” of the function $f(x)$: linear $f(x)$ s have $\int [f''(x)]^2 dx = 0$, while non-linear fs produce values bigger than zero. λ is a non-negative smoothing parameter that must be chosen by the data analyst. It governs the tradeoff between the goodness of fit to the data and (as measured by and wiggleness of the function. Larger values of λ force $f(x)$ to be smoother.

For any value of λ , the solution to (4.62) is a cubic spline, i.e., a piecewise cubic polynomial with pieces joined at the unique observed values of x in the dataset. Fast and stable numerical procedures are available for computation of the fitted curve. What value of λ did we use in practice? In fact it is not convenient to express the desired smoothness of $f(x)$ in terms of λ , as the meaning of λ depends on the units of the prognostic factor x . Instead, it is possible to define an **“effective number of parameters”** or **“degrees of freedom”** of a cubic spline smoother, and then use a numerical search to determine the value of λ to yield this number. In practice, if we chose the effective number of parameters to be 5, roughly speaking, this means that the complexity of the curve is about the same as a polynomial regression of degrees 4. However, the cubic spline smoother “spreads out” its parameters in

a more even manner, and hence is much more flexible than a polynomial regression. Note that the degrees of freedom of a smoother need not be an integer.

The above discussion tells how to fit a curve to a single prognostic factor. With multiple prognostic factors, if x_{ij} denotes the value of the j th prognostic factor for the i th observation, we fit the additive model

$$y_i = \sum_{j=1}^d f_j(x_{ij}) + \varepsilon_i.$$

A criterion like (4.62) can be specified for this problem, and a simple iterative procedure exists for estimating the f_j s. We apply a cubic spline smoother to the outcome $y_i - \sum_{j \neq k}^d \hat{f}_j(x_{ij})$ as a function of x_{ik} , for each prognostic factor in turn. The process is continued until the estimates $\hat{f}_j(x)$ stabilize. This procedure is known as “backfitting” and the resulting fit is analogous to a multiple regression for linear models.

To fit an additive model or a partially additive model in **R**, the function is **gam()** in the package **gam**. For details, please look at the help command **help(gam)** after loading the package **gam** [**library(gam)**]. Note that the function **gam()** allows to fit a **semi-parametric additive model** as

$$Y = \beta^T \mathbf{X} + \sum_{j=1}^p g_j(Z_j) + \varepsilon,$$

which can be done by specifying some components without smooth.

4.6.3 Projection Method

This section is devoted to a brief review of the projection method and discusses its merits and disadvantages.

It is assumed that all additive components have continuous second partial derivatives, so that $m(\mathbf{u}, \mathbf{v})$ can be locally approximated by a linear term in a neighborhood of (\mathbf{x}, \mathbf{y}) , namely, $m(\mathbf{u}, \mathbf{v}) \approx \beta_0 + \beta_1^T (\mathbf{u} - \mathbf{x}) + \beta_2^T (\mathbf{v} - \mathbf{y})$ with $\{\beta_j\}$ depending on \mathbf{x} and \mathbf{y} , where β_1^T denotes the transpose of β_1 .

Let $K(\cdot)$ and $L(\cdot)$ be symmetric kernel functions in \mathbb{R}^p and \mathbb{R}^q , respectively, and $h_{11} = h_{11}(n) > 0$ and $h_{12} = h_{12}(n) > 0$ be bandwidths in the step of estimating the regression surface. Here, to handle various degrees of smoothness, Cai and Fan (2000) propose using h_{11}

and h_{12} differently although the implementation may not be easy in practice. The reader is referred to the paper by Cai and Fan (2000) for details. Given observations $\{\mathbf{X}_t, \mathbf{Y}_t, Z_t\}_{t=1}^n$, let $\hat{\beta}_j$ be the minimizer of the following locally weighted least squares

$$\sum_{t=1}^n \left\{ Z_t - \beta_0 - \beta_1^T (\mathbf{X}_t - \mathbf{x}) - \beta_2^T (\mathbf{Y}_t - \mathbf{y}) \right\}^2 K_{h_{11}}(\mathbf{X}_t - \mathbf{x}) L_{h_{12}}(\mathbf{Y}_t - \mathbf{y}),$$

where $K_h(\cdot) = K(\cdot/h)/h^p$ and $L_h(\cdot) = L(\cdot/h)/h^q$. Then, the local linear estimator of the regression surface $m(\mathbf{x}, \mathbf{y})$ is $\hat{m}(\mathbf{x}, \mathbf{y}) = \hat{\beta}_0$. By computing the sample average of $\hat{m}(\cdot, \cdot)$ based on (4.61), the projection estimators of $g_1(\cdot)$ and $g_2(\cdot)$ are defined as, respectively,

$$\hat{g}_1(\mathbf{x}) = \frac{1}{n} \sum_{t=1}^n \hat{m}(\mathbf{x}, \mathbf{Y}_t) - \hat{\mu}, \quad \text{and} \quad \hat{g}_2(\mathbf{y}) = \frac{1}{n} \sum_{t=1}^n \hat{m}(\mathbf{X}_t, \mathbf{y}) - \hat{\mu},$$

where $\hat{\mu} = n^{-1} \sum_{t=1}^n Z_t$. Under some regularity conditions, by using the same arguments as those employed in the proof of Theorem 3 in Cai and Masry (2000), it can be shown (although not easy and tedious) that the asymptotic bias and asymptotic variance of $\hat{g}_1(\mathbf{x})$ are, respectively, $h_{11}^2 \text{tr}\{\mu_2(K) g_1''(\mathbf{x})\}/2$ and $v_1(\mathbf{x}) = \nu_0(K) A(\mathbf{x})$, where

$$A(\mathbf{x}) = \int p_2^2(\mathbf{y}) \sigma^2(\mathbf{x}, \mathbf{y}) p^{-1}(\mathbf{x}, \mathbf{y}) d\mathbf{y} \quad \text{and} \quad \sigma^2(\mathbf{x}, \mathbf{y}) = \text{Var}(Z_t | \mathbf{X}_t = \mathbf{x}, \mathbf{Y}_t = \mathbf{y}).$$

Here, $p(\mathbf{x}, \mathbf{y})$ stands for the joint density of \mathbf{X}_t and \mathbf{Y}_t , $p_1(\mathbf{x})$ denotes the marginal density of \mathbf{X}_t , $p_2(\mathbf{y})$ is the marginal density of \mathbf{Y}_t , $\nu_0(K) = \int K^2(\mathbf{u}) d\mathbf{u}$, and $\mu_2(K) = \int \mathbf{u} \mathbf{u}^T K(\mathbf{u}) d\mathbf{u}$.

The foregoing method has some advantages, such as it is easy to understand, it can make computation fast, and it allows an asymptotic analysis. However, it can be quite inefficient in an asymptotic sense. To demonstrate this idea, let us consider the ideal situation that $g_2(\cdot)$ and μ are known. In such a case, one can estimate $g_1(\cdot)$ by directly regressing the partial error $\tilde{Z}_t = Z_t - \mu - g_2(\mathbf{Y}_t)$ on \mathbf{X}_t and such an ideal estimator is optimal in an asymptotic minimax sense (see, e.g., Fan and Gijbels, 1996). The asymptotic bias for the ideal estimator is $h_{11}^2 \text{tr}\{\mu_2(K) g_1''(\mathbf{x})\}/2$ and the asymptotic variance is

$$v_0(\mathbf{x}) = \nu_0(K) B(\mathbf{x}) \quad \text{with} \quad B(\mathbf{x}) = p_1^{-1}(\mathbf{x}) E\{\sigma^2(\mathbf{X}_t, \mathbf{Y}_t) | \mathbf{X}_t = \mathbf{x}\} \quad (4.63)$$

(see, e.g., Masry and Fan, 1997). It is clear that $v_1(\mathbf{x}) = v_0(\mathbf{x})$ if \mathbf{X}_t and \mathbf{Y}_t are independent. If \mathbf{X}_t and \mathbf{Y}_t are correlated and when $\sigma^2(\mathbf{x}, \mathbf{y})$ is a constant, it follows from the Cauchy-Schwarz inequality that

$$B(\mathbf{x}) = \frac{\sigma^2}{p_1(\mathbf{x})} \int p^{1/2}(\mathbf{y} | \mathbf{x}) \frac{p_2(\mathbf{y})}{p^{1/2}(\mathbf{y} | \mathbf{x})} d\mathbf{y} \leq \frac{\sigma^2}{p_1(\mathbf{x})} \int \frac{p_2^2(\mathbf{y})}{p(\mathbf{y} | \mathbf{x})} d\mathbf{y} = A(\mathbf{x}),$$

which implies that the ideal estimator has always smaller asymptotic variance than the projection method although both have the same bias. This suggests that the projection method could lead to an inefficient estimation of $g_1(\cdot)$ and $g_2(\cdot)$ when \mathbf{X}_t and \mathbf{Y}_t are serially correlated, which is particularly relevant for autoregressive models. To alleviate this shortcoming, I propose the two-stage approach described next.

4.6.4 Two-Stage Procedure

The two-stage method due to Linton (1997, 2000) is introduced. The basic idea is to get an initial estimate for $\widehat{g}_2(\cdot)$ using a small bandwidth h_{12} . The initial estimate can be obtained by the projection method and h_{12} can be chosen so small that the bias of estimating $\widehat{g}_2(\cdot)$ can be asymptotically negligible. Then, using the partial residuals $Z_t^* = Z_t - \widehat{\mu} - \widehat{g}_2(\mathbf{Y}_t)$, we apply the local linear regression technique to the pseudo regression model

$$Z_t^* = g_1(\mathbf{X}_t) + \varepsilon_t^*$$

to estimate $g_1(\cdot)$. This leads naturally to the weighted least-squares problem

$$\sum_{t=1}^n \{Z_t^* - \beta_1 - \beta_2^T (\mathbf{X}_t - \mathbf{x})\}^2 J_{h_2}(\mathbf{X}_t - \mathbf{x}), \quad (4.64)$$

where $J(\cdot)$ is the kernel function in \mathbb{R}^p and $h_2 = h_2(n) > 0$ is the bandwidth at the second-stage. The advantage of this is twofold: the bandwidth h_2 can now be selected purposely for estimating $g_1(\cdot)$ only and any bandwidth selection technique for nonparametric regression can be applied here. Maximizing (4.64) with respect to β_1 and β_2 gives the two-stage estimate of $g_1(\mathbf{x})$, denoted by $\widetilde{g}_1(\mathbf{x}) = \widehat{\beta}_1$, where $\widehat{\beta}_1$ and $\widehat{\beta}_2$ are the minimizer of (4.64).

It is shown in Theorem 4.3, in which follows, that under some regularity conditions, the asymptotic bias and variance of the two-stage estimate $\widetilde{g}_1(\mathbf{x})$ are the same as those for the ideal estimator, provided that the initial bandwidth h_{12} satisfies $h_{12} = o(h_2)$.

Sampling Properties

To establish the asymptotic normality of the two-stage estimator, it is assumed that the initial estimator satisfies a linear approximation; namely,

$$\widehat{g}_2(\mathbf{Y}_t) - g_2(\mathbf{Y}_t) \approx \frac{1}{n} \sum_{i=1}^n L_{h_{12}}(\mathbf{Y}_i - \mathbf{Y}_t) \Gamma(\mathbf{X}_i, \mathbf{Y}_t) \delta_i + \frac{1}{2} h_{12}^2 \text{tr}\{\mu_2(L) g_2''(\mathbf{Y}_t)\}, \quad (4.65)$$

where $\delta_t = Z_t - m(\mathbf{X}_t, \mathbf{Y}_t)$ and $\Gamma(\mathbf{x}, \mathbf{y}) = p_1(\mathbf{x})/p(\mathbf{x}, \mathbf{y})$. Note that under some regularity conditions, by following the same arguments as in Masry (1996), one might show (although the proof is not easy, quite lengthy, and tedious) that (4.65) holds. Note that this assumption is also imposed in Linton (2000) for iid samples to simplify the proof of the asymptotic results of the two-stage estimator. Now, the asymptotic normality for the two-stage estimator is stated here and its proof can be found in Cai (2002).

THEOREM 4.3. *Under (4.65) and Assumptions A1 – A9 stated in Cai (2002), if bandwidths h_{12} and h_2 are chosen such that $h_{12} \rightarrow 0$, $n h_{12}^q \rightarrow \infty$, $h_2 \rightarrow 0$, and $n h_2^p \rightarrow \infty$ as $n \rightarrow \infty$, then,*

$$\sqrt{n h_2^p} [\tilde{g}_1(\mathbf{x}) - g_1(\mathbf{x}) - \text{bias}(\mathbf{x}) + o_p(h_{12}^2 + h_2^2)] \xrightarrow{\mathcal{D}} N\{0, v_0(\mathbf{x})\},$$

where the asymptotic bias is

$$\text{bias}(\mathbf{x}) = \frac{h_2^2}{2} \text{tr}\{\mu_2(J) g_1''(\mathbf{x})\} - \frac{h_{12}^2}{2} \text{tr}\{\mu_2(L) E(g_2''(\mathbf{Y}_t) | \mathbf{X}_t = \mathbf{x})\}$$

and the asymptotic variance is $v_0(\mathbf{x}) = \nu_0(J) B(\mathbf{x})$.

We remark that by Theorem 4.3, the asymptotic variance of the two-stage estimator is independent of the initial bandwidths. Thus, the initial bandwidths should be chosen as small as possible. This is another benefit of using the two-stage procedure: the bandwidth selection problem becomes relatively easy. In particular, when $h_{12} = o(h_2)$, the bias from the initial estimation can be asymptotically negligible. For the ideal situation that $g_2(\cdot)$ is known, Masry and Fan (1997) show that under some regularity conditions, the optimal estimate of $g_1(\mathbf{x})$, denoted by $\hat{g}_1^*(\mathbf{x})$, by using (4.64) in which the partial residual Z_t^* is replaced by the partial error $\tilde{Z}_t = \mathbf{Y}_t - \mu - g_2(\mathbf{Y}_t)$, is asymptotically normally distributed,

$$\sqrt{n h_2^p} \left[\hat{g}_1^*(\mathbf{x}) - g_1(\mathbf{x}) - \frac{h_2^2}{2} \text{tr}\{\mu_2(J) g_1''(\mathbf{x})\} + o_p(h_2^2) \right] \xrightarrow{\mathcal{D}} N\{0, v_0(\mathbf{x})\}.$$

This, in conjunction with Theorem 4.3, shows that the two-stage estimator and the ideal estimator share the same asymptotic bias and variance if $h_{12} = o(h_2)$.

4.6.5 Monte Carlo Simulations and Applications

See the paper by Cai (2002) for the detailed Monte Carlo simulation results and applications.

4.6.6 New Developments

See the paper by Mammen, Linton and Nielsen (1999).

4.6.7 Additive Model to to Boston House Price Data

There have been several papers devoted to the analysis of this dataset using some nonparametric methods. For example, Breiman and Friedman (1985), Pace (1993), Chaudhuri, Doksum and Samarov (1997), and Opsomer and Ruppert (1998) used four covariates: X_6 , X_{10} , X_{11} and X_{13} or their transformations (including the transformation on Y) to fit the data through a mean **additive regression model** such as

$$\log(Y) = \mu + g_1(X_6) + g_2(X_{10}) + g_3(X_{11}) + g_4(X_{13}) + \varepsilon, \quad (4.66)$$

where the additive components $\{g_j(\cdot)\}$ are unspecified smooth functions. Pace (1993) and Chaudhuri, Doksum and Samarov (1997) also considered the nonparametric estimation of the first derivative of each additive component which measures how much the response changes as one covariate is perturbed while the other covariates are held fixed; see Chaudhuri, Doksum and Samarov (1997). Let us use model (4.66) to fit the Boston house price data. The results are summarized in Figure 4.3 (the **R** code can be found in Section 4.7.2). Also, we fit a **semi-parametric additive model** as

$$\log(Y) = \mu + g_1(X_6) + \beta_2 X_{10} + \beta_3 X_{11} + \beta_4 X_{13} + \varepsilon. \quad (4.67)$$

The results are summarized in Figure 4.4 (the **R** code can be found in Section 4.7.2).

4.7 Computer Code

4.7.1 Example 4.1

```
# 04-28-2007
graphics.off() # clean the previous graphs on the screen

#####
# Example 4.1
#####
```

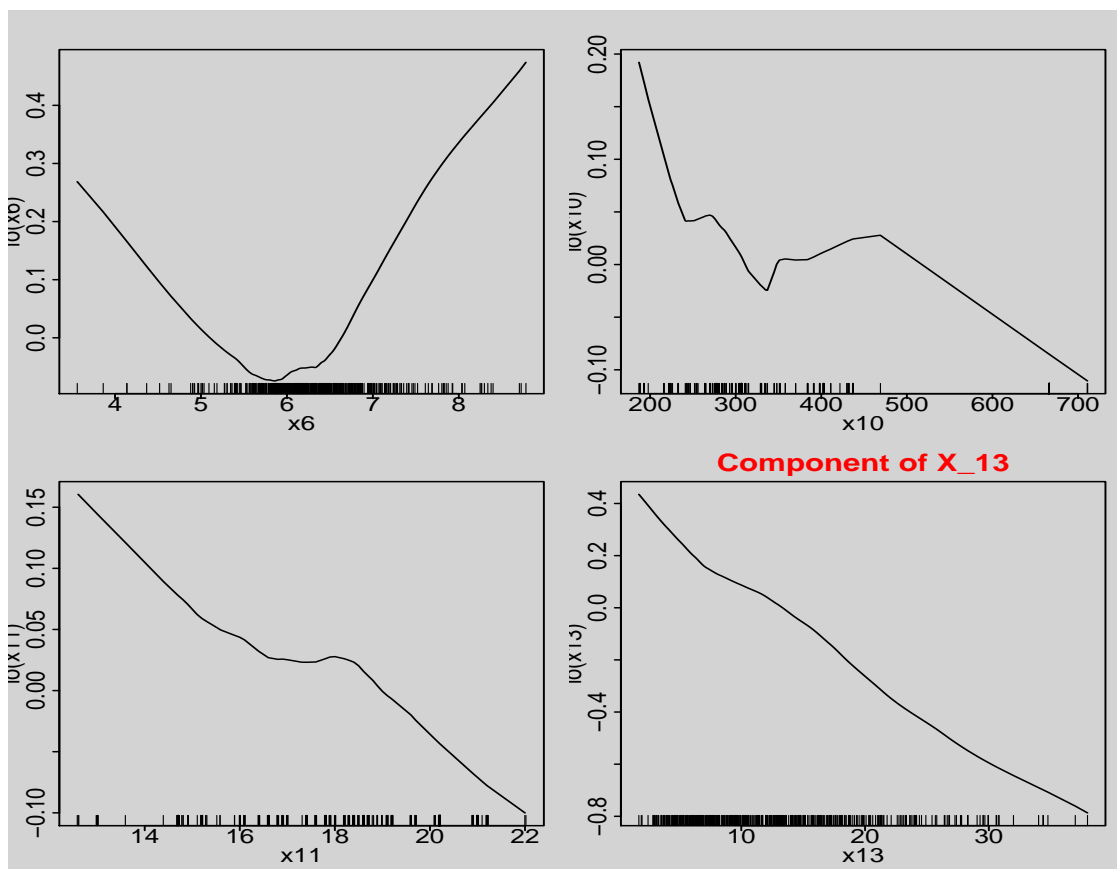



Figure 4.3: The results from model (4.66).

```
#####
z1=read.table(file="c:/res-teach/xiada/teaching05-07/data/ex4-1.txt")
# data: weekly 3-month Treasury bill from 1970 to 1997
x=z1[,4]/100
n=length(x)
y=diff(x)                # Delta x_t=x_t-x_{t-1}
x=x[1:(n-1)]
n=n-1
x_star=(x-mean(x))/sqrt(var(x))
z=seq(min(x),max(x),length=50)

win.graph()
#postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-4.1.eps",
```

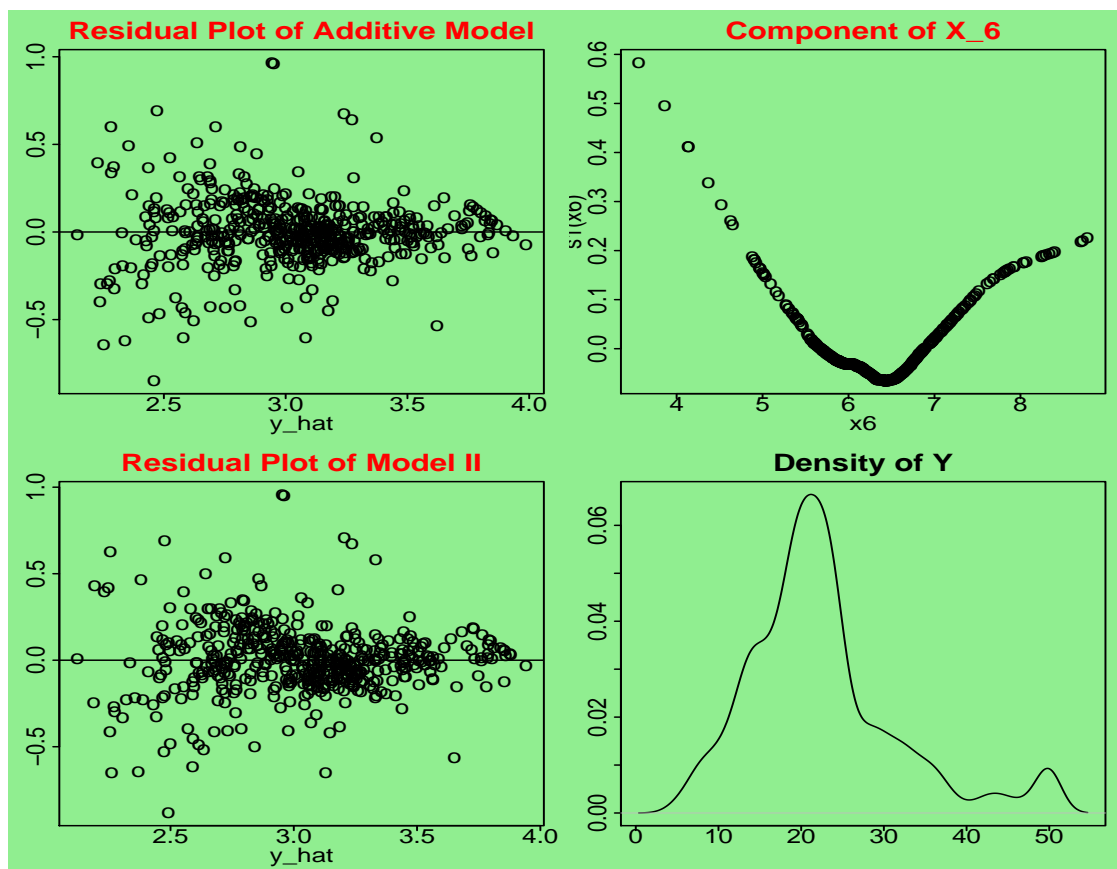


Figure 4.4: (a) Residual plot for model (4.66). (b) Plot of $g_1(x_6)$ versus x_6 . (c) Residual plot for model (4.67). (d) Density estimate of Y .

```
# horizontal=F,width=6,height=6)
par(mfrow=c(2,2),mex=0.4,bg="light blue")
scatter.smooth(x,y,span=1/10,ylab="",xlab="x(t-1)",evaluation=60)
  title(main="(a) y(t) vs x(t)",col.main="red")
scatter.smooth(x,abs(y),span=1/10,ylab="",xlab="x(t-1)",evaluation=60)
  title(main="(b) |y(t)| vs x(t)",col.main="red")
scatter.smooth(x,y^2,span=1/10,ylab="",xlab="x(t-1)",evaluation=60)
  title(main="(c) y(t)^2 vs x(t)",col.main="red")
#dev.off()

#####

#####

# Nonparametric Fitting #
```

```
#####

#####

# Define the Epanechnikov kernel function
kernel<-function(x){0.75*(1-x^2)*(abs(x)<=1)}

#####

# Define the kernel density estimator
kernden=function(x,z,h,ker){
  # parameters: x=variable; h=bandwidth; z=grid point; ker=kernel
  nz<-length(z)
  nx<-length(x)
  x0=rep(1,nx*nz)
  dim(x0)=c(nx,nz)
  x1=t(x0)
  x0=x*x0
  x1=z*x1
  x0=x0-t(x1)
  if(ker==1){x1=kernel(x0/h)}      # Epanechnikov kernel
  if(ker==0){x1=dnorm(x0/h)}      # normal kernel
  f1=apply(x1,2,mean)/h
  return(f1)
}

#####

# Define the local constant estimator
local.constant=function(y,x,z,h,ker){
  # parameters: x=variable; h=bandwidth; z=grid point; ker=kernel
  nz<-length(z)
  nx<-length(x)
  x0=rep(1,nx*nz)
  dim(x0)=c(nx,nz)
  x1=t(x0)
```

```

x0=x*x0
x1=z*x1
x0=x0-t(x1)
if(ker==1){x1=kernel(x0/h)}          # Epanechnikov kernel
if(ker==0){x1=dnorm(x0/h)}           # normal kernel
x2=y*x1
f1=apply(x1,2,mean)
f2=apply(x2,2,mean)
f3=f2/f1
return(f3)
}

#####

# Define the local linear estimator
local.linear<-function(y,x,z,h){
# parameters: y=response, x=design matrix; h=bandwidth; z=grid point
nz<-length(z)
ny<-length(y)
beta<-rep(0,nz*2)
dim(beta)<-c(nz,2)
for(k in 1:nz){
x0=x-z[k]
w0<-kernel(x0/h)
beta[k,]<-glm(y~x0,weight=w0)$coeff
}
return(beta)
}

#####

h=0.02

# Local constant estimate

```

```

mu_hat=local.constant(y,x,z,h,1)
sigma_hat=local.constant(abs(y),x,z,h,1)
sigma2_hat=local.constant(y^2,x,z,h,1)

#win.graph()
postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-4.1.eps",
  horizontal=F,width=6,height=6)
par(mfrow=c(2,2),mex=0.4,bg="light yellow")
scatter.smooth(x,y,span=1/10,ylab="",xlab="x(t-1)")
points(z,mu_hat,type="l",lty=1,lwd=3,col=2)
  title(main="(a) y(t) vs x(t)",col.main="red")
legend(0.04,0.0175,"Local Constant Estimate")
scatter.smooth(x,abs(y),span=1/10,ylab="",xlab="x(t-1)")
points(z,sigma_hat,type="l",lty=1,lwd=3,col=2)
  title(main="(b) |y(t)| vs x(t)",col.main="red")
scatter.smooth(x,y^2,span=1/10,ylab="",xlab="x(t-1)")
  title(main="(c) y(t)^2 vs x(t)",col.main="red")
points(z,sigma2_hat,type="l",lty=1,lwd=3,col=2)
dev.off()

# Local Linear Estimate

fit2=local.linear(y,x,z,h)
mu_hat=fit2[,1]
fit2=local.linear(abs(y),x,z,h)
sigma_hat=fit2[,1]
fit2=local.linear(y^2,x,z,h)
sigma2_hat=fit2[,1]

#win.graph()
postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-4.2.eps",

```

```

horizontal=F,width=6,height=6)
par(mfrow=c(2,2),mex=0.4,bg="light green")
scatter.smooth(x,y,span=1/10,ylab="",xlab="x(t-1)")
points(z,mu_hat,type="l",lty=1,lwd=3,col=2)
title(main="(a) y(t) vs x(t)",col.main="red")
legend(0.04,0.0175,"Local Linear Estimate")
scatter.smooth(x,abs(y),span=1/10,ylab="",xlab="x(t-1)")
points(z,sigma_hat,type="l",lty=1,lwd=3,col=2)
title(main="(b) |y(t)| vs x(t)",col.main="red")
scatter.smooth(x,y^2,span=1/10,ylab="",xlab="x(t-1)")
title(main="(c) y(t)^2 vs x(t)",col.main="red")
points(z,sigma2_hat,type="l",lty=1,lwd=3,col=2)
dev.off()
#####

```

4.7.2 Codes for Additive Modeling Analysis of Boston Data

The following is the R code for making Figures 4.3 and 4.4.

```

data=read.table("c:/res-teach/xiada/teaching05-07/data/ex4-2.txt")
y=data[,14]
x1=data[,1]
x6=data[,6]
x10=data[,10]
x11=data[,11]
x13=data[,13]
y_log=log(y)
library(gam)
fit_gam=gam(y_log~lo(x6)+lo(x10)+lo(x11)+lo(x13))
resid=fit_gam$residuals
y_hat=fit_gam$fitted
postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-boston1.eps",
horizontal=F,width=6,height=6,bg="light grey")

```

```

par(mfrow=c(2,2),mex=0.4)
plot(fit_gam)
title(main="Component of X_13",col.main="red",cex=0.6)
dev.off()
fit_gam1=gam(y_log~lo(x6)+x10+x11+x13)
s1=fit_gam1$smooth[,1]          # obtain the smoothed component
resid1=fit_gam1$residuals
y_hat1=fit_gam1$fitted

print(summary(fit_gam1))

postscript(file="c:/res-teach/xiada/teaching05-07/figs/fig-boston2.eps",
  horizontal=F,width=6,height=6,bg="light green")
par(mfrow=c(2,2),mex=0.4)
plot(y_hat,resid,type="p",pch="o",ylab="",xlab="y_hat")
title(main="Residual Plot of Additive Model",col.main="red",cex=0.6)
abline(0,0)
plot(x6,s1,type="p",pch="o",ylab="s1(x6)",xlab="x6")
title(main="Component of X_6",col.main="red",cex=0.6)
plot(y_hat1,resid1,type="p",pch="o",ylab="",xlab="y_hat")
title(main="Residual Plot of Model II",col.main="red",cex=0.5)
abline(0,0)
plot(density(y),ylab="",xlab="",main="Density of Y")
dev.off()

```

4.8 References

- Aït-Sahalia, Y. (1996). Nonparametric pricing of interest rate derivative securities. *Econometrica*, **64**, 527-560.
- Belsley, D.A., E. Kuh and R.E. Welsch (1980). *Regression Diagnostic: Identifying Influential Data and Sources of Collinearity*. New York: Wiley.
- Breiman, L. and J.H. Friedman (1985). Estimating optimal transformation for multiple

- regression and correlation. *Journal of the American Statistical Association*, **80**, 580-619.
- Cai, Z. (2002). A two-stage approach to additive time series models. *Statistica Neerlandica*, **56**, 415-433.
- Cai, Z., M. Das, H. Xiong and X. Wu (2006). Functional-Coefficient Instrumental Variables Models. *Journal of Econometrics*, **133**, 207-241.
- Cai, Z. and J. Fan (2000). Average regression surface for dependent data. *Journal of Multivariate Analysis*, **75**, 112-142.
- Cai, Z., J. Fan and Q. Yao (2000). Functional-coefficient regression models for nonlinear time series. *Journal of American Statistical Association*, **95**, 941-956.
- Cai, Z. and E. Masry (2000). Nonparametric estimation of additive nonlinear ARX time series: Local linear fitting and projection. *Econometric Theory*, **16**, 465-501.
- Cai, Z. and R.C. Tiwari (2000). Application of a local linear autoregressive model to BOD time series. *Environmetrics*, **11**, 341-350.
- Chaudhuri, P., K. Doksum and A. Samarov (1997). On average derivative quantile regression. *The Annals of Statistics*, **25**, 715-744.
- Chen, R. and R. Tsay (1993). Nonlinear additive ARX models. *Journal of the American Statistical Association*, **88**, 310-320.
- Engle, R.F., C.W.J. Grabger, J. Rice, and A. Weiss (1986). Semiparametric estimates of the relation between weather and electricity sales. *Journal of The American Statistical Association*, **81**, 310-320.
- Fan, J. (1993). Local linear regression smoothers and their minimax efficiency. *The Annals of Statistics*, **21**, 196-216.
- Fan, J., T. Gasser, I. Gijbels, M. Brockmann and J. Engel (1996). Local polynomial fitting: optimal kernel and asymptotic minimax efficiency. *Annals of the Institute of Statistical Mathematics*, **49**, 79-99.
- Fan, J. and I. Gijbels (1996). *Local Polynomial Modeling and Its Applications*. London: Chapman and Hall.
- Fan, J., N.E. Heckman, and M.P. Wand (1995). Local polynomial kernel regression for generalized linear models and quasi-likelihood functions. *Journal of the American Statistical Association*, **90**, 141-150.
- Fan, J. and T. Huang (2005). Profile likelihood inferences on semiparametric varying-coefficient partially linear models. *Bernoulli*, **11**, 1031-1057.
- Fan, J. and Q. Yao (2003). *Nonlinear Time Series: Nonparametric and Parametric Methods*. New York: Springer-Verlag.

- Fan, J., Q. Yao and Z. Cai (2003). Adaptive varying-coefficient linear models. *Journal of the Royal Statistical Society, Series B*, **65**, 57-80.
- Fan, J. and C. Zhang (2003). A re-examination of diffusion estimators with applications to financial model validation. *Journal of the American Statistical Association*, **98**, 118-134.
- Fan, J., C. Zhang and J. Zhang (2001). Generalized likelihood test statistic and Wilks phenomenon. *The Annals of Statistics*, **29**, 153-193.
- Gasser, T. and H.-G. Müller (1979). Kernel estimation of regression functions. In *Smoothing Techniques for Curve Estimation*, Lecture Notes in Mathematics, **757**, 23-28. Springer-Verlag, New York.
- Gilley, O.W. and R.K. Pace (1996). On the Harrison and Rubinfeld Data. *Journal of Environmental Economics and Management*, **31**, 403-405.
- Granger, C.W.J., and T. Teräsvirta (1993). *Modeling Nonlinear Economic Relationships*. Oxford University Press, Oxford, U.K..
- Hall, P., and C.C. Heyde (1980). *Martingale Limit Theory and Its Applications*. New York: Academic Press.
- Hall, P., and I. Johnstone (1992). Empirical functional and efficient smoothing parameter selection (with discussion). *Journal of the Royal Statistical Society, Series B*, **54**, 475-530.
- Harrison, D. and D.L. Rubinfeld (1978). Hedonic housing prices and demand for clean air. *Journal of Environmental Economics and Management*, **5**, 81-102.
- Hastie, T.J. and R.J. Tibshirani (1990). *Generalized Additive Models*. Chapman and Hall, London.
- Hong, Y. and Lee, T.-H. (2003). Inference on via generalized spectrum and nonlinear time series models. *The Review of Economics and Statistics*, **85**, 1048-1062.
- Hurvich, C.M., J.S. Simonoff and C.-L. Tsai (1998). Smoothing parameter selection in nonparametric regression using an improved Akaike information criterion. *Journal of the Royal Statistical, Society B*, **60**, 271-293.
- Juhl, T. (2005). Functional coefficient models under unit root behavior. *Econometrics Journal*, **8**, 197-213.
- Koenker, R. (2005). *Quantile Regression*. Cambridge University Press, New York.
- Koenker, R. and G.W. Bassett (1978). Regression quantiles. *Econometrica*, **46**, 33-50.
- Koenker, R. and G.W. Bassett (1982). Robust tests for heteroscedasticity based on regression quantiles. *Econometrica*, **50**, 43-61.

- Kreiss, J.P., M. Neumann and Q. Yao (1998). Bootstrap tests for simple structures in nonparametric time series regression. **Unpublished manuscript**.
- Li, Q., C. Huang, D. Li and T. Fu (2002). Semiparametric smooth coefficient models. *Journal of Business and Economic Statistics*, **20**, 412-422.
- Linton, O.B. (1997). Efficient estimation of additive nonparametric regression models. *Biometrika*, **84**, 469-473.
- Linton, O.B. (2000). Efficient estimation of generalized additive nonparametric regression models. *Econometric Theory*, **16**, 502-523.
- Linton, O.B. and J.P. Nielsen (1995). A kernel method of estimating structured nonparametric regression based on marginal integration. *Biometrika*, **82**, 93-100.
- Mammen, E., O.B. Linton, and J.P. Nielsen (1999). The existence and asymptotic properties of a backfitting projection algorithm under weak conditions. *The Annals of Statistics*, **27**, 1443-1490.
- Masry, E. and J. Fan (1997). Local polynomial estimation of regression functions for mixing processes. *Scandinavian Journal of Statistics*, **24**, 165-179.
- Masry, E. and D. Tjøstheim (1997). Additive nonlinear ARX time series and projection estimates. *Econometric Theory*, **13**, 214-252.
- Şentürk, D. and H.-G. Müller (2003). Inference for covariate adjusted regression via varying coefficient models. Forthcoming in *Scandinavian Journal of Statistics*.
- Nadaraya, E.A. (1964). On estimating regression. *Theory of Probability and Its Applications*, **9**, 141-142.
- Øksendal, B. (1985). *Stochastic Differential Equations: An Introduction with Applications*, 3th edition. New York: Springer-Verlag.
- Opsomer, J.D. and D. Ruppert (1998). A fully automated bandwidth selection for additive regression model. *Journal of The American Statistical Association*, **93**, 605-618.
- Pace, R.K. (1993). Nonparametric methods with applications hedonic models. *Journal of Real Estate Finance and Economics*, **7**, 185-204.
- Pace, R.K. and O.W. Gilley (1997). Using the spatial configuration of the data to improve estimation. *Journal of the Real Estate Finance and Economics*, **14**, 333-340.
- Priestley, M.B. and M.T. Chao (1972). Nonparametric function fitting. *Journal of the Royal Statistical Society, Series B*, **34**, 384-392.
- Rice, J. (1984). Bandwidth selection for nonparametric regression. *The Annals of Statistics*, **12**, 1215-1230.

- Ruppert, D., S.J. Sheather and M.P. Wand (1995). An effective bandwidth selector for local least squares regression. *Journal of American Statistical Association*, **90**, 1257-1270.
- Ruppert, D. and M.P. Wand (1994). Multivariate weighted least squares regression. *The Annals of Statistics*, **22**, 1346-1370.
- Rousseeuw, R.J. and A.M. Leroy (1987). *Robust Regression and Outlier Detection*. New York: Wiley.
- Shao, Q. and H. Yu (1996). Weak convergence for weighted empirical processes of dependent sequences. *The Annals of Probability*, **24**, 2098-2127.
- Sperlich, S., D. Tjøstheim, and L. Yang (2002). Nonparametric estimation and testing of interaction in additive models. *Econometric Theory*, **18**, 197-251.
- Stanton, R. (1997). A nonparametric model of term structure dynamics and the market price of interest rate risk. *Journal of Finance*, **52**, 1973-2002.
- Sun, Z. (1984). Asymptotic unbiased and strong consistency for density function estimator. *Acta Mathematica Sinica*, **27**, 769-782.
- Tjøstheim, D. and B. Auestad (1994a). Nonparametric identification of nonlinear time series: Projections. *Journal of the American Statistical Association*, **89**, 1398-1409.
- Tjøstheim, D. and B. Auestad (1994b). Nonparametric identification of nonlinear time series: Selecting significant lags. *Journal of the American Statistical Association*, **89**, 1410-1419.
- van Dijk, D., T. Teräsvirta, and P.H. Franses (2002). Smooth transition autoregressive models - a survey of recent developments. *Econometric Reviews*, **21**, 1-47.
- Watson, G.S. (1964). Smooth regression analysis. *Sankhyā, Series A*, **26**, 359-372.

Chapter 5

Nonparametric Quantile Models

For details, see the paper by Cai and Xu (2005). If you like to read the whole paper, you can download it from the web site at <http://www.wise.xmu.edu.cn/> at **WORKING PAPER** column. Next we present only a part of the whole paper of Cai and Xu (2005).

5.1 Introduction

Over the last three decades, quantile regression, also called conditional quantile or regression quantile, introduced by Koenker and Bassett (1978), has been used widely in various disciplines, such as finance, economics, medicine, and biology. It is well-known that when the distribution of data is typically skewed or data contains some outliers, the median regression, a special case of quantile regression, is more explicable and robust than the mean regression. Also, regression quantiles can be used to test heteroscedasticity formally or graphically (Koenker and Bassett, 1982; Efron, 1991; Koenker and Zhao, 1996; Koenker and Xiao, 2002). Although some individual quantiles, such as the conditional median, are sometimes of interest in practice, more often one wishes to obtain a collection of conditional quantiles which can characterize the entire conditional distribution. More importantly, another application of conditional quantiles is the construction of prediction intervals for the next value given a small section of the recent past values in a stationary time series (Granger, White, and Kamstra, 1989; Koenker, 1994; Zhou and Portnoy, 1996; Koenker and Zhao, 1996; Taylor and Bunn, 1999). Also, Granger, White, and Kamstra (1989), Koenker and Zhao (1996), and Taylor and Bunn (1999) considered an interval forecasting for parametric autoregressive conditional heteroscedastic (ARCH) type models. For more details about the historical and recent developments of quantile regression with applications for time series data, particularly

in finance, see, for example, the papers and books by J.P. Morgan (1995), Duffie and Pan (1997), Jorin (2000), Koenker (2000), Koenker and Hallock (2001), Tsay (2000, 2002), Khindanova and Rachev (2000), and Bao, Lee and Saltoğlu (2001), and the references therein.

Recently, the quantile regression technique has been successfully applied to politics. For example, in the 1992 presidential selection, the Democrats used the yearly Current Population Survey data to show that between 1980 and 1992 there was an increase in the number of people in the high-salary category as well as an increase in the number of people in the low-salary category. This phenomena could be illustrated by using the quantile regression method as follows: computing 90% and 10% quantile regression functions of salary as a function of time. An increasing 90% quantile regression function and a decreasing 10% quantile regression function corresponded to the Democrats' claim that "the rich got richer and the poor got poorer" during the Republican administrations; see Figure 6.4 in Fan and Gijbels (1996, p. 229).

More importantly, by following the regulations of the Bank for International Settlements, many of financial institutions have begun to use a uniform measure of risk to measure the market risks called Value-at-Risk (VaR), which can be defined as the maximum potential loss of a specific portfolio for a given horizon in finance. In essence, the interest is to compute an estimate of the lower tail quantile (with a small probability) of future portfolio returns, conditional on current information. Therefore, the VaR can be regarded as a special application of the quantile regression. There is a vast amount of literature in this area; see, to name just a few, J.P. Morgan (1995), Duffie and Pan (1997), Engle and Manganelli (2004), Jorin (2000), Tsay (2000, 2002), Khindanova and Rachev (2000), and Bao, Lee and Saltoğlu (2001), and references therein.

In this chapter, we assume that $\{\mathbf{X}_t, Y_t\}_{t=-\infty}^{\infty}$ is a stationary sequence. Denote $F(y|\mathbf{x})$ the conditional distribution of Y given $\mathbf{X} = \mathbf{x}$, where $\mathbf{X}_t = (X_{t1}, \dots, X_{td})'$ with $'$ denoting the transpose of a matrix or vector, is the associated covariate vector in \mathfrak{R}^d with $d \geq 1$, which might be a function of exogenous (covariate) variables or some lagged (endogenous) variables or time t . The regression (conditional) quantile function $q_\tau(\mathbf{x})$ is defined as, for

any $0 < \tau < 1$,

$$q_\tau(\mathbf{x}) = \inf\{y \in \mathbb{R}^1 : F(y|\mathbf{x}) \geq \tau\}, \quad \text{or} \quad q_\tau(\mathbf{x}) = \operatorname{argmin}_{a \in \mathbb{R}^1} E\{\rho_\tau(Y_t - a) | \mathbf{X}_t = \mathbf{x}\}, \quad (5.1)$$

where $\rho_\tau(y) = y(\tau - I_{\{y < 0\}})$ with $y \in \mathbb{R}^1$ is called the loss (“check”) function, and I_A is the indicator function of any set A . There are several advantages of using a quantile regression:

- **A quantile regression does not require knowing the distribution of the dependent variable.**
- **It does not require the symmetry of the measurement error.**
- **It can characterize the heterogeneity.**
- **It can estimate the mean and variance simultaneously.**
- **It is a robust procedure.**
- **There are a lot more.**

Having conditioned on the observed characteristics $\mathbf{X}_t = \mathbf{x}$, based on the Skorohod representation, Y_t and the quantile function $q_\tau(\mathbf{x})$ have a following relationship as

$$Y_t = q(\mathbf{X}_t, U_t), \quad (5.2)$$

where $U_t | \mathbf{X}_t \sim U(0, 1)$. We will refer to U_t as the rank variable, and note that representation (5.2) is essential to what follows. The rank variable U_t is responsible for heterogeneity of outcomes among individuals with the same observed characteristics \mathbf{X}_t . It also determines their relative ranking in terms of potential outcomes; hence one may think of rank U_t as representing some unobserved characteristic. This interpretation makes quantile analysis an interesting tool for describing and learning the structure of heterogeneous effects and controlling for unobserved heterogeneity.

Clearly, the simplest form of model (5.1) is $q_\tau(\mathbf{x}) = \beta'_\tau \mathbf{x}$, which is called the linear quantile regression model well studied by many authors. For details, see the papers by Duffie and Pan (1997), Koenker (2000), Tsay (2000), Koenker and Hallock (2001), Khindanova and Rachev (2000), and Bao, Lee and Saltoğlu (2001), Engle and Manganelli (2004), and references therein.

In many practical applications, however, the linear quantile regression model might not be “rich” enough to capture the underlying relationship between the quantile of response variable and its covariates. Indeed, some components may be highly nonlinear or some covariates may be interactive. To make the quantile regression model more flexible, there is a swiftly growing literature on nonparametric quantile regression. Various smoothing techniques, such as kernel methods, splines, and their variants, have been used to estimate the nonparametric quantile regression for both the independent and time series data. For the recent developments and the detailed discussions on theory, methodologies, and applications, see, for example, the papers by He, Ng, and Portony (1998), Yu and Jones (1998), He and Ng (1999), He and Portony (2000), Honda (2000, 2004), Tsay (2000, 2002), Lu, Hui and Zhao (2000), Khindanova and Rachev (2000), Bao, Lee and Saltoğlu (2001), Cai (2002a), De Gooijer, and Gannoun (2003), Horowitz and Lee (2005), Yu and Lu (2004), and Li and Racine (2004), and references therein. In particular, for the univariate case, recently, Honda (2000) and Lu, Hui and Zhao (2000) derived the asymptotic properties of the local linear estimator of the quantile regression function under α -mixing condition. For the high dimensional case, however, the aforementioned methods encounter some difficulties such as the so-called “curse of dimensionality” and their implementation in practice is not easy as well as the visual display is not so useful for the exploratory purposes.

To attenuate the above problems, De Gooijer and Zerom (2003), Horowitz and Lee (2005), and Yu and Lu (2004) considered an additive quantile regression model $q_\tau(\mathbf{X}_t) = \sum_{k=1}^d g_k(X_{tk})$. To estimate each component, for the time series case, De Gooijer and Zerom (2003) first estimated a high dimensional quantile function by inverting the conditional distribution function estimated by using a weighted Nadaraya-Watson approach, proposed by Cai (2002a), and then used a projection method to estimate each component, as discussed in Cai and Masry (2000), while Yu and Lu (2004) focused on the independent data and used a back-fitting algorithm method to estimate each component. On the other hand, to estimate each additive component for the independent data, Horowitz and Lee (2005) used a two-stage approach consisting of the series estimation at the first step and a local polynomial fitting at the second step. For the independent data, the above model was extended by He, Ng and Portony (1998), He and Ng (1999), and He and Portony (2000) to include interaction terms by using spline methods.

In this chapter, we adapt another dimension reduction modelling method to analyze

dynamic time series data, termed as the smooth (functional or varying) coefficient modelling approach. This approach allows appreciable flexibility on the structure of fitted models. It allows for linearity in some continuous or discrete variables which can be exogenous or lagged and nonlinear in other variables in the coefficients. In such a way, the model has the ability of capturing the individual variations. More importantly, it can ease the so-called “curse of dimensionality” and combines both additivity and interactivity. A smooth coefficient quantile regression model for time series data takes the following form

$$q_\tau(\mathbf{U}_t, \mathbf{X}_t) = \sum_{k=0}^d a_k(\mathbf{U}_t) X_{tk} = \mathbf{X}_t' \mathbf{a}_\tau(\mathbf{U}_t), \quad (5.3)$$

where \mathbf{U}_t is called the smoothing variable, which might be one part of X_{t1}, \dots, X_{td} or just time or other exogenous variables or the lagged variables, $\mathbf{X}_t = (X_{t0}, X_{t1}, \dots, X_{td})'$ with $X_{t0} \equiv 1$, $\{a_k(\cdot)\}$ are smooth coefficient functions, and $\mathbf{a}_\tau(\cdot) = (a_{0,\tau}(\cdot), \dots, a_{d,\tau}(\cdot))'$. Here, some of $\{a_{k,\tau}(\cdot)\}$ are allowed to depend on τ . For simplicity, we drop τ from $\{a_{k,\tau}(\cdot)\}$ in what follows. It is our interest here to estimate the coefficient functions $\mathbf{a}(\cdot)$ rather than the quantile regression surface $q_\tau(\cdot, \cdot)$ itself. Note that model (5.3) was studied by Honda (2004) for the independent sample, but our focus here is on the dynamic model for nonlinear time series, which is more appropriate for economic and financial applications.

The general setting in (5.3) covers many familiar quantile regression models, including the quantile autoregressive model (QAR) proposed by Koenker and Xiao (2004) who applied the QAR model for the unit root inference. In particular, it includes a specific class of ARCH models, such as heteroscedastic linear models considered by Koenker and Zhao (1996). Also, if there is no \mathbf{X}_t in the model ($d = 0$), $q_\tau(\mathbf{U}_t, \mathbf{X}_t)$ becomes $q_\tau(\mathbf{U}_t)$ so that model (5.3) reduces to the ordinary nonparametric quantile regression model which has been studied extensively. For the recent developments, refer to the papers by He, Ng and Portony (1998), Yu and Jones (1998), He and Ng (1999), He and Portony (2000), Honda (2000), Lu, Hui and Zhao (2000), Cai (2002a), De Gooijer and Zerom (2003), Horowitz and Lee (2005), Yu and Lu (2004), and Li and Racine (2004). If \mathbf{U}_t is just time, then the model is called the time-varying coefficient quantile regression model, which is potentially useful to see whether the quantile regression changes over time and in a case with a practical interest is, for example, the aforementioned illustrative example for the 1992 presidential election and the analysis of the reference growth data by Cole (1994), Wei, Pere, Koenker and He (2006), and Wei and He (2006), and the references therein. However, if \mathbf{U}_t is time, the observed time series might

not be stationary. Therefore, the treatment for non-stationary case would require a different approach so that it is beyond the scope of this chapter and deserves a further investigation. For more applications, see the work in Xu (2005). Finally, note that the smooth coefficient mean regression model is one of the most popular nonlinear time series models in mean regression and has various applications. For more discussions, refer to the papers by Chen and Tsay (1993), Cai, Fan, and Yao (2000), Cai and Tiwari (2000), Cai (2007), Hong and Lee (2003), and Wang (2003), and the book by Tsay (2002), and references therein.

The motivation of this study comes from an analysis of the well known Boston housing price data, consisting of several variables collected on each of 506 different houses from a variety of locations. The interest is to identify the factors affecting the house price in Boston area. As argued by Şentürk and Müller (2005), the correlation between the house price and the crime rate can be adjusted by the confounding variable which is the proportion of population of lower educational status through a varying coefficient model and the expected effect of increasing crime rate on declining house prices seems to be only observed for lower educational status neighborhoods in Boston. The interesting features of this dataset are that the response variable is the median price of a home in a given area and the distributions of the price and the major covariate (the confounding variable) are left skewed. Therefore, quantile methods are suitable for the analysis of this dataset. Therefore, such a problem can be tackled by using model (5.3). In another example, one is interested in exploring the possible nonlinearity feature, heteroscedasticity, and predictability of the exchange rates such as the Japanese Yen per US dollar. The detailed analysis of these data sets is reported in Section 3.

5.2 Modeling Procedures

5.2.1 Local Linear Quantile Estimate

Now, we apply the local polynomial method to the smooth coefficient quantile regression model as follows. For the sake of brevity, we only consider the case where \mathbf{U}_t in (5.3) is one-dimensional, denoted by U_t in what follows. Extension to multivariate \mathbf{U}_t involves fundamentally no new ideas although the theory and procedure continue to hold. Note that the models with high dimension might not be practically useful due to the curse of dimensionality. A local polynomial fitting has several nice properties such as high statistical

efficiency in an asymptotic minimax sense, design-adaptation, and automatic edge correction (see, e.g., Fan and Gijbels, 1996).

We estimate the functions $\{a_k(\cdot)\}$ using the local polynomial regression method from observations $\{(U_t, \mathbf{X}_t, Y_t)\}_{t=1}^n$. We assume throughout the chapter that the coefficient functions $\mathbf{a}(\cdot)$ have the $(q+1)$ th derivative, so that for any given grid point u_0 , $a_k(\cdot)$ can be approximated by a polynomial function in a neighborhood of the given grid point u_0 as $\mathbf{a}(U_t) \approx \mathbf{a}(u_0) + \mathbf{a}'(u_0)(U_t - u_0) + \cdots + \mathbf{a}^{(q)}(u_0)(U_t - u_0)^q/q!$ and

$$q_\tau(U_t, \mathbf{X}_t) \approx \sum_{j=0}^q \mathbf{X}_t' \boldsymbol{\beta}_j (U_t - u_0)^j,$$

where $\boldsymbol{\beta}_j = \mathbf{a}^{(j)}(u_0)/j!$. Then, the locally weighted loss function is

$$\sum_{t=1}^n \rho_\tau \left(Y_t - \sum_{j=0}^q \mathbf{X}_t' \boldsymbol{\beta}_j (U_t - u_0)^j \right) K_h(U_t - u_0), \quad (5.4)$$

where $K(\cdot)$ is a kernel function, $K_h(x) = K(x/h)/h$, and $h = h_n$ is a sequence of positive numbers tending to zero, which controls the amount of smoothing used in estimation. Solving the minimization problem in (5.4) gives $\hat{\mathbf{a}}(u_0) = \hat{\boldsymbol{\beta}}_0$, the local polynomial estimate of $\mathbf{a}(u_0)$, and $\hat{\mathbf{a}}^{(j)}(u_0) = j! \hat{\boldsymbol{\beta}}_j$ ($j \geq 1$), the local polynomial estimate of the j th derivative $\mathbf{a}^{(j)}(u_0)$ of $\mathbf{a}(u_0)$. By moving u_0 along with the real line, one obtains the estimate for the entire curve. For various practical applications, Fan and Gijbels (1996) recommended using the local linear fit ($q = 1$). Therefore, for the expositional purpose, in what follows, we only consider the case $q = 1$ (local linear fitting).

The programming involved in the local (polynomial) linear quantile estimation is relatively simple and can be modified with few efforts from the existing programs for a linear quantile model. For example, for each grid point u_0 , the local linear quantile estimation can be implemented in the **R** package **quantreg**, of Koenker (2004) by setting covariates as \mathbf{X}_t and $\mathbf{X}_t(U_t - u_0)$ and the weight as $K_h(U_t - u_0)$.

Although some modifications are needed, the method developed here for the local linear quantile estimation is applicable to a general local polynomial quantile estimation. In particular, we note that the local constant (Nadaraya-Watson type) quantile estimation of $\mathbf{a}(u_0)$, denoted by $\tilde{\mathbf{a}}(u_0)$, is $\tilde{\boldsymbol{\beta}}$ minimizing the following subjective function

$$\sum_{t=1}^n \rho_\tau(Y_t - \mathbf{X}_t' \boldsymbol{\beta}) K_h(U_t - u_0), \quad (5.5)$$

which is a special case of (5.4) with $q = 0$. We compare $\widehat{\mathbf{a}}(u_0)$ and $\widetilde{\mathbf{a}}(u_0)$ theoretically at the end of Section 2.2 and empirically in Section 3.1 and the comparison leads to suggest that one should use the local linear approach in practice.

5.2.2 Asymptotic Results

We first give some regularity conditions that are sufficient for the consistency and asymptotic normality of the proposed estimators, although they might not be the weakest possible. We introduce the following notations. Denote

$$\Omega(u_0) \equiv E[\mathbf{X}_t \mathbf{X}_t' | \mathbf{U}_t = u_0] \quad \text{and} \quad \Omega^*(u_0) \equiv E[\mathbf{X}_t \mathbf{X}_t' f_{y|u,x}(q_\tau(u_0, \mathbf{X}_t)) | \mathbf{U}_t = u_0],$$

where $f_{y|u,x}(y)$ is the conditional density of Y given U and \mathbf{X} . Let $f_u(u)$ present the marginal density of U .

Assumptions:

- (C1) $\mathbf{a}(u)$ is twice continuously differentiable in a neighborhood of u_0 for any u_0 .
- (C2) $f_u(u)$ is continuous and $f_u(u_0) > 0$.
- (C3) $f_{y|u,x}(y)$ is bounded and satisfies the Lipschitz condition.
- (C4) The kernel function $K(\cdot)$ is symmetric and has a compact support, say $[-1, 1]$.
- (C5) $\{(\mathbf{X}_t, Y_t, \mathbf{U}_t)\}$ is a strictly α -mixing stationary process with mixing coefficient $\alpha(t)$ satisfies $\sum_{t \geq 1}^\infty t^l \alpha^{(\delta-2)/\delta}(t) < \infty$ for some positive real number $\delta > 2$ and $l > (\delta - 2)/\delta$.
- (C6) $E\|\mathbf{X}_t\|^{2\delta^*} < \infty$ with $\delta^* > \delta$.
- (C7) $\Omega(u_0)$ is positive-definite and continuous in a neighborhood of u_0 .
- (C8) $\Omega^*(u_0)$ is continuous and positive-definite in a neighborhood of u_0 .
- (C9) The bandwidth h satisfies $h \rightarrow 0$ and $nh \rightarrow \infty$.
- (C10) $f(u, v | \mathbf{x}_0, \mathbf{x}_s; s) \leq M < \infty$ for $s \geq 1$, where $f(u, v | \mathbf{x}_0, \mathbf{x}_s; s)$ is the conditional density of (U_0, U_s) given $(\mathbf{X}_0 = \mathbf{x}_0, \mathbf{X}_s = \mathbf{x}_s)$.
- (C11) $n^{1/2-\delta/4} h^{\delta/\delta^*-1/2-\delta/4} = O(1)$.

Remark 1: (*Discussion of Conditions*) Assumptions (C1) - (C3) include some smoothness conditions on functionals involved. The requirement in (C4) that $K(\cdot)$ be compactly supported is imposed for the sake of brevity of proofs, and can be removed at the cost of lengthier arguments. In particular, the Gaussian kernel is allowed. The α -mixing is one of the weakest mixing conditions for weakly dependent stochastic processes. Stationary time series or Markov chains fulfilling certain (mild) conditions are α -mixing with exponentially decaying coefficients; see the discussions in Section 1 and Cai (2002a) for more examples. On the other hand, the assumption on the convergence rate of $\alpha(\cdot)$ in (C5) might not be the weakest possible and is imposed to simplify the proof. Further, (C10) is just a technical assumption, which is also imposed by Cai (2002a). (C6) - (C8) require some standard moments. Clearly, (C11) allows the choice of a wide range of smoothing parameter values and is slightly stronger than the usual condition of $nh \rightarrow \infty$. However, for the bandwidths of optimal size (i.e., $h = O(n^{-1/5})$), (C11) is automatically satisfied for $\delta \geq 3$ and it is still fulfilled for $2 < \delta < 3$ if δ^* satisfies $\delta < \delta^* \leq 1 + 1/(3 - \delta)$, so that we do not concern ourselves with such refinements. Indeed, this assumption is also imposed by Cai, Fan and Yao (2000) for the mean regression. Finally, if there is no \mathbf{X}_t in model (5.3), (C5) can be replaced by (C5)': $\alpha(t) = O(t^{-\delta})$ for some $\delta > 2$ and (C11) can be substituted by (C11)': $nh^{\delta/(\delta-2)} \rightarrow \infty$; see Cai (2002a) for details.

Remark 2: (*Identification*) It is clear from (5.3) that

$$\Omega(u_0) \mathbf{a}(u_0) = E[q_\tau(u_0, \mathbf{X}_t) \mathbf{X}_t | U_t = u_0].$$

Then, $\mathbf{a}(u_0)$ is identified (uniquely determined) if and only if $\Omega(u_0)$ is positive definite for any u_0 . Therefore, Assumption (C7) is the necessary and sufficient condition for the model identification.

To establish the asymptotic normality of the proposed estimator, similar to Chaudhuri (1991), we first derive the local Bahadur representation for the local linear estimator. To this end, our analysis follows the approach of Koenker and Zhao (1996), which can simplify the theoretical proofs. Define, $\mu_j = \int u^j K(u) du$ and $\nu_j = \int u^j K^2(u) du$. Also, set $\psi_\tau(x) = \tau - I_{\{x < 0\}}$, $U_{th} = (U_t - u_0)/h$, $\mathbf{X}_t^* = \begin{pmatrix} \mathbf{X}_t \\ U_{th} \mathbf{X}_t \end{pmatrix}$, $Y_t^* = Y_t - \mathbf{X}_t'[\mathbf{a}(u_0) + \mathbf{a}'(u_0)(U_t - u_0)]$, and $\boldsymbol{\theta} = \sqrt{nh} \mathbf{H} \begin{pmatrix} \beta_0 - \mathbf{a}(u_0) \\ \beta_1 - \mathbf{a}'(u_0) \end{pmatrix}$ with $\mathbf{H} = \text{diag}\{\mathbf{I}, h \mathbf{I}\}$.

Theorem 5.1: (*Local Bahadur Representation*) Under assumptions (C1)- (C9), we have

$$\widehat{\boldsymbol{\theta}} = \frac{[\Omega_1^*(u_0)]^{-1}}{\sqrt{n h} f_u(u_0)} \sum_{t=1}^n \psi_{\tau}(Y_t^*) \mathbf{X}_t^* K(U_{th}) + o_p(1), \quad (5.6)$$

where $\Omega_1^*(u_0) = \text{diag}\{\Omega^*(u_0), \mu_2 \Omega^*(u_0)\}$.

Remark 3: From Theorem 5.1 and Lemma 5.5 (in Section 5.4), it is easy to see that the local linear estimator $\widehat{\mathbf{a}}(u_0)$ is consistent with the optimal nonparametric convergence rate $\sqrt{n h}$.

Theorem 5.2: (*Asymptotic Normality*) Under assumptions (C1)- (C11), we have the following asymptotic normality

$$\sqrt{n h} \left[\mathbf{H} \begin{pmatrix} \widehat{\mathbf{a}}(u_0) - \mathbf{a}(u_0) \\ \widehat{\mathbf{a}}'(u_0) - \mathbf{a}'(u_0) \end{pmatrix} - \frac{h^2}{2} \begin{pmatrix} \mathbf{a}''(u_0) \mu_2 \\ \mathbf{0} \end{pmatrix} + o_p(h^2) \right] \rightarrow N\{0, \boldsymbol{\Sigma}(u_0)\},$$

where $\boldsymbol{\Sigma}(u_0) = \text{diag}\{\tau(1 - \tau) \nu_0 \boldsymbol{\Sigma}_a(u_0), \tau(1 - \tau) \nu_2 \boldsymbol{\Sigma}_a(u_0)\}$ with

$$\boldsymbol{\Sigma}_a(u_0) = [\Omega^*(u_0)]^{-1} \Omega(u_0) [\Omega^*(u_0)]^{-1} / f_u(u_0). \quad (5.7)$$

In particular,

$$\sqrt{n h} \left[\widehat{\mathbf{a}}(u_0) - \mathbf{a}(u_0) - \frac{h^2 \mu_2}{2} \mathbf{a}''(u_0) + o_p(h^2) \right] \rightarrow N\{0, \tau(1 - \tau) \nu_0 \boldsymbol{\Sigma}_a(u_0)\}.$$

Remark 4: From Theorem 5.2, the asymptotic mean squares error (AMSE) of $\widehat{\mathbf{a}}(u_0)$ is given by

$$AMSE = \frac{h^4 \mu_2^2}{4} \|\mathbf{a}''(u_0)\|^2 + \frac{\tau(1 - \tau) \nu_0}{n h f_u(u_0)} \text{tr}(\boldsymbol{\Sigma}_a(u_0)),$$

which gives the optimal bandwidth h_{opt} by minimizing the AMSE

$$h_{opt} = \left(\frac{\tau(1 - \tau) \nu_0 \text{tr}(\boldsymbol{\Sigma}_a(u_0))}{f_u(u_0) \|\mathbf{a}''(u_0)\|^2} \right)^{1/5} n^{-1/5}$$

and the optimal AMSE is

$$AMSE_{opt} = \frac{5}{4} \left(\frac{\tau(1 - \tau) \nu_0 \text{tr}(\boldsymbol{\Sigma}_a(u_0))}{f_u(u_0)} \right)^{4/5} \|\mathbf{a}''(u_0)\|^{2/5} n^{-4/5}.$$

Further, notice that the similar results in Theorem 5.2 were obtained by Honda (2004) for the independent data. Finally, it is interesting to note that the asymptotic bias in Theorem

5.2 is the same as that for the mean regression case but the two asymptotic variances are different; see, for example, Cai, Fan and Yao (2000).

If model (5.3) does not have \mathbf{X} ($d = 0$), it becomes the nonparametric quantile regression model $q_\tau(\cdot)$. Then, we have the following asymptotic normality for the local linear estimator of the nonparametric quantile regression function $q_\tau(\cdot)$, which covers the results in Yu and Jones (1998), Honda (2000), Lu, Hui and Zhao (2000), and Cai (2002a) for both the independent and time series data.

Corollary 5.1: *If there is no \mathbf{X}_t in (5.3), then,*

$$\sqrt{n h} \left[\hat{q}_\tau(u_0) - q_\tau(u_0) - \frac{h^2 \mu_2}{2} q_\tau''(u_0) + o_p(h^2) \right] \rightarrow N \{0, \sigma_\tau^2(u_0)\},$$

where $\sigma_\tau^2(u_0) = \tau(1 - \tau) \nu_0 f_u^{-1}(u_0) f_{y|u}^{-2}(q_\tau(u_0))$.

Now we consider the comparison of the performance of the local linear estimation $\hat{\mathbf{a}}(u_0)$ obtained in (5.4) with that of the local constant estimation $\tilde{\mathbf{a}}(u_0)$ given in (5.5). To this effect, first, we derive the asymptotic results for the local constant estimator but the proof is omitted since it is along the same line with the proof of Theorems 5.1 and 5.2; see Xu (2005) for details. Under some regularity conditions, it can be shown that

$$\sqrt{n h} \left[\tilde{\mathbf{a}}(u_0) - \mathbf{a}(u_0) - \tilde{\mathbf{b}} + o_p(h^2) \right] \rightarrow N \{0, \tau(1 - \tau) \nu_0 \Sigma_a(u_0)\},$$

where

$$\tilde{\mathbf{b}} = \frac{h^2 \mu_2}{2} \left[\mathbf{a}''(u_0) + 2 \mathbf{a}'(u_0) f_u'(u_0)/f_u(u_0) + 2 \{\Omega^*(u_0)\}^{-1} \Omega^{*'}(u_0) \mathbf{a}'(u_0) \right],$$

which implies that the asymptotic bias for $\tilde{\mathbf{a}}(u_0)$ is different from that for $\hat{\mathbf{a}}(u_0)$ but both have the same asymptotic variance. Therefore, the local constant quantile estimator does not adapt to nonuniform designs: the bias can be large when $f_u'(u_0)/f_u(u_0)$ or $\{\Omega^*(u_0)\}^{-1} \Omega^{*'}(u_0)$ is large even when the true coefficient functions are linear. It is surprising that to the best of our knowledge, this finding seems to be new for the nonparametric quantile regression setting although it is well documented in literature for the ordinary regression case; see Fan and Gijbels (1996) for details.

Finally, to examine the asymptotic behaviors of the local linear and local constant quantile estimators at the boundaries, we offer Theorem 5.3 below but its proofs are omitted due

to their similarity to those for Theorem 5.2 with some modifications and for the ordinary regression setting (Fan and Gijbels, 1996); see Xu (2005) for the detailed proofs. Without loss of generality, we consider only the left boundary point $u_0 = ch$, $0 < c < 1$, if U_t takes values only from $[0, 1]$. A similar result in Theorem 5.3 holds for the right boundary point $u_0 = 1 - ch$. Define $\mu_{j,c} = \int_{-c}^1 u^j K(u) du$ and $\nu_{j,c} = \int_{-c}^1 u^j K^2(u) du$.

Theorem 5.3: (Asymptotic Normality) *Under assumptions of Theorem 5.2, we have the following asymptotic normality of the local linear quantile estimator at the left boundary point,*

$$\sqrt{n}h \left[\hat{\mathbf{a}}(ch) - \mathbf{a}(ch) - \frac{h^2 b_c}{2} \mathbf{a}''(0+) + o_p(h^2) \right] \rightarrow N \{0, \tau(1-\tau) v_c \Sigma_a(0+)\},$$

where

$$b_c = \frac{\mu_{2,c}^2 - \mu_{1,c} \mu_{3,c}}{\mu_{2,c} \mu_{0,c} - \mu_{1,c}^2} \quad \text{and} \quad v_c = \frac{\mu_{2,c}^2 \nu_{0,c} - 2 \mu_{1,c} \mu_{2,c} \nu_{1,c} + \mu_{1,c}^2 \nu_{2,c}}{[\mu_{2,c} \mu_{0,c} - \mu_{1,c}^2]^2}.$$

Further, we have the following asymptotic normality of the local constant quantile estimator at the left boundary point $u_0 = ch$ for $0 < c < 1$,

$$\sqrt{n}h \left[\tilde{\mathbf{a}}(ch) - \mathbf{a}(ch) - \tilde{\mathbf{b}}_c + o_p(h^2) \right] \rightarrow N \{0, \tau(1-\tau) \nu_{0,c} \Sigma_a(0+)/\mu_{0,c}^2\}.$$

where

$$\tilde{\mathbf{b}}_c = \left[h \mu_{1,c} \mathbf{a}'(0+) + \frac{h^2 \mu_{2,c}}{2} \left\{ \mathbf{a}''(0+) + \frac{2 \mathbf{a}'(0+) f'_u(0+)}{f_u(0+)} + 2 \Omega^{*-1}(0+) \Omega^{*'}(0+) \mathbf{a}'(0+) \right\} \right] / \mu_{0,c}.$$

Similar results hold for the right boundary point $u_0 = 1 - ch$.

Remark 5: We remark that if the point 0 were an interior point, then, Theorem 5.3 would hold with $c = 1$, which becomes Theorem 5.2. Also, as $c \rightarrow 1$, $b_c \rightarrow \mu_2$, and $v_c \rightarrow \nu_0$ and these limits are exactly the constant factors appearing respectively in the asymptotic bias and variance for an interior point. Therefore, Theorem 5.3 shows that the local linear estimation has the automatic good behavior at boundaries without the need of boundary correction. Further, one can see from Theorem 5.3 that at the boundaries, the asymptotic bias term for the local constant quantile estimate is of the order h by comparing to the order h^2 for the local linear quantile estimate. This shows that the local linear quantile estimate does not suffer from boundary effects but the local constant quantile estimate does, which is another advantage of the local linear quantile estimator over the local constant quantile estimator. This suggests that one should use the local linear approach in practice.

As a special case, Theorem 5.3 includes the asymptotic properties for the local constant quantile estimator of the nonparametric quantile function $q_\tau(\cdot)$ at both the interior and boundary points, stated as follows.

Corollary 5.2: *If there is no \mathbf{X}_t in (5.3), then, the asymptotic normality of the local constant quantile estimator is given by*

$$\sqrt{n}h \left[\tilde{q}_\tau(u_0) - q_\tau(u_0) - \frac{h^2\mu_2}{2} \{q''_\tau(u_0) + 2q'_\tau(u_0)f'_u(u_0)/f_u(u_0)\} + o_p(h^2) \right] \rightarrow N \{0, \sigma_\tau^2(u_0)\}.$$

Further, at the left boundary point, we have

$$\sqrt{n}h \left[\tilde{q}_\tau(ch) - q_\tau(ch) - \tilde{b}_c^* + o_p(h^2) \right] \rightarrow N \{0, \sigma_c^2\},$$

where

$$\tilde{b}_c^* = \left[h\mu_{1,c}q'_\tau(0+) + \frac{h^2\mu_{2,c}}{2} \{q''_\tau(0+) + 2q'_\tau(0+)f'_u(0+)/f_u(0+)\} \right] / \mu_{0,c}$$

and $\sigma_c^2 = \tau(1-\tau)\nu_{0,c}f_u^{-1}(0+)f_{y|u}^{-2}(q_\tau(0+))/\mu_{0,c}^2$.

5.2.3 Bandwidth Selection

It is well known that the bandwidth plays an essential role in the trade-off between reducing bias and variance. To the best of our knowledge, there has been almost nothing done about selecting the bandwidth in the context of estimating the coefficient functions in the quantile regression even though there is a rich amount of literature on this issue in the mean regression setting; see, for example, Cai, Fan and Yao (2000). In practice, it is desirable to have a quick and easily implemented data-driven fashioned method. Based on this spirit, Yu and Jones (1998) or Yu and Lu (2004) proposed a simple and convenient method for the nonparametric quantile estimation. Their approach assumes that the second derivatives of the quantile function are parallel. However, this assumption might not be valid for many applications in economics and finance due to (nonlinear) heteroscedasticity. Further, the mean regression approach can not directly estimate the variance function. To attenuate these problems, we propose a method of selecting bandwidth for the foregoing estimation procedure, based on the nonparametric version of the Akaike information criterion (AIC), which can attend to the structure of time series data and the over-fitting or under-fitting tendency. This idea is motivated by its analogue of Cai and Tiwari (2000) and Cai (2002b) for nonlinear time series models. The basic idea is described below.

By recalling the classical AIC for linear models under the likelihood setting

$$-2 (\text{maximized log likelihood}) + 2 (\text{number of estimated parameters}),$$

we propose the following nonparametric version of the bias-corrected AIC, due to Hurvich and Tsai (1989) for parametric models and Hurvich, Simonoff and Tsai (1998) for nonparametric regression models, to select h by minimizing

$$\text{AIC}(h) = \log \{\hat{\sigma}_\tau^2\} + 2(p_h + 1)/[n - (p_h + 2)], \quad (5.8)$$

where $\hat{\sigma}_\tau^2$ and p_h are defined later. This criterion may be interpreted as the AIC for the local quantile smoothing problem and seems to perform well in some limited applications. Note that similar to (5.8), Koenker, Ng and Portnoy (1994) considered the Schwarz information criterion (SIC) of Schwarz (1978) with the second term on the right-hand side of (5.8) replayed by $2n^{-1}p_h \log n$, where p_h is the number of “active knots” for the smoothing spline quantile setting, and Machado (1993) studied similar criteria for parametric quantile regression models and more general M-estimators of regression.

Now the question is how to define $\hat{\sigma}_\tau^2$ and p_h in this setting. In the mean regression setting, $\hat{\sigma}_\tau^2$ is just the estimate of the variance σ^2 . In the quantile regression, we define $\hat{\sigma}_\tau^2$ as $n^{-1} \sum_{t=1}^t \rho_\tau(Y_t - \mathbf{X}_t' \hat{\mathbf{a}}(U_t))$, which may be interpreted as the mean square error in the least square setting and was also used by Koenker, Ng and Portnoy (1994). In nonparametric models, p_h is the nonparametric version of degrees of freedom, called the effective number of parameters, and it is usually based on the trace of various quasi-projection (hat) matrices in the least square theory (linear estimators); see, for example, Hastie and Tibshirani (1990), Cai and Tiwari (2000), and Cai (2002b) for a cogent discussion for nonparametric regression models and nonlinear time series models. For the quantile smoothing setting, the explicit expression for the quasi-projection matrix does not exist due to its nonlinearity. However, we can use the first order approximation (the local Bahadur representation) given in (5.6) to derive an explicit expression, which may be interpreted as the quasi-projection matrix in this setting. To this end, define

$$\mathbf{S}_n = \mathbf{S}_n(u_0) = a_n \sum_{t=1}^n \xi_t \mathbf{X}_t^* \mathbf{X}_t^{*'} K(U_{th}),$$

where $\xi_t = I(Y_t \leq \mathbf{X}_t' \mathbf{a}(u_0) + a_n) - I(Y_t \leq \mathbf{X}_t' \mathbf{a}(u_0))$ and $a_n = (nh)^{-1/2}$. It is shown in Section 5.5 that

$$\mathbf{S}_n(u_0) = f_u(u_0) \Omega_1^*(u_0) + o_p(1). \quad (5.9)$$

From (5.6), it is easy to verify that $\widehat{\boldsymbol{\theta}} \approx a_n \mathbf{S}_n^{-1} \sum_{t=1}^n \psi_\tau(Y_t^*) \mathbf{X}_t^* K(U_{th})$. Then, we have

$$\widehat{q}_\tau(U_t, \mathbf{X}_t) - q_\tau(U_t, \mathbf{X}_t) \approx \frac{1}{n} \sum_{s=1}^n \psi_\tau(Y_s^*(U_t)) K_h((U_s - U_t)/h) \mathbf{X}_t^{0'} \mathbf{S}_n^{-1}(U_t) \mathbf{X}_s^*$$

where $\mathbf{X}_t^0 = \begin{pmatrix} \mathbf{X}_t \\ \mathbf{0} \end{pmatrix}$. The coefficient of $\psi_\tau(Y_s^*(U_t))$ on the right-hand side of the above expression is $\gamma_s = a_n^2 K(0) \mathbf{X}_s^{0'} \mathbf{S}_n^{-1}(U_s) \mathbf{X}_s^0$. Now, we have that $p_h = \sum_{s=1}^n \gamma_s$, which can be regarded as an approximation to the trace of the quasi-projection (hat) matrix for linear estimators. In the practical implementation, we need to estimate $\mathbf{a}(u_0)$ first since $\mathbf{S}_n(u_0)$ involves $\mathbf{a}(u_0)$. We recommend using a pilot bandwidth which can be chosen as the one proposed by Yu and Jones (1998). Similar to the least square theory, as expected, the criterion proposed in (5.8) counteracts the over-fitting tendency of the generalized cross-validation due to its relatively weak penalty and the under-fitting of the SIC of Schwarz (1978) studied by Koenker, Ng and Portnoy (1994) because of the heavy penalty.

5.2.4 Covariance Estimate

For the purpose of statistical inference, we next consider the estimation of the asymptotic covariance matrix to construct the pointwise confidence intervals. In practice, a quick and simple way to estimate the asymptotic covariance matrix is desirable. In view of (5.7), the explicit expression of the asymptotic covariance provides a direct estimator. Therefore, we can use the so-called “sandwich” method. In other words, we need to obtain a consistent estimate for both $\Omega(u_0)$ and $\Omega^*(u_0)$. To this effect, define,

$$\widehat{\Omega}_{n,0} = \frac{1}{n} \sum_{t=1}^n \mathbf{X}_t \mathbf{X}_t' K_h(U_t - u_0) \quad \text{and} \quad \widehat{\Omega}_{n,1} = \frac{1}{n} \sum_{t=1}^n w_t \mathbf{X}_t \mathbf{X}_t' K_h(U_t - u_0),$$

where $w_t = I(\mathbf{X}_t' \widehat{\mathbf{a}}(u_0) - \delta_n < Y_t \leq \mathbf{X}_t' \widehat{\mathbf{a}}(u_0) + \delta_n) / (2\delta_n)$ for any $\delta_n \rightarrow 0$ as $n \rightarrow \infty$. It is shown in Section 5.5 that

$$\widehat{\Omega}_{n,0} = f_u(u_0) \Omega(u_0) + o_p(1) \quad \text{and} \quad \widehat{\Omega}_{n,1} = f_u(u_0) \Omega^*(u_0) + o_p(1). \quad (5.10)$$

Therefore, the consistent estimate of $\Sigma_a(u_0)$ is given by

$$\widehat{\Sigma}_a(u_0) = \left[\widehat{\Omega}_{n,1}(u_0) \right]^{-1} \widehat{\Omega}_{n,0}(u_0) \left[\widehat{\Omega}_{n,1}(u_0) \right]^{-1}.$$

Note that $\widehat{\Omega}_{n,1}(u_0)$ might be close to singular for some sparse regions. To avoid this computational difficulty, there are two alternative ways to construct a consistent estimate of

$f_u(u_0)\Omega^*(u_0)$ through estimating the conditional density of Y , $f_{y|u,x}(q_\tau(u, \mathbf{x}))$. The first method is the Nadaraya-Watson type (or local linear) double kernel method of Fan, Yao and Tong (1996) defined as,

$$\hat{f}_{y|u,x}(q_\tau(u, \mathbf{x})) = \sum_{t=1}^n K_{h_2}(U_t - u, \mathbf{X}_t - \mathbf{x}) L_{h_1}(Y_t - q_\tau(u, \mathbf{x})) / \sum_{t=1}^n K_{h_2}(U_t - u, \mathbf{X}_t - \mathbf{x}),$$

where $L(\cdot)$ is a kernel function, and the second one is the difference quotients method of Koenker and Xiao (2004) such as

$$\hat{f}_{y|u,x}(q_\tau(u, \mathbf{x})) = (\tau_j - \tau_{j-1}) / [q_{\tau_j}(u, \mathbf{x}) - q_{\tau_{j-1}}(u, \mathbf{x})],$$

for some appropriately chosen sequence of $\{\tau_j\}$; see Koenker and Xiao (2004) for more discussions. Then, in view of the definition of $f_u(u_0)\Omega^*(u_0)$, the estimator $\tilde{\Omega}_{n,1}$ can be constructed as,

$$\tilde{\Omega}_{n,1} = \frac{1}{n} \sum_{t=1}^n \hat{f}_{y|u,x}(\hat{q}_\tau(U_t, \mathbf{X}_t)) \mathbf{X}_t \mathbf{X}_t' K_h(U_t - u_0).$$

By an analogue of (5.10), one can show that under some regularity conditions, both estimators are consistent.

5.3 Empirical Examples

In this section we report a Monte Carlo simulation to examine the finite sample property of the proposed estimator and to further explore the possible nonlinearity feature, heteroscedasticity, and predictability of the exchange rate of the Japanese Yen per US dollar and to identify the factors affecting the house price in Boston. In our computation, we use the Epanechnikov kernel $K(u) = 0.75(1 - u^2)I(|u| \leq 1)$ and construct the pointwise confidence intervals based on the consistent estimate of the asymptotic covariance described in Section 2.4 without the bias correction. For a predetermined sequence of h 's from a wide range, say from h_a to h_b with an increment h_δ , based on the AIC bandwidth selector described in Section 2.3, we compute $\text{AIC}(h)$ for each h and choose h_{opt} to minimize $\text{AIC}(h)$.

5.3.1 A Simulated Example

Example 5.1: We consider the following data generating process

$$Y_t = a_1(U_t) Y_{t-1} + a_2(U_t) Y_{t-2} + \sigma(U_t) e_t, \quad t = 1, \dots, n, \quad (5.11)$$

where $a_1(U_t) = \sin(\sqrt{2}\pi U_t)$, $a_2(U_t) = \cos(\sqrt{2}\pi U_t)$, and $\sigma(U_t) = 3 \exp(-4(U_t - 1)^2) + 2 \exp(-5(U_t - 2)^2)$. U_t is generated from uniform $(0, 3)$ independently and $e_t \sim N(0, 1)$. The quantile regression is

$$q_\tau(U_t, Y_{t-1}, Y_{t-2}) = a_0(U_t) + a_1(U_t)Y_{t-1} + a_2(U_t)Y_{t-2},$$

where $a_0(U_t) = \Phi^{-1}(\tau)\sigma(U_t)$ and $\Phi^{-1}(\tau)$ is the τ -th quantile of the standard normal. Therefore, only $a_0(\cdot)$ is a function of τ . Note that $a_0(\cdot) = 0$ when $\tau = 0.5$. To assess the performance of finite samples, we compute the mean absolute deviation errors (MADE) for $\hat{a}_j(\cdot)$, which is defined as

$$MADE_j = n_0^{-1} \sum_{k=1}^{n_0} |\hat{a}_j(u_k) - a_j(u_k)|,$$

where $\hat{a}_j(\cdot)$ is either the local linear or local constant quantile estimate of $a_j(\cdot)$ and $\{z_k = 0.1(k-1) + 0.2 : 1 \leq k \leq n_0 = 27\}$ are the grid points. The Monte Carlo simulation is repeated 500 times for each sample size $n = 200, 500$, and 1000 and for each $\tau = 0.05, 0.50$ and 0.95. We compute the optimal bandwidth for each replication, sample size, and τ . We compute the median and standard deviation (in parentheses) of 500 MADE values for each scenario and summarize the results in Table 5.3.1.

From Table 5.3.1, we can observe that the MADE values for both the local linear and local constant quantile estimates decrease when n increases for all three values of τ and the local linear estimate outperforms the local constant estimate. This is another example to show that the local linear method is superior over the local constant even in the quantile setting. Also, the performance for the median quantile estimate is slightly better than that for two tails ($\tau = 0.05$ and 0.95). This observation is not surprising because of the sparsity of data in the tailed regions. Moreover, another benefit of using the quantile method is that we can obtain the estimate of $a_0(\cdot)$ (conditional standard deviation) simultaneously with the estimation of $a_1(\cdot)$ and $a_2(\cdot)$ (functions in the conditional mean), which, in contrast, avoids a two-stage approach needed to estimate the variance function in the mean regression; see Fan and Yao (1998) for details. However, it is interesting to see that due to the larger variation, the performance for $a_0(\cdot)$, although it is reasonably good, is not as good as that of $a_1(\cdot)$ and $a_2(\cdot)$. This can be further evidenced from Figure 1. The results in this simulated experiment show that the proposed procedure is reliable and they are along the line of our asymptotic theory.

Table 5.1: The Median and Standard Deviation of 500 MADE Values

The Local Linear Estimator									
	$\tau = 0.05$			$\tau = 0.5$			$\tau = 0.95$		
n	MADE ₀	MADE ₁	MADE ₂	MADE ₀	MADE ₁	MADE ₂	MADE ₀	MADE ₁	MADE ₂
200	0.911 (0.520)	0.186 (0.041)	0.177 (0.041)	0.401 (0.091)	0.092 (0.032)	0.089 (0.032)	0.920 (0.517)	0.187 (0.042)	0.175 (0.039)
500	0.510 (0.414)	0.085 (0.023)	0.083 (0.02)	0.311 (0.056)	0.055 (0.019)	0.055 (0.018)	0.517 (0.390)	0.085 (0.023)	0.083 (0.023)
1000	0.419 (0.071)	0.060 (0.018)	0.059 (0.017)	0.311 (0.051)	0.050 (0.014)	0.049 (0.014)	0.416 (0.072)	0.060 (0.017)	0.059 (0.017)

The Local Constant Estimator									
	$\tau = 0.05$			$\tau = 0.5$			$\tau = 0.95$		
n	MADE ₀	MADE ₁	MADE ₂	MADE ₀	MADE ₁	MADE ₂	MADE ₀	MADE ₁	MADE ₂
200	3.753 (2.937)	0.285 (0.050)	0.290 (0.051)	0.501 (0.115)	0.144 (0.027)	0.147 (0.028)	3.763 (3.188)	0.287 (0.052)	0.287 (0.051)
500	2.201 (3.025)	0.147 (0.024)	0.146 (0.025)	0.355 (0.062)	0.084 (0.016)	0.085 (0.015)	2.223 (3.320)	0.147 (0.025)	0.147 (0.025)
1000	0.883 (0.462)	0.086 (0.015)	0.086 (0.014)	0.322 (0.054)	0.060 (0.012)	0.061 (0.011)	0.882 (0.427)	0.086 (0.015)	0.087 (0.015)

Finally, Figure 5.1 plots the local linear estimates for all three coefficient functions with their true values (solid line): $\sigma(\cdot)$ in Figure 5.1(a), $a_1(\cdot)$ in Figure 5.1(b), and $a_2(\cdot)$ in Figure 5.1(c), for three quantiles $\tau = 0.05$ (dashed line), 0.50 (dotted line) and 0.95 (dotted-dashed line), for $n = 500$ based on a typical sample which is chosen based on its MADE value equal to the median of the 500 MADE values. The selected optimal bandwidths are $h_{opt} = 0.10$ for $\tau = 0.05$, 0.075 for $\tau = 0.50$, and 0.10 for $\tau = 0.95$. Note that the estimate of $\sigma(\cdot)$ for $\tau = 0.50$ can not be recovered from the estimate of $a_0(\cdot) = 0$ and it is not presented in Figure 5.1(a). The 95% point-wise confidence intervals without the bias correction are depicted in Figure 1 in thick lines for the $\tau = 0.05$ quantile estimate. By the same token, we can compute the point-wise confidence intervals (not shown here) for the rest. Basically, all confidence intervals cover the true values. Also, we can see that the confidence interval for $\hat{a}_0(\cdot)$ is wider than that for $\hat{a}_1(\cdot)$ and $\hat{a}_2(\cdot)$ due to the larger variation. Similar plots

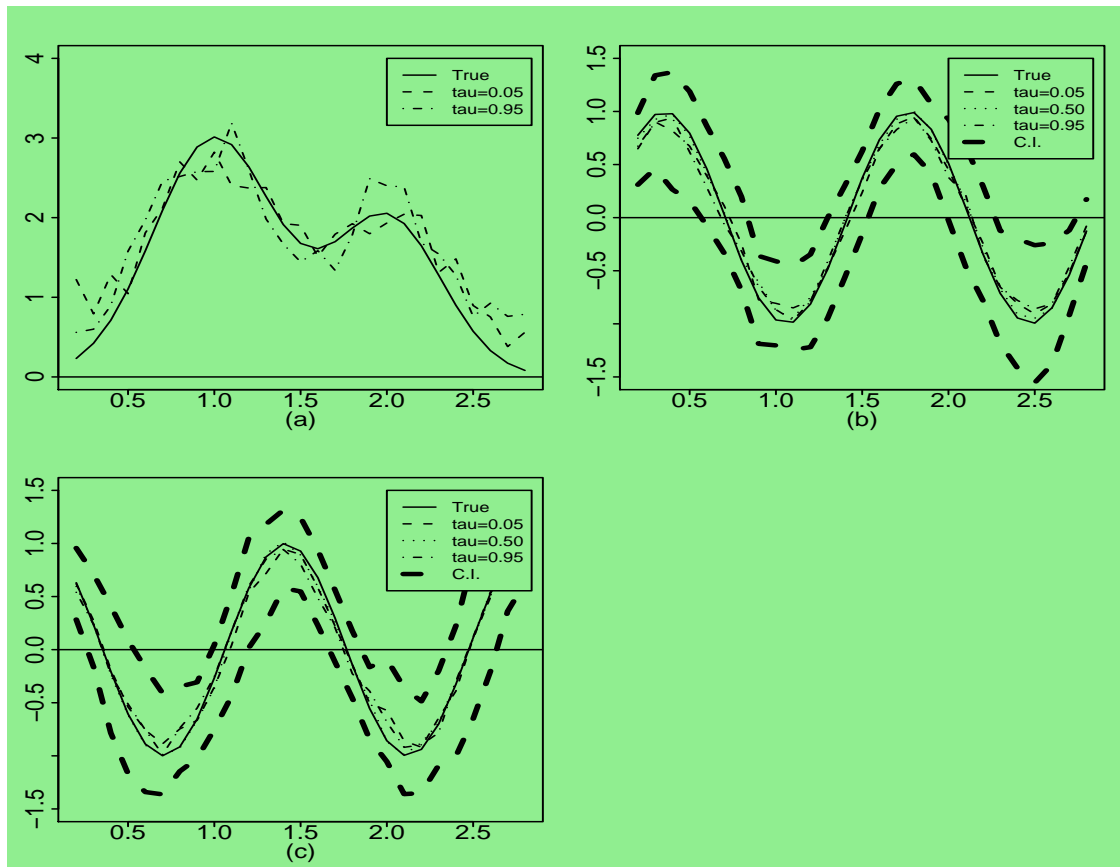


Figure 5.1: *Simulated Example*: The plots of the estimated coefficient functions for three quantiles $\tau = 0.05$ (dashed line), $\tau = 0.50$ (dotted line), and $\tau = 0.95$ (dot-dashed line) with their true functions (solid line): $\sigma(u)$ versus u in (a), $a_1(u)$ versus u in (b), and $a_2(u)$ versus u in (c), together with the 95% point-wise confidence interval (thick line) with the bias ignored for the $\tau = 0.5$ quantile estimate.

are obtained (not shown here) for the local constant estimates due to the space limitations. Overall, the proposed modeling procedure performs fairly well.

5.3.2 Real Data Examples

Example 5.2: (*Boston House Price Data*) We analyze a subset of the Boston house price data (available at <http://lib.stat.cmu.edu/datasets/boston>) of Harrison and Rubinfeld (1978). This dataset consists of 14 variables collected on each of 506 different houses from a variety of locations. The dependent variable is Y , the median value of owner-occupied homes in \$1,000's (house price); some major factors affecting the house prices used are: proportion of population of lower educational status (i.e. proportion of adults with high

school education and proportion of male workers classified as labors), denoted by U , the average number of rooms per house in the area, denoted by X_1 , the per capita crime rate by town, denoted by X_2 , the full property tax rate per \$10,000, denoted by X_3 , and the pupil/teacher ratio by town school district, denoted by X_4 . For the complete description of all 14 variables, see Harrison and Rubinfeld (1978). Gilley and Pace (1996) provided corrections and examined censoring. Recently, there have been several papers devoted to the analysis of this dataset. For example, Breiman and Friedman (1985), Chaudhuri, Doksum and Samarov (1997), and Opsomer and Ruppert (1998) used four covariates: X_1 , X_3 , X_4 and U or their transformations to fit the data through a mean additive regression model whereas Yu and Lu (2004) employed the additive quantile technique to analyze the data. Further, Pace and Gilley (1997) added the georeferencing factor to improve estimation by a spatial approach. Recently, Şentürk and Müller (2005) studied the correlation between the house price Y and the crime rate X_2 adjusted by the confounding variable U through a varying coefficient model and they concluded that the expected effect of increasing crime rate on declining house prices seems to be only observed for lower educational status neighborhoods in Boston. Some existing analyses (e.g., Breiman and Friedman, 1985; Yu and Lu, 2004) in both mean and quantile regressions concluded that most of the variation seen in housing prices in the restricted data set can be explained by two major variables: X_1 and U . Indeed, the correlation coefficients between Y and U and X_1 are -0.7377 and 0.6954 respectively. The scatter plots of Y versus U and X_1 are displayed in Figures 5.2(a) and 5.2(b) respectively. The interesting features of this data set are that the response variable is the median price of a home in a given area and the distributions of Y and the major covariate U are left skewed (the density estimates are not presented). Therefore, quantile methods are particularly well suited to the analysis of this dataset. Finally, it is surprising that all the existing nonparametric models aforementioned above did not include the crime rate X_2 , which may be an important factor affecting the housing price, and did not consider the interaction terms such as U and X_2 .

Based on the above discussions, it concludes that the model studied in this chapter might be well suitable to the analysis of this dataset. Therefore, we analyze this dataset by the

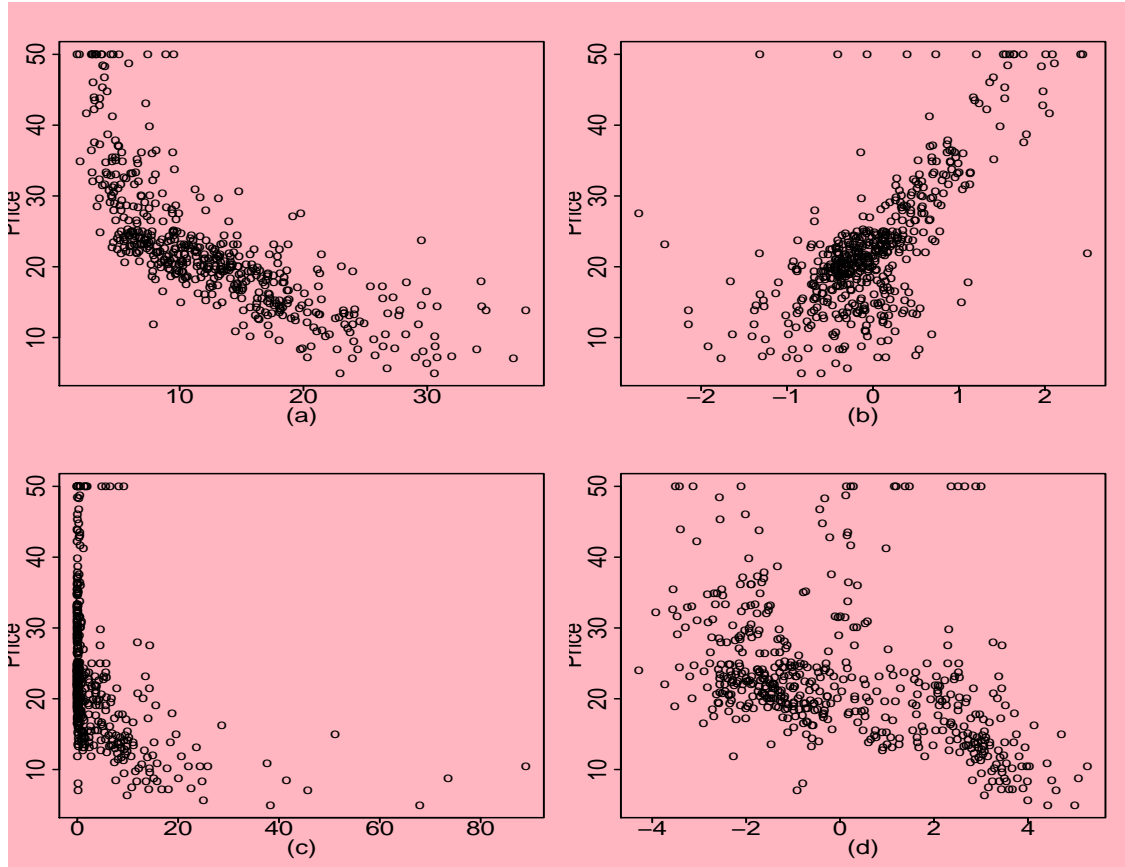


Figure 5.2: *Boston Housing Price Data*: Displayed in (a)-(d) are the scatter plots of the house price versus the covariates U , X_1 , X_2 and $\log(X_2)$, respectively.

following quantile smooth coefficient model¹

$$q_\tau(U_t, \mathbf{X}_t) = a_{0,\tau}(U_t) + a_{1,\tau}(U_t) X_{t1} + a_{2,\tau}(U_t) X_{t2}^*, \quad 1 \leq t \leq n = 506, \quad (5.12)$$

where $X_{t2}^* = \log(X_{t2})$. The reason for using the logarithm of X_{t2} in (5.12), instead of X_{t2} itself, is that the correlation between Y_t and X_{t2}^* (the correlation coefficient is -0.4543) is slightly stronger than that for Y_t and X_{t2} (-0.3883), which can be witnessed as well from Figures 5.2(c) and 5.2(d). In the model fitting, covariates X_1 and X_2 are centralized. For the purpose of comparison, we also consider the following functional coefficient model in the mean regression

$$Y_t = a_0(U_t) + a_1(U_t) X_{t1} + a_2(U_t) X_{t2}^* + e_t \quad (5.13)$$

¹We do not include the other variables such as X_3 and X_4 in model (5.12), since we found that the coefficient functions for these variables seem to be constant. Therefore, a semiparametric model would be appropriate if the model includes these variables. It of course deserves a further investigation.

and we employ the local linear fitting technique to estimate the coefficient functions $\{a_j(\cdot)\}$, denoted by $\{\hat{a}_j(\cdot)\}$; see Cai, Fan and Yao (2000) for details.

The coefficient functions are estimated through the local linear quantile approach by using the bandwidth selector described in Section 2.3. The selected optimal bandwidths are

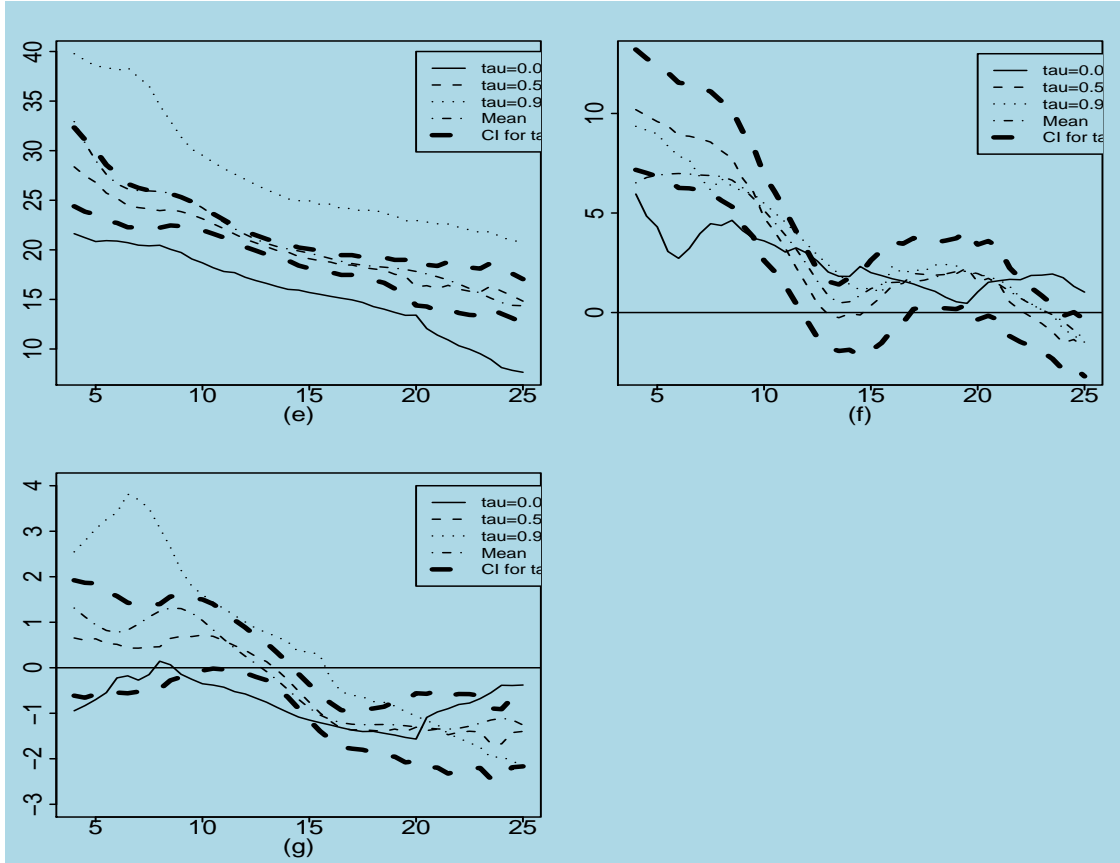


Figure 5.3: *Boston Housing Price Data*: The plots of the estimated coefficient functions for three quantiles $\tau = 0.05$ (solid line), $\tau = 0.50$ (dashed line), and $\tau = 0.95$ (dotted line), and the mean regression (dot-dashed line): $\hat{a}_{0,\tau}(u)$ and $\hat{a}_0(u)$ versus u in (e), $\hat{a}_{1,\tau}(u)$ and $\hat{a}_1(u)$ versus u in (f), and $\hat{a}_{2,\tau}(u)$ and $\hat{a}_2(u)$ versus u in (g). The thick dashed lines indicate the 95% point-wise confidence interval for the median estimate with the bias ignored.

$h_{opt} = 2.0$ for $\tau = 0.05$, 1.5 for $\tau = 0.50$, and 3.5 for $\tau = 0.95$. Figures (5.3(e), (5.3(f) and (5.3(g) present the estimated coefficient functions $\hat{a}_{0,\tau}(\cdot)$, $\hat{a}_{1,\tau}(\cdot)$, and $\hat{a}_{2,\tau}(\cdot)$ respectively, for three quantiles $\tau = 0.05$ (solid line), 0.50 (dashed line) and 0.95 (dotted line), together with the estimates $\{\hat{a}_j(\cdot)\}$ from the mean regression model (dot-dashed line). Also, the 95% point-wise confidence intervals for the median estimate are displayed by the thick dashed lines without the bias correction. First, from these three figures, one can see that the

median estimates are quite close to the mean estimates and the estimates based on the mean regression are always within the 95% confidence interval of the median estimates. It can be concluded that the distribution of the measurement error e_t in (5.13) might be symmetric and $\hat{a}_{j,0.5}(\cdot)$ in (5.12) is almost same as $\hat{a}_j(\cdot)$ in (5.13). Also, one can observe from Figure 5.3(e) that three quantile curves are parallel, which implies that the intercept in $\hat{a}_{0,\tau}(\cdot)$ depends on τ , and they decrease exponentially, which can support that the logarithm transformation may be needed as argued in Yu and Lu (2004). More importantly, one can observe from Figures 5.3(f) and 5.3(g) that three quantile estimated coefficient curves are intersect. This reveals that the structure of quantiles is complex and the lower and upper quantiles have different behaviors and the heteroscedasticity might exist. But unfortunately, this phenomenon was not observed in any previous analyses in the aforementioned papers.

From Figure 5.3(f), first, we can observe that $\hat{a}_{1,0.50}(\cdot)$ and $\hat{a}_{1,0.95}(\cdot)$ are almost same but $\hat{a}_{1,0.05}(\cdot)$ is different. Secondly, we can see that the correlation between the house price and the number of rooms per house is almost positive except for houses with the median price and/or higher than ($\tau = 0.50$ and 0.95) in very low educational status neighborhoods ($U > 23$). Thirdly, for the low price houses ($\tau = 0.05$), the correlation is always positive and it decreases when U is between 0 and 14 and then keeps almost constant afterwards. This implies that the expected effect of increasing the number of rooms can make the house price slightly higher in any low educational status neighborhoods but much higher in relatively high educational status neighborhoods. Finally, for the median and/or higher price houses, the correlation decreases when U is between 0 and 14 and then keeps almost constant until U up to 20 and finally decreases again afterwards, and it becomes negative for U larger than 23. This means that the number of room has a positive effect on the median and/or higher price houses in relatively high and low educational status neighborhoods but increasing the number of rooms might not increase the house price in very low educational status neighborhoods. In other words, it is very difficult to sell high price houses with high number of rooms at a reasonable price in very low educational status neighborhoods.

From Figure 5.3(g), first, one can conclude that the overall trend for all curves is decreasing with $\hat{a}_{3,0.95}(\cdot)$ decreasing faster than the others, and that $\hat{a}_{3,0.05}(\cdot)$ and $\hat{a}_{3,0.50}(\cdot)$ tend to be constant for U larger than 16. Secondly, the correlation between the housing prices ($\tau = 0.50$ and 0.95) and the crime rate seems to be positive for smaller U values (about $U \leq 13$) and becomes negative afterwards. This positive correlation between the housing prices ($\tau = 0.50$

and 0.95) and the crime rate for relatively high educational status neighborhoods seems against intuitive. However, the reason for this positive correlation is the existence of high educational status neighborhoods close to central Boston where high house prices and crime rate occur simultaneously. Therefore, the expected effect of increasing crime rate on declining house prices for $\tau = 0.50$ and 0.95 seems to be observed only for lower educational status neighborhoods in Boston. Finally, it can be seen that the correlation between the housing prices for $\tau = 0.05$ and the crime rate is almost negative although the degree depends on the value of U . This implies that increasing crime rate slightly decreases relatively the house prices for the cheap houses ($\tau = 0.05$).

In summary, it concludes that there is a nonlinear relationship between the conditional quantiles of the housing price and the affecting factors. It seems that the factors U , X_1 and X_2 do have different effects on the different quantiles of the conditional distribution of the housing price. Overall, the housing price and the proportion of population of lower educational status have a strong negative correlation, and the number of rooms has a mostly positive effect on the housing price whereas the crime rate has the most negative effect on the housing price. In particular, by using the proportion of population of lower educational status U as the confounding variable, we demonstrate the substantial benefits obtained by characterizing the affecting factors X_1 and X_2 on the housing price based on the neighborhoods.

Example 5.3: (*Exchange Rate Data*) This example concerns the closing bid prices of the Japanese Yen (JPY) in terms of US dollar. There is a vast amount of literature devoted to the study of the exchange rate time series; see Sercu and Uppal (2000) and the references therein for details. Here we use the proposed model and its modeling approaches to explore the possible nonlinearity feature, heteroscedasticity, and predictability of the exchange rate series. The data is a weekly series from January 1, 1974 to December 31, 2003. The daily noon buying rates in New York City certified by the Federal Reserve Bank of New York for customs and cable transfers purposes were obtained from the Chicago Federal Reserve Board (www.frbchi.org). The weekly series is generated by selecting the Wednesdays series (if a Wednesday is a holiday then the following Thursday is used), which has 1566 observations. The use of weekly data avoids the so-called weekend effect as well as other biases associated with nontrading, bid-ask spread, asynchronous rates and so on, which are often present in higher frequency data. The previous analysis of this “particularly difficult”

data set can be found in Gallant, Hsieh and Tauchen (1991), Fan, Yao and Cai (2003),

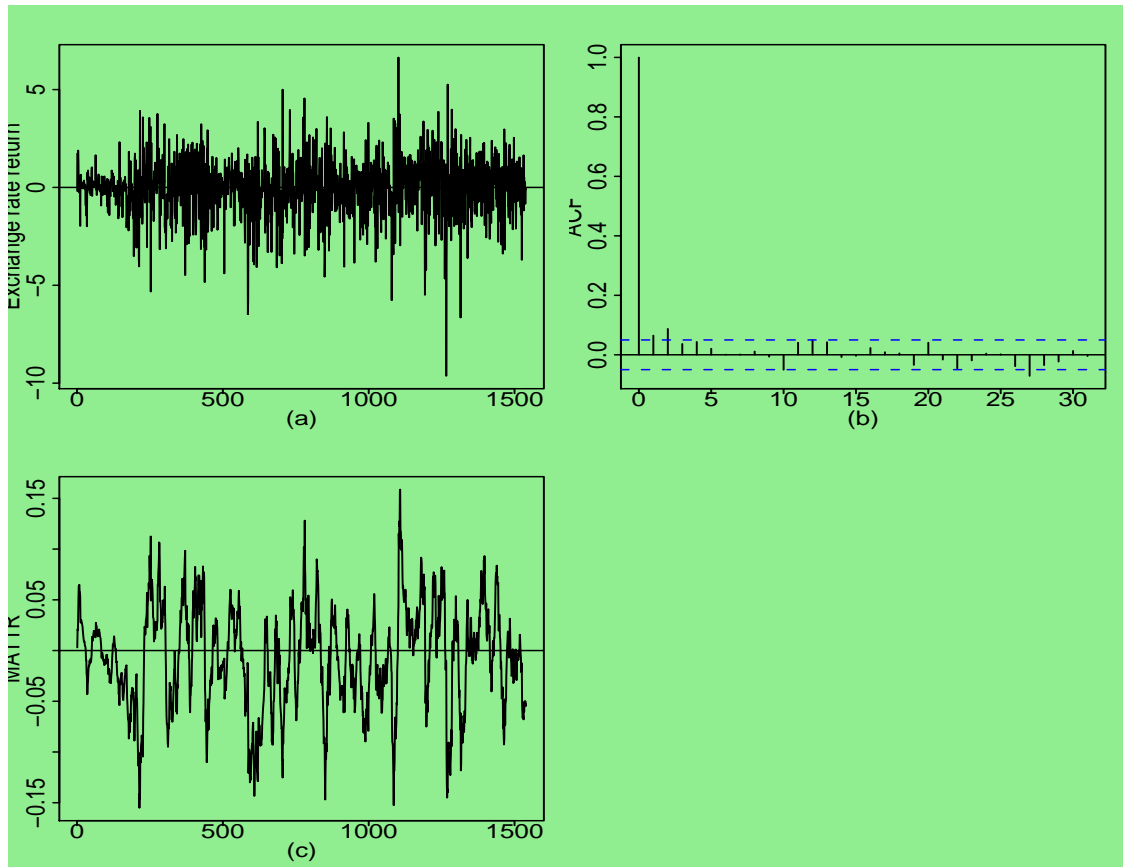


Figure 5.4: *Exchange Rate Series*: (a) Japanese-dollar exchange rate return series $\{Y_t\}$; (b) autocorrelation function of $\{Y_t\}$; (c) moving average trading technique rule.

and Hong and Lee (2003), and the references within. We model the return series $Y_t = 100 \log(\xi_t/\xi_{t-1})$, plotted in Figure 5.4(a), using the techniques developed in this chapter, where ξ_t is an exchange rate level on the t -th week. Typically the classical financial theory would treat $\{Y_t\}$ as a martingale difference process. Therefore, Y_t would be unpredictable. But this assumption was strongly rejected by Hong and Lee (2003) by examining five major currencies and applying several testing procedures. Note that the return series $\{Y_t\}$ has 1565 observations. Figure 5.4(b) shows that there exists almost no significant autocorrelation in $\{Y_t\}$, which also was confirmed by Tsay (2002) and Hong and Lee (2003) by using several statistical testing procedures.

Based on the evidence from Fan, Yao and Cai (2003) and Hong and Lee (2003), the

exchange rate series is predictable by using the functional coefficient autoregressive model

$$Y_t = a_0(U_t) + \sum_{j=1}^d a_j(U_t) Y_{t-j} + \sigma_t e_t, \quad (5.14)$$

where U_t is the smooth variable defined later and σ_t is a function of U_t and the lagged variables. If $\{U_t\}$ is observable, $a_j(\cdot)$ can be estimated by a local linear fitting; see Cai, Fan and Yao (2000) for details, denoted by $\hat{a}_j(\cdot)$. Here, σ_t is the stochastic volatility which may depend on U_t and the lagged variables $\{Y_{t-j}\}$. Now the question is how to choose U_t . Usually, U_t can be chosen based on the knowledge of data or economic theory. However, if no prior information is available, U_t may be chosen as a function of explanatory vector $\{\xi_{t-j}\}$ or through the use of data-driven methods such as AIC or cross-validation. Recently, Fan, Yao and Cai (2003) proposed a data-driven method to the choice of U_t by a linear combination of $\{\xi_{t-j}\}$ and the lagged variables $\{Y_{t-j}\}$. By following the analysis of Fan, Yao and Cai (2003) and Hong and Lee (2003), we choose the smooth variable U_t as an moving average technical trading rule (MATTR) in finance so that the autoregressive coefficients vary with investment positions. U_t is defined as $U_t = \xi_{t-1}/M_t - 1$, where $M_t = \sum_{j=1}^L \xi_{t-j}/L$, which is the moving average and can be regarded as a proxy for the trend at the time $t - 1$. Similar to Hong and Lee (2003), We choose $L = 26$ (half a year). $U_t + 1$ is the ratio of the exchange rate at the time $t - 1$ to the average rate of the most recent L periods of exchange rates at time $t - 1$. The time series plot of $\{U_t\}$ is given in Figure 5.4(c). As pointed out by Hong and Lee (2003), U_t is expected to reveal some useful information on the direction of changes. The MATTR signals 1 (the position to *buy* JPY) when $U_t > 0$ and -1 (the position to *sell* JPY) when $U_t < 0$. For the detailed discussions of the MATTR, see (for example) the papers by LeBaron (1997, 1999), Hong and Lee (2003), Fan, Yao and Cai (2003), and the reference therein. Note that model (5.12) was studied by Fan, Yao and Cai (2003) for the daily data and Hong and Lee (2003) for the weekly data under the homogenous assumption (assume that $\sigma_t = \sigma$) based on the least square theory. In particular, Hong and Lee (2003) provided some empirical evidences to conclude that model (5.14) outperforms the martingale model and autoregressive models.

We analyze this exchange rate series by using the smooth coefficient model under the

quantile regression framework with only two lagged variables² as follows

$$q_\tau(U_t, Y_{t-1}, Y_{t-2}) = a_{0,\tau}(U_t) + a_{1,\tau}(U_t) Y_{t-1} + a_{2,\tau}(U_t) Y_{t-2}. \quad (5.15)$$

The first 1540 observations of $\{Y_t\}$ are used for estimation and the last 25 observations are left for prediction. The coefficient functions $\{a_{j,\tau}(\cdot)\}$ are estimated through the local linear quantile approach, denoted by $\{\hat{a}_{j,\tau}(\cdot)\}$. The previous analysis of this “particularly difficult” data set can be found in optimal bandwidths are $h_{opt} = 0.03$ for $\tau = 0.05, 0.025$ for $\tau = 0.50$, and 0.03 for $\tau = 0.95$. Figures 5.5(d) - 5.5(g) depict the estimated coefficient functions $\hat{a}_{0,\tau}(\cdot)$, $\hat{a}_{1,\tau}(\cdot)$, and $\hat{a}_{2,\tau}(\cdot)$ respectively, for three quantiles $\tau = 0.05$ (solid line), 0.50 (dashed line) and 0.95 (dotted line), together with the estimates $\{\hat{a}_j(\cdot)\}$ (dot-dashed line) from the mean regression model in (5.14). Also, the 95% point-wise confidence intervals for the median estimate are displayed by the thick dashed lines without the bias correction.

First, from Figures 5.5(d), 5.5(f) and 5.5(g), we see clearly that the median estimates $\hat{a}_{j,0.50}(\cdot)$ in (5.15) are almost parallel with or close to the mean estimates $\hat{a}_j(\cdot)$ in (5.14) and the mean estimates are almost within the 95% confidence interval of the median estimates. Secondly, $\hat{a}_{0,0.50}(\cdot)$ in Figure 3(d) shows a nonlinear pattern (increasing and then decreasing) and $\hat{a}_{0,0.05}(\cdot)$ and $\hat{a}_{0,0.95}(\cdot)$ in Figure 5.5(e) exhibit nonlinearly (slightly *U*-shape) and symmetrically. More importantly, one can observe from Figures 5.5(f) and 5.5(g) that the lower and upper quantile estimated coefficient curves are intersect and they behave slightly differently. Particularly, from Figure 5.5(g), we observe that $\hat{a}_{2,0.05}(U_t)$ seems to be nonlinear but $\hat{a}_{2,0.95}(U_t)$ looks like constant when $U_t < 0.06$, and both $\hat{a}_{2,0.05}(U_t)$ and $\hat{a}_{2,0.95}(U_t)$ decrease when $U_t > 0.06$. One might conclude that the distribution of the measurement error e_t in (5.14) might not be symmetric about 0 and there exists a nonlinearity in $a_{j,\tau}(\cdot)$. This supports the nonlinearity test of Hong and Lee (2003). Also, our findings lead to the conclusions that the quantile has a complex structure and the heteroscedasticity exists. This observation supports the existing conclusion in literature that the GARCH (generalized ARCH) effects occur in the exchange rate time series; see Engle, Ito and Lin (1990) and Tsay (2002).

Finally, we consider the post-sample forecasting for the last 25 observations based on the local linear quantile estimators which are computed by using the same bandwidths as those used in the model fitting. The 95% nonparametric prediction interval is constructed

²We also considered the models with more than two lagged variables and we found that the conclusions are similar and not reported here.

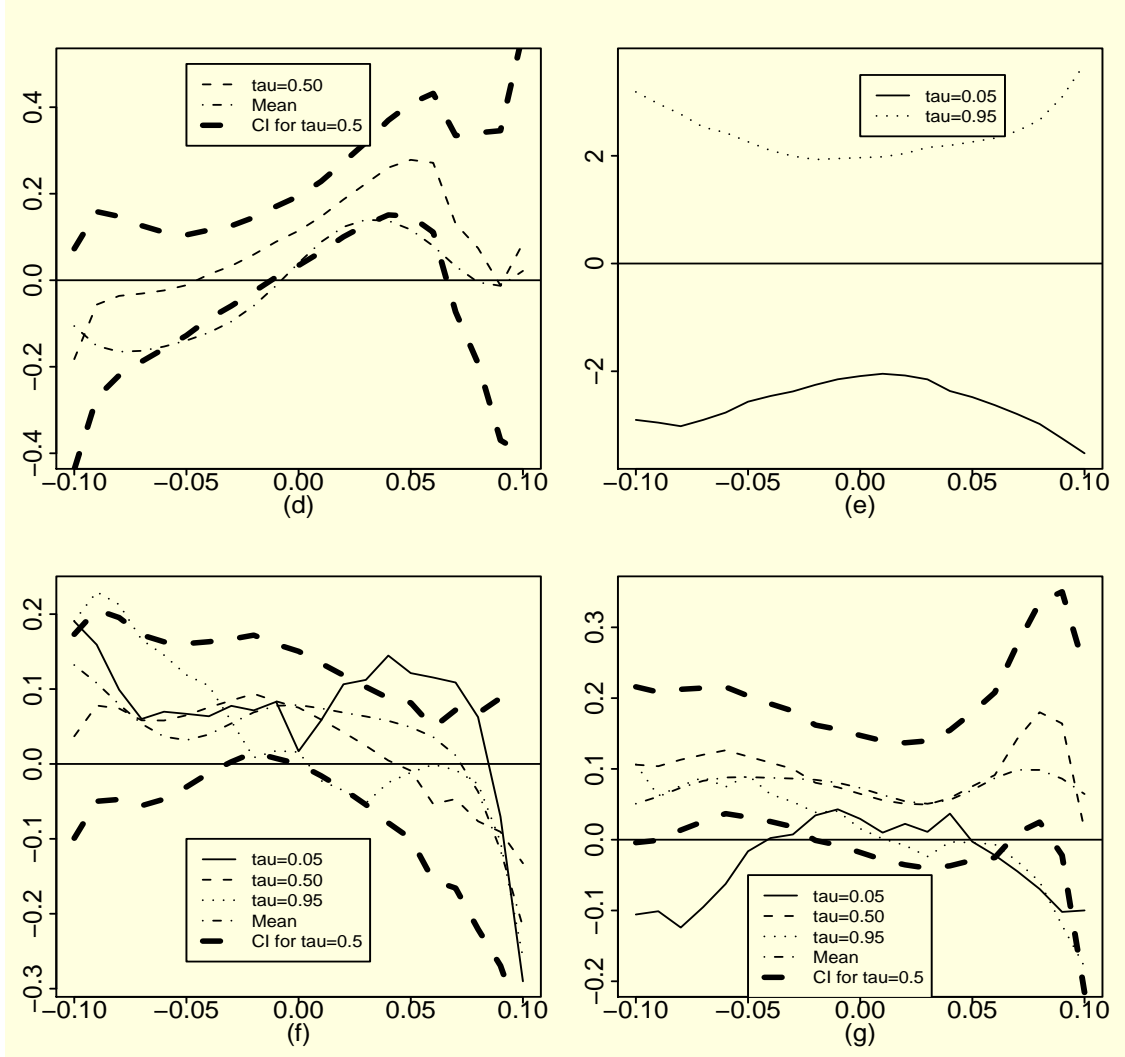


Figure 5.5: *Exchange Rate Series*: The plots of the estimated coefficient functions for three quantiles $\tau = 0.05$ (solid line), $\tau = 0.50$ (dashed line), and $\tau = 0.95$ (dotted line), and the mean regression (dot-dashed line): $\hat{a}_{0,0.50}(u)$ and $\hat{a}_0(u)$ versus u in (d), $\hat{a}_{0,0.05}(u)$ and $\hat{a}_{0,0.95}(u)$ versus u in (e), $\hat{a}_{1,\tau}(u)$ and $\hat{a}_1(u)$ versus u in (f), and $\hat{a}_{2,\tau}(u)$ and $\hat{a}_2(u)$ versus u in (g). The thick dashed lines indicate the 95% point-wise confidence interval for the median estimate with the bias ignored.

as $(\hat{q}_{0.025}(\cdot), \hat{q}_{0.975}(\cdot))$ and the prediction results are reported in Table 2, which shows that 24 out of 25 predictive intervals contain the corresponding true values. The average length of the intervals is 5.77, which is about 35.5% of the range of the data. Therefore, we can conclude that under the dynamic smooth coefficient quantile regression model assumption, the prediction intervals based on the proposed method work reasonably well.

Table 2: The Post-Sample Predictive Intervals For Exchange Rate Data

Observation	True Value	Prediction Interval
Y_{1541}	0.392	(-2.891, 2.412)
Y_{1542}	0.509	(-3.099, 2.405)
Y_{1543}	1.549	(-2.943, 2.446)
Y_{1544}	-0.121	(-2.684, 2.525)
Y_{1545}	-0.991	(-2.677, 2.530)
Y_{1546}	-0.646	(-3.110, 2.401)
Y_{1547}	-0.354	(-3.178, 2.365)
Y_{1548}	-1.393	(-3.083, 2.372)
Y_{1549}	0.997	(-3.110, 2.230)
Y_{1550}	-0.916	(-3.033, 2.431)
Y_{1551}	-3.707	(-3.021, 2.286)
Y_{1552}	-0.919	(-3.841, 2.094)
Y_{1553}	-0.901	(-3.603, 2.770)
Y_{1554}	0.071	(-3.583, 2.821)
Y_{1555}	-0.497	(-3.351, 2.899)
Y_{1556}	-0.648	(-3.436, 2.783)
Y_{1557}	1.648	(-3.524, 2.866)
Y_{1558}	-1.184	(-3.121, 2.810)
Y_{1559}	0.530	(-3.529, 2.531)
Y_{1560}	0.107	(-3.222, 2.648)
Y_{1561}	-0.804	(-3.294, 2.651)
Y_{1562}	0.274	(-3.419, 2.534)
Y_{1563}	-0.847	(-3.242, 2.640)
Y_{1564}	-0.060	(-3.426, 2.532)
Y_{1565}	-0.088	(-3.300, 2.576)

5.4 Derivations

In this section, we give the derivations of the theorems and present certain lemmas with their detailed proofs relegated to Section 5.5. First, we need the following two lemmas.

Lemma 5.1: *Let $V_n(\Delta)$ be a vector function that satisfies*

$$(i) \quad -\Delta' V_n(\lambda \Delta) \geq -\Delta' V_n(\Delta) \text{ for } \lambda \geq 1$$

and

(ii) $\sup_{\|\Delta\| \leq M} \|V_n(\Delta) + \mathbf{D} \Delta - \mathbf{A}_n\| = o_p(1)$, where $\|\mathbf{A}_n\| = O_p(1)$, $0 < M < \infty$, and \mathbf{D} is a positive-definite matrix. Suppose that Δ_n is a vector such that $\|V_n(\Delta_n)\| = o_p(1)$, then, we have

$$(1) \quad \|\Delta_n\| = O_p(1) \quad \text{and} \quad (2) \quad \Delta_n = \mathbf{D}^{-1} \mathbf{A}_n + o_p(1).$$

Proof: The proof follows from Jurečková (1977) and Koenker and Zhao (1996).

Lemma 5.2: Let $\widehat{\boldsymbol{\beta}}$ be the minimizer of the function

$$\sum_{t=1}^n w_t \rho_\tau(y_t - \mathbf{X}_t' \boldsymbol{\beta}),$$

where $w_t > 0$. Then,

$$\|\sum_{t=1}^n w_t \mathbf{X}_t \psi_\tau(y_t - \mathbf{X}_t' \widehat{\boldsymbol{\beta}})\| \leq \dim(\mathbf{X}) \max_{t \leq n} \|w_t \mathbf{X}_t\|.$$

Proof: The proof follows from Ruppert and Carroll (1980).

From the definition of $\boldsymbol{\theta}$, we have

$$\boldsymbol{\beta} = \begin{pmatrix} \mathbf{a}(u_0) \\ \mathbf{a}'(u_0) \end{pmatrix} + a_n \mathbf{H}^{-1} \boldsymbol{\theta},$$

where a_n is defined in (5.10). Then, $Y_t - \sum_{j=0}^q \mathbf{X}_t' \boldsymbol{\beta}_j (U_t - u_0)^j = Y_t^* - a_n \boldsymbol{\theta}' \mathbf{X}_t^*$. Therefore,

$$\widehat{\boldsymbol{\theta}} = \operatorname{argmin} \sum_{t=1}^n \rho_\tau[Y_t^* - a_n \boldsymbol{\theta}' \mathbf{X}_t^*] K(U_{th}) \equiv \operatorname{argmin} G(\boldsymbol{\theta}).$$

Now, define $V_n(\boldsymbol{\theta})$ as

$$V_n(\boldsymbol{\theta}) = a_n \sum_{t=1}^n \psi_\tau[Y_t^* - a_n \boldsymbol{\theta}' \mathbf{X}_t^*] \mathbf{X}_t^* K(U_{th}). \quad (5.16)$$

To establish the asymptotic properties of $\widehat{\boldsymbol{\theta}}$, in the next three lemmas, we show that $V_n(\boldsymbol{\theta})$ satisfies Lemma 5.1 so that we can derive the local Bahadur representation for $\widehat{\boldsymbol{\theta}}$. The results are stated here and their detailed proofs are given in Section 5.5. For the notational convenience define $A_m = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\| \leq M\}$ for some $0 < M < \infty$.

Lemma 5.3: Under assumptions of Theorem 5.1, we have

$$\sup_{\boldsymbol{\theta} \in A_m} \|V_n(\boldsymbol{\theta}) - V_n(0) - E[V_n(\boldsymbol{\theta}) - V_n(0)]\| = o_p(1).$$

Lemma 5.4: *Under assumptions of Theorem 5.1, we have*

$$\sup_{\boldsymbol{\theta} \in A_m} \|E[V_n(\boldsymbol{\theta}) - V_n(0)] + f(u_0)\Omega_1^*(u_0)\boldsymbol{\theta}\| = o(1).$$

Lemma 5.5: *Let $\mathbf{Z}_t = \psi_\tau(Y_t^*) \mathbf{X}_t^* K(U_{th})$. Under assumptions of Theorem 5.1, we have*

$$E[\mathbf{Z}_1] = \frac{h^3 f(u_0)}{2} \begin{pmatrix} \mu_2 \Omega^*(u_0) \mathbf{a}''(u_0) \\ \mathbf{0} \end{pmatrix} \{1 + o(1)\}$$

and

$$\text{Var}[\mathbf{Z}_1] = h \tau(1 - \tau) f(u_0) \Omega_1(u_0) \{1 + o(1)\},$$

where

$$\Omega_1(u_0) = \begin{pmatrix} \nu_0 \Omega(u_0) & \mathbf{0} \\ \mathbf{0} & \nu_2 \Omega(u_0) \end{pmatrix}.$$

Further,

$$\text{Var}[V_n(0)] \rightarrow \tau(1 - \tau) f(u_0) \Omega_1(u_0).$$

Therefore, $\|V_n(0)\| = O_p(1)$.

Now we can embrace the proofs of the theorems.

Proof of Theorem 5.1: By Lemmas 5.5, 5.3, and 5.4, $V_n(\boldsymbol{\theta})$ satisfies the condition (ii) of Lemma 5.1; that is, $\|\mathbf{A}_n\| = O_p(1)$ and $\sup_{\boldsymbol{\theta} \in A_m} \|V_n(\boldsymbol{\theta}) + \mathbf{D}\boldsymbol{\theta} - \mathbf{A}_n\| = o_p(1)$ with $\mathbf{D} = f_u(u_0) \Omega_1^*(u_0)$ and $\mathbf{A}_n = V_n(0)$. It follows Lemma 5.2 that $\|V_n(\hat{\boldsymbol{\theta}})\| = o_p(1)$, where $\hat{\boldsymbol{\theta}}$ is the minimizer of $G(\boldsymbol{\theta})$. Finally, since $\psi_\tau(x)$ is an increasing function of x , then,

$$\begin{aligned} -\boldsymbol{\theta}' V_n(\lambda \boldsymbol{\theta}) &= a_n \sum_{t=1}^n (-\boldsymbol{\theta}') (\psi_\tau(Y_t^* - \lambda a_n \boldsymbol{\theta}' \mathbf{X}_t^*) \mathbf{X}_t^* K(U_{th})) \\ &= a_n \sum_{t=1}^n \psi_\tau[Y_t^* + \lambda a_n (-\boldsymbol{\theta}' \mathbf{X}_t^*)] (-\boldsymbol{\theta}' \mathbf{X}_t^*) K(U_{th}) \end{aligned}$$

is an increasing function of λ . Thus, the condition (i) of Lemma 5.1 is satisfied. Therefore, it follows that

$$\hat{\boldsymbol{\theta}} = \mathbf{D}^{-1} \mathbf{A}_n + o_p(1) = \frac{(\Omega_1^*)^{-1}}{\sqrt{n} h f_u(u_0)} \sum_{t=1}^n \psi_\tau(Y_t^*) \mathbf{X}_t^* K(U_{th}) + o_p(1). \quad (5.17)$$

This proves (5.6). □

Proof of Theorem 5.2: Let $\varepsilon_t = \psi_\tau(Y_t - \mathbf{X}_t' \mathbf{a}(U_t))$. Then, $E(\varepsilon_t) = 0$ and $\text{Var}(\varepsilon_t) = \tau(1 - \tau)$. From (5.17),

$$\hat{\boldsymbol{\theta}} \approx \frac{(\Omega_1^*)^{-1}}{\sqrt{n} h f_u(u_0)} \sum_{t=1}^n [\psi_\tau(Y_t^*) - \varepsilon_t] \mathbf{X}_t^* K(U_{th}) + \frac{(\Omega_1^*)^{-1}}{\sqrt{n} h f_u(u_0)} \sum_{t=1}^n \varepsilon_t \mathbf{X}_t^* K(U_{th}) \equiv \mathbf{B}_n + \boldsymbol{\xi}_n.$$

Similar to the proof of Theorem 2 in Cai, Fan and Yao (2000), by using the small-block and large-block technique and the Cramér-Wold device, one can show that

$$\boldsymbol{\xi}_n \rightarrow N(\mathbf{0}, \boldsymbol{\Sigma}(u_0)). \quad (5.18)$$

By the stationarity and Lemma 5.5,

$$E[\mathbf{B}_n] = \frac{(\Omega_1^*)^{-1}}{\sqrt{n} h f_u(u_0)} n E[\mathbf{Z}_1] \{1 + o(1)\} = a_n^{-1} \frac{h^2}{2} \begin{pmatrix} \mathbf{a}''(u_0) \mu_2 \\ \mathbf{0} \end{pmatrix} \{1 + o(1)\}. \quad (5.19)$$

Since $\psi_\tau(Y_t^*) - \varepsilon_t = I(Y_t \leq \mathbf{X}_t' \mathbf{a}(U_t)) - I(Y_t \leq \mathbf{X}_t' (\mathbf{a}(u_0) + \mathbf{a}'(u_0)(U_t - u_0)))$, then,

$$[\psi_\tau(Y_t^*) - \varepsilon_t]^2 = I(d_{1t} < Y_t \leq d_{2t}), \quad (5.20)$$

where $d_{1t} = \min(c_{1t}, c_{2t})$ and $d_{2t} = \max(c_{1t}, c_{2t})$ with $c_{1t} = \mathbf{X}_t' \mathbf{a}(U_t)$ and $c_{2t} = \mathbf{X}_t' [\mathbf{a}(u_0) + \mathbf{a}'(u_0)(U_t - u_0)]$. Further,

$$E[\{\psi_\tau(Y_t^*) - \varepsilon_t\}^2 K^2(U_{th}) \mathbf{X}_t^* \mathbf{X}_t^{*'}] = E[\{F_{y|u,x}(d_{2t}) - F_{y|u,x}(d_{1t})\} K^2(U_{th}) \mathbf{X}_t^* \mathbf{X}_t^{*'}] = O(h^3).$$

Thus, $\text{Var}(\mathbf{B}_n) = o(1)$. This, in conjunction with (5.18) and (5.19) and the Slutsky Theorem, proves the theorem. \square

5.5 Proofs of Lemmas

Note that the same notations in Sections 5.2 and 5.4 are used here. Throughout this section, we denote a generic constant by C , which may take different values at different appearances. Let $F_{y|u,x}(y)$ denote the conditional distribution of Y given U and \mathbf{X} .

Proof of Lemma 5.3: First, for any $\boldsymbol{\theta} \in A_m$, we consider the following term

$$V_n(\boldsymbol{\theta}) - V_n(0) = a_n \sum_{t=1}^n [\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)] \mathbf{X}_t^* K(U_{th}) \equiv a_n \sum_{i=1}^n V_{nt}(\boldsymbol{\theta}),$$

where $Y_{nt}^* = Y_t^* - a_n \boldsymbol{\theta}' \mathbf{X}_t^*$ and $V_{nt}(\boldsymbol{\theta}) = V_{nt} = [\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)] \mathbf{X}_t^* K(U_{th}) = (V'_{nt1}, V'_{nt2})'$ with

$$V_{nt1} = [\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)] \mathbf{X}_t K(U_{th}) \quad \text{and} \quad V_{nt2} = [\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)] \mathbf{X}_t U_{th} K(U_{th}).$$

Thus,

$$\begin{aligned} & \|V_n(\boldsymbol{\theta}) - V_n(0) - E[V_n(\boldsymbol{\theta}) - V_n(0)]\| \\ & \leq a_n \left\| \sum_{t=1}^n (V_{nt1} - EV_{nt1}) \right\| + a_n \left\| \sum_{t=1}^n (V_{nt2} - EV_{nt2}) \right\| \equiv V_n^{(1)} + V_n^{(2)}. \end{aligned}$$

Clearly,

$$V_n^{(1)} \equiv a_n \left\| \sum_{t=1}^n (V_{nt1} - EV_{nt1}) \right\| \leq \sum_{i=0}^d \|V_n^{(1i)}\|,$$

where $V_n^{(1i)} = a_n \sum_{t=1}^n (V_{nt1}^{(i)} - EV_{nt1}^{(i)})$ and $V_{nt1}^{(i)} = [\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)] X_{ti} K(U_{th})$, which is the i -th component of V_{nt1} . Then,

$$\begin{aligned} \text{Var}(V_n^{(1i)}) &= a_n^2 E \left\{ \sum_{t=1}^n (V_{nt1}^{(i)} - EV_{nt1}^{(i)}) \right\}^2 \\ &= a_n^2 \left[\sum_{t=1}^n \text{Var}(V_{nt1}^{(i)}) + 2 \sum_{s=1}^{n-1} \left(1 - \frac{s}{n}\right) \text{Cov}(V_{n11}^{(i)}, V_{n(s+1)1}^{(i)}) \right] \\ &\leq \frac{1}{h} \left[\text{Var}(V_{n11}^{(i)}) + 2 \sum_{s=1}^{d_n-1} |\text{Cov}(V_{n11}^{(i)}, V_{n(s+1)1}^{(i)})| + 2 \sum_{s=d_n}^{\infty} |\text{Cov}(V_{n11}^{(i)}, V_{n(s+1)1}^{(i)})| \right] \\ &\equiv J_1 + J_2 + J_3 \end{aligned}$$

for some $d_n \rightarrow \infty$ specified later. For J_3 , use the Davydov's inequality (see, e.g., Corollary A.2 of Hall and Heyde, 1980) to obtain

$$|\text{Cov}(V_{n11}^{(i)}, V_{n(s+1)1}^{(i)})| \leq C \alpha^{1-2/\delta}(s) [E|V_{n11}^{(i)}|^\delta]^{2/\delta}.$$

Similar to (5.20), for any $k > 0$,

$$|\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)|^k = I(d_{3t} < Y_t \leq d_{4t}),$$

where $d_{3t} = \min(c_{2t}, c_{2t} + c_{3t})$ and $d_{4t} = \max(c_{2t}, c_{2t} + c_{3t})$ with $c_{3t} = a_n \boldsymbol{\theta}' \mathbf{X}_t^*$. Therefore, by Assumption (C3), there exists a $C > 0$ independent of $\boldsymbol{\theta}$ such that

$$E \{ |\psi_\tau(Y_{nt}^*) - \psi_\tau(Y_t^*)|^k \mid \mathbf{U}_t, \mathbf{X}_t \} = F_{y|u,x}(c_{4t}) - F_{y|u,x}(c_{3t}) \leq C a_n |\boldsymbol{\theta}' \mathbf{X}_t^*|,$$

which implies that

$$\begin{aligned} E|V_{n11}^{(i)}|^\delta &= E [|\psi_\tau(Y_{n1}^*) - \psi_\tau(Y_1^*)|^\delta |X_{1i}|^\delta K^\delta(U_{1h})] \\ &\leq C a_n E [|\boldsymbol{\theta}' \mathbf{X}_t^*| |X_{1i}|^\delta K^\delta(U_{1h})] \leq C a_n h \end{aligned}$$

uniformly in $\boldsymbol{\theta}$ over A_m by Assumption (C6). Then,

$$J_3 \leq C a_n^{2/\delta} h^{2/\delta-1} \sum_{s=d_n}^{\infty} [\alpha(s)]^{1-2/\delta} \leq C a_n^{2/\delta} h^{2/\delta-1} d_n^{-l} \sum_{s=d_n}^{\infty} s^l [\alpha(s)]^{1-2/\delta} = o(a_n^{2/\delta} h^{2/\delta-1} d_n^{-l})$$

uniformly in $\boldsymbol{\theta}$ over A_m . As for J_2 , we use Assumption (C10) to get

$$|\text{Cov}(V_{n11}^{(i)}, V_{n(s+1)1}^{(i)})| \leq C [E\{|X_{1i} X_{(s+1)i}| K(U_{1h}) K(U_{(s+1)h})\} + a_n^2 h^2] = O(h^2)$$

uniformly in $\boldsymbol{\theta}$ over A_m . It follows that $J_2 = O(d_n h)$ uniformly in $\boldsymbol{\theta}$ over A_m . Analogously,

$$J_1 = h^{-1} \text{Var}(V_{n11}^{(i)}) \leq h^{-1} E(V_{n11}^{(i)})^2 = O(a_n)$$

uniformly in $\boldsymbol{\theta}$ over A_m . By choosing d_n such that $d_n^l h^{1-2/\delta} = c$, then, $d_n h \rightarrow 0$ and $\text{Var}(V_n^{(1i)}) = o(1)$. Therefore, $V_n^{(1i)} = o_p(1)$ so that $V_n^{(1)} = o_p(1)$ uniformly in $\boldsymbol{\theta}$ over A_m . By the same token, we can show that $V_n^{(2)} = o_p(1)$ uniformly in $\boldsymbol{\theta}$ over A_m . This completes the proof of the lemma. \square

Proof of Lemma 5.4: It is easy to justify that

$$\begin{aligned} E[V_n(\boldsymbol{\theta}) - V_n(0)] &= n a_n E[(\psi_\tau(Y_t^* - a_n \boldsymbol{\theta}' \mathbf{X}_t^*) - \psi_\tau(Y_t^*)) \mathbf{X}_t^* K(U_{th})] \\ &= n a_n E[\{F_{y|u,x}(c_{2t}) - F_{y|u,x}(c_{2t} + a_n \boldsymbol{\theta}' \mathbf{X}_t^*)\} \mathbf{X}_t^* K(U_{th})] \\ &\approx -\frac{1}{h} E[f_{y|u,x}(c_{2t}) \mathbf{X}_t^* \mathbf{X}_t^{*'} K(U_{th})] \boldsymbol{\theta} \\ &\approx -f_u(u_0) \Omega_1^*(u_0) \boldsymbol{\theta} \end{aligned}$$

uniformly in $\boldsymbol{\theta}$ over A_m by Assumption (C3). The proof of the lemma is complete. \square

Proof of Lemma 5.5: Observe by Taylor expansions and Assumption (C3) that

$$\begin{aligned} E[\mathbf{Z}_t] &= E[\{\tau - F_{y|u,x}(c_{2t})\} \mathbf{X}_t^* K(U_{th})] \\ &\approx E[\{F_{y|u,x}(c_{2t} + \mathbf{X}_t' \mathbf{a}''(u_0) h^2 U_{th}^2/2) - F_{y|u,x}(c_{2t})\} \mathbf{X}_t^* K(U_{th})] \\ &\approx \frac{h^2}{2} E[f_{y|u,x}(c_{2t}) \mathbf{X}_t^* \mathbf{X}_t' \mathbf{a}''(u_0) U_{th}^2 K(U_{th})] \\ &\approx \frac{h^2}{2} E[f_{y|u,x}(q_\tau(u_0, \mathbf{X}_t)) \mathbf{X}_t^* \mathbf{X}_t' \mathbf{a}''(u_0) U_{th}^2 K(U_{th})] \\ &\approx \frac{h^3 f_u(u_0)}{2} \begin{pmatrix} \mu_2 \Omega^*(u_0) \mathbf{a}''(u_0) \\ \mathbf{0} \end{pmatrix}. \end{aligned} \tag{5.21}$$

Also, we have

$$\begin{aligned}
\text{Var}[\mathbf{Z}_t] &= E[\{\tau - I(Y_t < c_{2t})\}^2 \mathbf{X}_t^* \mathbf{X}_t^{*'} K^2(U_{th})] \\
&\approx E[\{\tau^2 - 2\tau F_{y|u,x}(c_{2t}) + F_{y|u,x}(c_{2t})\} \mathbf{X}_t^* \mathbf{X}_t^{*'} K^2(U_{th})] \\
&\approx \tau(1 - \tau) E[\mathbf{X}_t^* \mathbf{X}_t^{*'} K^2(U_{th})] \\
&\approx \tau(1 - \tau) h f_u(u_0) \Omega_1(u_0).
\end{aligned} \tag{5.22}$$

Next, we show that the last part of lemma holds true. Clearly, $V_n(0) = a_n \sum_{t=1}^n \mathbf{Z}_t$. Similar to the proof of Lemma 5.3, we have

$$\begin{aligned}
\text{Var}[V_n(0)] &= \frac{1}{h} \text{Var}(\mathbf{Z}_1) + \frac{2}{h} \sum_{s=1}^{d_n-1} \left(1 - \frac{s}{n}\right) \text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{s+l}) + \frac{2}{h} \sum_{s=d_n}^n \left(1 - \frac{s}{n}\right) \text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{s+l}) \\
&\equiv J_4 + J_5 + J_6
\end{aligned}$$

for some $d_n \rightarrow \infty$ specified later. By (5.22),

$$J_4 \rightarrow \tau(1 - \tau) f_u(u_0) \Omega_1(u_0).$$

Therefore, it suffices to show that $|J_5| = o(1)$ and $|J_6| = o(1)$. For J_6 , using the Davydov's inequality (see, e.g., Corollary A.2 of Hall and Heyde, 1980) and the boundedness of $\psi_\tau(\cdot)$ to obtain

$$|\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{s+1})| \leq C \alpha^{1-2/\delta}(s) [E|\mathbf{Z}_1|^\delta]^{2/\delta} \leq C h^{2/\delta} \alpha^{1-2/\delta}(s),$$

which gives

$$J_6 \leq C h^{2/\delta-1} \sum_{s=d_n}^{\infty} [\alpha(s)]^{1-2/\delta} \leq C h^{2/\delta-1} d_n^{-l} \sum_{s=d_n}^{\infty} s^l [\alpha(s)]^{1-2/\delta} = o(h^{2/\delta-1} d_n^{-l}) = o(1)$$

by choosing d_n to satisfy $d_n^l h^{1-2/\delta} = c$. As for J_5 , we use Assumption (C10) and (5.21) to get

$$|\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_{s+1})| \leq C [E\{|\mathbf{X}_1^* \mathbf{X}_{s+1}^{*'}| K(U_{1h}) K(U_{(s+1)h})\} + h^6] = O(h^2)$$

so that $J_5 = O(d_n h) = o(1)$ by the choice of d_n . We finish the proof of this lemma. \square

Proof of (5.9) and (5.10): By the Taylor expansion,

$$E[\xi_t | U_t, \mathbf{X}_t] = F_{y|u,x}(\mathbf{X}_t' \mathbf{a}(u_0) + a_n) - F_{y|u,x}(\mathbf{X}_t' \mathbf{a}(u_0)) \approx f_{y|u,x}(\mathbf{X}_t' \mathbf{a}(u_0)) a_n.$$

Therefore,

$$E[\mathbf{S}_n] \approx h^{-1} E[f_{y|u,x}(\mathbf{X}_t' \mathbf{a}(u_0)) \mathbf{X}_t^* \mathbf{X}_t^{*'} K(U_{th})] \approx f_u(u_0) \Omega_1^*(u_0).$$

Similar to the proof of $\text{Var}[V_n(0)]$ in Lemma 5.5, one can show that $\text{Var}(\mathbf{S}_n) \rightarrow 0$. Therefore, $\mathbf{S}_n \rightarrow f_u(u_0) \Omega_1^*(u_0)$ in probability. This proves (5.9). Clearly,

$$E[\widehat{\Omega}_{n,0}] = E[\mathbf{X}_t \mathbf{X}_t' K_h(U_t - u_0)] = \int \Omega(u_0 + h v) f_u(u_0 + h v) K(v) dv \approx f_u(u_0) \Omega(u_0).$$

Similarly, one can show that $\text{Var}(\widehat{\Omega}_{n,0}) \rightarrow 0$. This proves the first part of (5.10). By the same token, one can show that $E[\widehat{\Omega}_{n,1}] \approx f_u(u_0) \Omega^*(u_0)$ and $\text{Var}(\widehat{\Omega}_{n,1}) \rightarrow 0$. Thus, $\widehat{\Omega}_{n,1} = f_u(u_0) \Omega^*(u_0) + o_p(1)$. We prove (5.10). \square

5.6 Computer Codes

Please see the files [chapter5-1.r](#), [chapter5-2.r](#), and [chapter5-3.r](#) for making figures. If you want to learn the codes for computation, they are available upon request.

5.7 References

- An, H.Z. and Chen, S.G. (1997). A Note on the Ergodicity of Nonlinear Autoregressive Models. *Statistics and Probability Letters*, **34**, 365-372.
- An, H.Z. and Huang, F.C. (1996). The Geometrical Ergodicity of Nonlinear Autoregressive Models. *Statistica Sinica*, **6**, 943-956.
- Auestad, B. and Tjøstheim, D. (1990). Identification of nonlinear time series: First order characterization and order determination. *Biometrika*, **77**, 669-687.
- Bao, Y., Lee, T.-H. and Saltoğlu, B. (2001). Evaluating predictive performance of value-at-risk models in emerging markets: a reality check. *Journal of Forecasting*, forthcoming.
- Breiman, L. and Friedman, J. H. (1985). Estimating optimal transformation for multiple regression and correlation. *Journal of the American Statistical Association*, **80**, 580-619.
- Cai, Z. (2002a). Regression quantile for time series. *Econometric Theory*, **18**, 169-192.
- Cai, Z. (2002b). A two-stage approach to additive time series models. *Statistica Neerlandica*, **56**, 415-433.
- Cai, Z. (2007). Trending time-varying coefficient time series models with serially correlated errors. *Journal of Econometrics*, **137**, 163-188.
- Cai, Z., Fan, J. and Yao, Q. (2000). Functional-coefficient regression models for nonlinear time series. *Journal of the American Statistical Association*, **95**, 941-956.
- Cai, Z. and Masry, E. (2000). Nonparametric estimation in nonlinear ARX time series models: Projection and linear fitting. *Econometric Theory*, **16**, 465-501.

- Cai, Z. and Tiwari, R.C. (2000). Application of a local linear autoregressive model to BOD time series. *Environmetrics*, **11**, 341-350.
- Cai, Z. and X. Xu (2005). Nonparametric quantile estimations for dynamic smooth coefficient models. Forthcoming in *Journal of the American Statistical Association*.
- Chaudhuri, P. (1991). Nonparametric estimates of regression quantiles and their local Bahadur representation. *The Annals of Statistics*, **19**, 760-777.
- Chaudhuri, P., Doksum, K. and Samarov, A. (1997). On average derivative quantile regression. *The Annals of Statistics*, **25**, 715-744.
- Chen, R. and Tsay, R.S. (1993). Functional-coefficient autoregressive models. *Journal of the American Statistical Association*, **88**, 298-308.
- Cole, T.J. (1994). Growth charts for both cross-sectional and longitudinal data. *Statistics in Medicine*, **13**, 2477-2492.
- De Gooijer, J. and Zerom, D. (2003). On additive conditional quantiles with high dimensional covariates. *Journal of American Statistical Association*, **98**, 135-146.
- Duffie, D. and Pan, J. (1997). An overview of value at risk. *Journal of Derivatives*, **4**, 7-49.
- Engle, R.F., Ito, T. and Lin, W. (1990). Meteor showers or heat waves? Heteroskedastic intra-daily volatility in the foreign exchange market. *Econometrica*, **58**, 525-542.
- Engle, R.F. and Manganelli, S. (2004). CAViaR: conditional autoregressive value at risk by regression quantile. *Journal of Business and Economics Statistics*, **22**, 367-381.
- Efron, B. (1991). Regression percentiles using asymmetric squared error loss. *Statistica Sinica*, **1**, 93-125.
- Fan, J. and Gijbels, I. (1996). *Local Polynomial Modeling and Its Applications*. Chapman and Hall, London.
- Fan, J. and Yao, Q. (1998). Efficient estimation of conditional variance functions in stochastic regression. *Biometrika*, **85**, 645-660.
- Fan, J., Yao, Q. and Cai, Z. (2003). Adaptive varying-coefficient linear models. *Journal of the Royal Statistical Society, series B*, **65**, 57-80.
- Fan, J., Yao, Q. and Tong, H. (1996). Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika*, **83**, 189-206.
- Gallant, A.R., Hsieh, D.A. and Tauchen, G.E. (1991). On fitting a recalcitrant series: the pound/dollar exchange rate, 1974-1983. In *Nonparametric And Semiparametric Methods in Econometrics and Statistics* (W.A. Barnett, J. Powell and G.E. Tauchen, eds.), pp.199-240. Cambridge: Cambridge University Press.

- Gilley, O.W. and Pace, R.K. (1996). On the Harrison and Rubinfeld Data. *Journal of Environmental Economics and Management*, **31**, 403-405.
- Gorodetskii, V.V. (1977). On the strong mixing property for linear sequences. *Theory of Probability and Its Applications*, **22**, 411-413.
- Granger, C.W.J., White, H. and Kamstra, M. (1989). Interval forecasting: an analysis based upon ARCH-quantile estimators. *Journal of Econometrics*, **40**, 87-96.
- Hall, P. and Heyde, C.C. (1980). *Martingale Limit Theory and its Applications*. Academic Press, New York.
- Harrison, D. and Rubinfeld, D.L. (1978). Hedonic housing prices and demand for clean air. *Journal of Environmental Economics and Management*, **5**, 81-102.
- Hastie, T.J. and Tibshirani, R. (1990). *Generalized Additive Models*. Chapman and Hall, London.
- He, X. and Ng, P. (1999). Quantile splines with several covariates. *Journal of Statistical Planning and Inference*, **75**, 343-352.
- He, X., Ng, P. and Portony, S. (1998). Bivariate quantile smoothing splines. *Journal of the Royal Statistical Society, Series B*, **60**, 537-550.
- He, X. and Portnoy, S. (2000). Some asymptotic results on bivariate quantile splines. *Journal of Statistical Planning and Inference*, **91**, 341-349.
- Honda, T. (2000). Nonparametric estimation of a conditional quantile for α -mixing processes. *Annals of the Institute of Statistical Mathematics*, **52**, 459-470.
- Honda, T. (2004). Quantile regression in varying coefficient models. *Journal of Statistical Planning and Inferences*, **121**, 113-125.
- Hong, Y. and Lee, T.-H. (2003). Inference on via generalized spectrum and nonlinear time series models. *The Review of Economics and Statistics*, **85**, 1048-1062.
- Horowitz, J.L. and Lee, S. (2005). Nonparametric Estimation of an Additive Quantile Regression Model. *Journal of the American Statistical Association*, **100**, 1238-1249.
- Hurvich, C.M., Simonoff, J.S. and Tsai, C.-L. (1998). Smoothing parameter selection in nonparametric regression using an improved Akaike information criterion. *Journal of the Royal Statistical Society, Series B*, **60**, 271-293.
- Hurvich, C.M. and Tsai, C.-L. (1989). Regression and time series model selection in small samples. *Biometrika*, **76**, 297-307.
- Jorion, P. (2000). *Value at Risk*, 2ed. McGraw Hill, New York.
- Jurečkoá, J. (1977). Asymptotic relations of M -estimates and R -estimates in linear regression model. *The Annals of Statistics*, **5**, 464-472.

- Khindanova, I.N. and Rachev, S.T. (2000). Value at risk: Recent advances. *Handbook on Analytic-Computational Methods in Applied Mathematics*, CRC Press LLC.
- Koenker, R. (1994). Confidence intervals for regression quantiles. In *Proceedings of the Fifth Prague Symposium on Asymptotic Statistics* (P. Mandl and M. Huskova, eds.), 349-359. Physica, Heidelberg.
- Koenker, R. (2004). *Quantreg: An R package for quantile regression and related methods* <http://cran.r-project.org>.
- Koenker R. (2000). Galton, Edgeworth, Frisch, and prospects for quantile regression in econometrics. *Journal of Econometrics*, **95**, 347-374.
- Koenker, R. and Bassett, G.W. (1978). Regression quantiles. *Econometrica*, **46**, 33-50.
- Koenker, R. and Bassett, G.W. (1982). Robust tests for heteroscedasticity based on regression quantiles. *Econometrica*, **50**, 43-61.
- Koenker, R. and Hallock, K.F. (2001). Quantile regression: An introduction. *Journal of Economic Perspectives*, **15**, 143-157.
- Koenker, R., Ng, P. and Portnoy, S. (1994). Quantile smoothing splines. *Biometrika*, **81**, 673-680.
- Koenker, R. and Xiao, Z. (2002). Inference on the quantile regression process. *Econometrica*, **70**, 1583-1612.
- Koenker, R. and Xiao, Z. (2004). Unit root quantile autoregression inference. *Journal of American Statistical Association*, **99**, 775-787.
- Koenker, R. and Zhao, Q. (1996). Conditional quantile estimation and inference for ARCH models. *Econometric Theory*, **12**, 793-813.
- LeBaron, B. (1997). Technical trading rule and regime shifts in foreign exchange. In *Advances in Trading Rules* (E. Acar and S. Satchell, eds.). Butterworth-Heinemann.
- LeBaron, B. (1999). Technical trading rule profitability and foreign exchange intervention. *Journal of International Economics*, **49**, 125-143.
- Li, Q. and Racine, J. (2004). Nonparametric estimation of conditional CDF and quantile functions with mixed categorical and continuous data. *Journal of Business and Economic Statistics*, forthcoming.
- Lu, Z. (1998). On the ergodicity of non-linear autoregressive model with an autoregressive conditional heteroscedastic term. *Statistica Sinica*, **8**, 1205-1217.
- Lu, Z., Hui, Y.V. and Zhao, Q. (2000). Local linear quantile regression under dependence: Bahadur representation and application. *Working Paper*, Department of Management Sciences, City University of Hong Kong.

- Masry, E. and Tjøstheim, D. (1995). Nonparametric estimation and identification of nonlinear ARCH time series: Strong convergence and asymptotic normality. *Econometric Theory*, **11**, 258-289.
- Masry, E. and Tjøstheim, D. (1997). Additive nonlinear ARX time series and projection estimates. *Econometric Theory*, **13**, 214-252.
- Machado, J.A.F. (1993). Robust model selection and M -estimation. *Econometric Theory*, **9**, 478-493.
- Morgan, J.P. (1995). *Riskmetrics Technical Manual*, 3ed.
- Opsomer, J.D. and Ruppert, D. (1998). A fully automated bandwidth selection for additive regression model. *Journal of The American Statistical Association*, **93**, 605-618.
- Pace, R.K. and Gilley, O.W. (1997). Using the spatial configuration of the data to improve estimation. *The Journal of Real Estate Finance and Economics*, **14**, 333-340.
- Ruppert, D. and Carroll, R.J. (1980). Trimmed least squares estimation in the linear model. *Journal of The American Statistical Association*, **75**, 828-838.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, **6**, 461-464.
- Sercu, P. and Uppal, R. (2000). *Exchange Rate Volatility, Trade, and Capital Flows under Alternative Rate Regimes*. Cambridge: Cambridge University Press.
- Şentürk, D. and Müller, H.G. (2005). Covariate adjusted correlation analysis. Forthcoming in *Scandinavian Journal of Statistics*.
- Taylor, J.W. and Bunn, D.W. (1999). A quantile regression approach to generating prediction intervals. *Management Science*, **45**, 225-237.
- Tsay, R.S. (2000). Extreme values, quantile estimation and value at risk. *Working paper*, Graduate School of Business, University of Chicago.
- Tsay, R.S. (2002). *Analysis of Financial Time Series*. John Wiley & Sons, New York.
- Wang, K. (2003). Asset pricing with conditioning information: A new test. *Journal of Finance*, **58**, 161-196.
- Wei, Y. and He, X. (2006). Conditional growth charts (with discussion). *The Annals of Statistics*, **34**, 2069-2097.
- Wei, Y., Pere, A., Koenker, R. and He, X. (2006). Quantile regression methods for reference growth charts. *Statistics in Medicine*, **25**, 1369-1382.
- Withers, C.S. (1981). Conditions for linear processes to be strong mixing. *Zeitschrift für Wahrscheinlichkeitstheorie verwandte Gebiete*, **57**, 477-480.

- Yu, K. and Jones, M.C. (1998). Local linear quantile regression. *Journal of the American Statistical Association*, **93**, 228-237.
- Yu, K. and Lu, Z. (2004). Local linear additive quantile regression. *Scandinavian Journal of Statistics*, **31**, 333-346.
- Xu, X. (2005). Semiparametric Quantile Dynamic Time Series Models and Their Applications. *Ph.D. Dissertation*, University of North Carolina at Charlotte.
- Zhou, K.Q. and Portnoy, S.L. (1996). Direct use of regression quantiles to construct confidence sets in linear models. *The Annals of Statistics*, **24**, 287-306.

Chapter 6

Conditional VaR and Expected Shortfall

For details, see the paper by Cai and Wang (2006). If you like to read the whole paper, you can download it from the web site at <http://www.wise.xmu.edu.cn/> at **WORKING PAPER** column. Next we present only a part of the whole paper of Cai and Wang (2006).

6.1 Introduction

The **value-at-risk** (hereafter, VaR) and **expected shortfall** (ES) have become two popular measures on market risk associated with an asset or a portfolio of assets during the last decade. In particular, VaR has been chosen by the Basle Committee on Banking Supervision as the benchmark of risk measures for capital requirements and both of them have been used by financial institutions for asset managements and minimization of risk as well as have been developed rapidly as analytic tools to assess riskiness of trading activities. See, to name just a few, Morgan (1996), Duffie and Pan (1997), Jorion (2001, 2003), and Duffie and Singleton (2003) for the financial background, statistical inferences, and various applications. In terms of the formal definition, **VaR is simply a quantile of the loss distribution (future portfolio values) over a prescribed holding period (e.g., 2 weeks) at a given confidence level, while ES is the expected loss, given that the loss is at least as large as some given quantile of the loss distribution** (e.g., VaR). It is well known from Artzner, Delbaen, Eber and Heath (1999) that ES is a **coherent risk measure** such as it satisfies the following **four axioms**:

- **homogeneity**: increasing the size of a portfolio by a factor should scale its risk measure by the same factor,

- **monotonicity**: a portfolio must have greater risk if it has systematically lower values than another,
- **risk-free condition or translation invariance**: adding some amount of cash to a portfolio should reduce its risk by the same amount, and
- **subadditivity**: the risk of a portfolio must be less than the sum of separate risks or merging portfolios cannot increase risk.

VaR satisfies homogeneity, monotonicity, and risk-free condition but is not sub-additive. See Artzner, *et al.* (1999) for details. As advocated by Artzner, *et al.* (1999), ES is preferred due to its better properties although VaR is widely used in applications.

Measures of risk might depend on the state of the economy since economic and market conditions vary from time to time. This requires risk managers should focus on the conditional distributions of profit and loss, which take full account of current information about the investment environment (macroeconomic and financial as well as political) in forecasting future market values, volatilities, and correlations. As pointed out by Duffie and Singleton (2003), not only are the prices of the underlying market indices changing randomly over time, the portfolio itself is changing, as are the volatilities of prices, the credit qualities of counterparties, and so on. On the other hand, one would expect the VaR to increase as the past returns become very negative, because one bad day makes the probability of the next somewhat greater. Similarly, very good days also increase the VaR, as would be the case for volatility models. Therefore, VaR could depend on the past returns in some way. Hence, an appropriate risk analytical tool or methodology should be allowed to adapt to varying market conditions and to reflect the latest available information in a time series setting rather than the iid framework. Most of the existing risk management literature has concentrated on unconditional distributions and the iid setting although there have been some studies on the conditional distributions and time series data. For more background, see Chernozhukov and Umanstev (2001), Cai (2002), Fan and Gu (2003), Engle and Manganelli (2004), Cai and Xu (2005), Scaillet (2005), and Cosma, Scaillet and von Sachs (2006), and references therein for conditional models, and Duffie and Pan (1997), Artzner, *et al.* (1999), Rockafellar and Uryasev (2000), Acerbi and Tasche (2002), Frey and McNeil (2002), Scaillet (2004), Chen and Tang (2005), Chen (2006), and among others for unconditional models. Also, most of studies in the literature and applications are limited to parametric models, such as all stan-

dard industry models like CreditRisk⁺, CreditMetrics, CreditPortfolio View and the model proposed by the KMV corporation. See Chernozhukov and Umanstev (2001), Frey and McNeil (2002), Engle and Manganelli (2004), and references therein on parametric models in practice and Fan and Gu (2003) and references therein for semiparametric models.

The main focus of this chapter is on studying the **conditional value-at-risk (CVaR)** and **conditional expected shortfall (CES)** and proposing a new nonparametric estimation procedure to estimate CVaR and CES functions where the conditional information is allowed to contain economic and market (exogenous) variables and the past observed returns. Parametric models for CVaR and CES can be most efficient if the underlying functions are correctly specified. See Chernozhukov and Umanstev (2001) for a polynomial type regression model and Engle and Manganelli (2004) for a GARCH type parametric model for CVaR based on regression quantile. However, a misspecification may cause serious bias and model constraints may distort the underlying distributions. A nonparametric modeling is appealing in several aspects. One of the advantages for nonparametric modeling is that little or no restrictive prior information on functionals is needed. Further, it may provide a useful insight for further parametric fitting.

The approach proposed by Cai and Wang (2006) has several advantages. The first one is to propose a new nonparametric approach to estimate CVaR and CES. In essence, our estimator for CVaR is based on inverting a newly proposed estimator of the conditional distribution function for time series data and the estimator for CES is by a plugging-in method based on plugging in the estimated conditional probability density function and the estimated CVaR function. Note that they are analogous to the estimators studied by Scaillet (2005) by using the Nadaraya-Watson (NW) type double kernel (smoothing in both the y and x directions) estimation, and Cai (2002) by utilizing the weighted Nadaraya-Watson (WNW) kernel type technique to avoid the so-called boundary effects as well as Yu and Jones (1998) by employing the double kernel local linear method. More precisely, our newly proposed estimator combines the WNW method of Cai (2002) and the double kernel local linear technique of Yu and Jones (1998), termed as *weighted double kernel local linear* (WDKLL) estimator.

The second merit is to establish the asymptotic properties for the WDKLL estimators of the conditional probability density function (PDF) and cumulative distribution function

(CDF) for the α -mixing time series at both boundary and interior points. It is therefore shown that the WDKLL method enjoys the same convergence rates as those of the double kernel local linear estimator of Yu and Jones (1998) and the WNW estimator of Cai (2002). It is also shown that the WDKLL estimators have desired sampling properties at both boundary and interior points of the support of the design density, which seems to be seminal. Finally, we derive the WDKLL estimator of CVaR by inverting the WDKLL conditional distribution estimator and the WDKLL estimator of CES by plugging in the WDKLL estimators of PDF and CVaR. We show that the WDKLL estimator of CVaR exists always due to the WDKLL estimator of CDF being a distribution function itself, and that it inherits all better properties from the WDKLL estimator of CDF; that is, the WDKLL estimator of CDF is a CDF and differentiable, and it possess the asymptotic properties such as design adaption, avoiding boundary effects, and mathematical efficiency. Note that to preserve shape constraints, recently, Cosma, Scaillet and von Sachs (2006) used a wavelet method to estimate conditional probability density and cumulative distribution functions and then to estimate conditional quantiles.

Note that CVaR defined here is essentially the conditional quantile or quantile regression of Koenker and Bassett (1978), based on the conditional distribution, rather than CVaR defined in some risk management literature (see, e.g., Rockafellar and Uryasev, 2000; Jorion, 2001, 2003) which is what we call ES here. Also, note that the ES here is called TailVaR in Artzner, *et al.* (1999). Moreover, as aforementioned, CVaR can be regarded as a special case of quantile regression. See Cai and Xu (2005) for the state-of-the-art about current research on nonparametric quantile regression, including CVaR. Further, note that both ES and CES have been known for decades among actuary sciences and they are very popular in insurance industry. Indeed, they have been used to assess risk on a portfolio of potential claims, and to design reinsurance treaties. See the book by Embrechts, Kluppelberg, and Mikosch (1997) for the excellent review on this subject and the papers by McNeil (1997), Hürlimann (2003), Scaillet (2005), and Chen (2006). Finally, ES or CES is also closely related to other applied fields such as the mean residual life function in reliability and the biometric function in biostatistics. See Oakes and Dasu (1990) and Cai and Qian (2000) and references therein.

6.2 Setup

Assume that the observed data $\{(X_t, Y_t); 1 \leq t \leq n\}$, $X_t \in \mathbb{R}^d$, are available and they are observed from a stationary time series model. Here Y_t is the risk or loss variable which can be the negative logarithm of return (log loss) and X_t is allowed to include both economic and market (exogenous) variables and the lagged variables of Y_t and also it can be a vector. But, for the expositional purpose, we consider only the case when X_t is a scalar ($d = 1$). Note that the proposed methodologies and their theory for the univariate case ($d = 1$) continue to hold for multivariate situations ($d > 1$). Extension to the case $d > 1$ involves no fundamentally new ideas. Note that models with large d are often not practically useful due to “curse of dimensionality”.

We now turn to considering the nonparametric estimation of the conditional expected shortfall $\mu_p(x)$, which is defined as

$$\mu_p(x) = E[Y_t | Y_t \geq \nu_p(x), X_t = x],$$

where $\nu_p(x)$ is the conditional value-at-risk, which is defined as the solution of

$$P(Y_t \geq \nu_p(x) | X_t = x) = S(\nu_p(x) | x) = p$$

or expressed as $\nu_p(x) = S^{-1}(p | x)$, where $S(y | x)$ is the conditional survival function of Y_t given $X_t = x$; $S(y | x) = 1 - F(y | x)$, and $F(y | x)$ is the conditional cumulative distribution function. It is easy to see that

$$\mu_p(x) = \int_{\nu_p(x)}^{\infty} y f(y | x) dy / p,$$

where $f(y | x)$ is the conditional probability density function of Y_t given $X_t = x$. To estimate $\mu_p(x)$, one can use the plugging-in method as

$$\hat{\mu}_p(x) = \int_{\hat{\nu}_p(x)}^{\infty} y \hat{f}(y | x) dy / p, \quad (6.1)$$

where $\hat{\nu}_p(x)$ is a nonparametric estimation of $\nu_p(x)$ and $\hat{f}(y | x)$ is a nonparametric estimation of $f(y | x)$. But the bandwidths for $\hat{\nu}_p(x)$ and $\hat{f}(y | x)$ are not necessary to be same.

Note that Scaillet (2005) used the NW type double kernel method to estimate $f(y | x)$ first, due to Roussas (1969), denoted by $\tilde{f}(y | x)$, and then estimated $\nu_p(x)$ by inverting

the estimated conditional survival function, denoted by $\tilde{\nu}_p(x)$, and finally estimated $\mu_p(x)$ by plugging $\tilde{f}(y|x)$ and $\tilde{\nu}_p(x)$ into (6.1), denoted by $\tilde{\mu}_p(x)$, where $\tilde{\nu}_p(x) = \tilde{S}^{-1}(y|x)$ and $\tilde{S}(y|x) = \int_y^\infty \tilde{f}(u|x)du$. But, it is well documented (see, e.g., Fan and Gijbels, 1996) that the NW kernel type procedures have serious drawbacks: the asymptotic bias involves the design density so that they can not be adaptive, and boundary effects exist so that they require boundary modifications. In particular, boundary effects might cause a serious problem for estimating $\nu_p(x)$ since it is only concerned with the tail probability. The question is now how to provide a better estimate for $f(y|x)$ and $\nu_p(x)$ so that we have a good estimate for $\mu_p(x)$. Therefore, we address this issue in the next section.

6.3 Nonparametric Estimating Procedures

We start with the nonparametric estimators for the conditional density function and its distribution function first and then turn to discussing the nonparametric estimators for the conditional VaR and ES functions.

There are several methods available for estimating $\nu_p(x)$, $f(y|x)$, and $F(y|x)$ in the literature, such as kernel and nearest-neighbor.¹ To attenuate these drawbacks of the kernel type estimators mentioned in Section 6.2, recently, some new methods have been proposed to estimate conditional quantiles. The first one, a more direct approach, by using the “check” function such as the robustified local linear smoother, was provided by Fan, Hu, and Troung (1994) and further extended by Yu and Jones (1997, 1998) for iid data. A more general nonparametric setting was explored by Cai and Xu (2005) for time series data. This modeling idea was initialed by Koenker and Bassett (1978) for linear regression quantiles and Fan, Hu, and Troung (1994) for nonparametric models. See Cai and Xu (2005) and references therein for more discussions on models and applications. An alternative procedure is first to estimate the conditional distribution function by using double kernel local linear technique of Fan, Yao, and Tong (1996) and then to invert the conditional distribution estimator to produce an estimator of a conditional quantile or CVaR. Yu and Jones (1997, 1998) compared these two methods theoretically and empirically and suggested that the double kernel local linear would be better.

¹ To name just a few, see Lejeune and Sarda (1988), Troung (1989), Samanta (1989), and Chaudhuri (1991) for iid errors, Roussas (1969) and Roussas (1991) for Markovian processes, and Troung and Stone (1992) and Boente and Fraiman (1995) for mixing sequences.

6.3.1 Estimation of Conditional PDF and CDF

To make a connection between the conditional density (distribution) function and nonparametric regression problem, it is noted by the standard kernel estimation theory (see, e.g., Fan and Gijbels, 1996) that for a given symmetric density function $K(\cdot)$,

$$E\{K_{h_0}(y - Y_t) | X_t = x\} = f(y | x) + \frac{h_0^2}{2} \mu_2(K) f^{2,0}(y | x) + o(h_0^2) \approx f(y | x), \quad \text{as } h_0 \rightarrow 0, \quad (6.2)$$

where $K_{h_0}(u) = K(u/h_0)/h_0$, $\mu_2(K) = \int_{-\infty}^{\infty} u^2 K(u) du$, $f^{2,0}(y | x) = \partial^2 / \partial y^2 f(y | x)$, and \approx denotes an approximation by ignoring the higher terms. Note that $Y_t^*(y) = K_{h_0}(y - Y_t)$ can be regarded as an initial estimate of $f(y | x)$ smoothing in the y direction. Also, note that this approximation ignores the higher order terms $O(h_0^j)$ for $j \geq 2$, since they are negligible if $h_0 = o(h)$, where h is the bandwidth used in smoothing in the x direction (see (6.3) below). Therefore, the smoothing in the y direction is not important in the context of this subject so that intuitively, it should be under-smoothed. Thus, the left hand side of (6.2) can be regraded as a nonparametric regression of the observed variable $Y_t^*(y)$ versus X_t and the local linear (or polynomial) fitting scheme of Fan and Gijbels (1996) can be applied to here. This leads us to consider the following locally weighted least squares regression problem:

$$\sum_{t=1}^n \{Y_t^*(y) - a - b(X_t - x)\}^2 W_h(x - X_t), \quad (6.3)$$

where $W(\cdot)$ is a kernel function and $h = h(n) > 0$ is the bandwidth satisfying $h \rightarrow 0$ and $nh \rightarrow \infty$ as $n \rightarrow \infty$, which controls the amount of smoothing used in the estimation. Note that (6.3) involves two kernels $K(\cdot)$ and $W(\cdot)$. This is the reason of calling “double kernel”.

Minimizing the above locally weighted least squares in (6.3) with respect to a and b , we obtain the locally weighted least squares estimator of $f(y | x)$, denoted by $\hat{f}(y | x)$, which is \hat{a} . From Fan and Gijbels (1996) or Fan, Yao and Tong (1996), $\hat{f}(y | x)$ can be re-expressed as a linear estimator form as

$$\hat{f}_u(y | x) = \sum_{t=1}^n W_{u,t}(x, h) Y_t^*(y),$$

where with $S_{n,j}(x) = \sum_{t=1}^n W_h(x - X_t) (X_t - x)^j$, the weights $\{W_{u,t}(x, h)\}$ are given by

$$W_{u,t}(x, h) = \frac{[S_{n,2}(x) - (x - X_t) S_{n,1}(x)] W_h(x - X_t)}{S_{n,0}(x) S_{n,2}(x) - S_{n,1}^2(x)}.$$

Clearly, $\{W_{u,t}(x, h)\}$ satisfy the so-called discrete moments conditions as follows: for $0 \leq j \leq 1$,

$$\sum_{t=1}^n W_{u,t}(x, h) (X_t - x)^j = \delta_{0,j} = \begin{cases} 1 & \text{if } j = 0 \\ 0 & \text{otherwise} \end{cases} \quad (6.4)$$

based on the least squares theory; see (3.12) of Fan and Gijbels (1996, p.63). Note that the estimator $\hat{f}_u(y|x)$ can range outside $[0, \infty)$. The double kernel local linear estimator of $F(y|x)$ is constructed (see (8) of Yu and Jones (1998)) by integrating $\hat{f}_u(y|x)$

$$\hat{F}_u(y|x) = \int_{-\infty}^y \hat{f}_u(y|x) dy = \sum_{t=1}^n W_{u,t}(x, h) G_{h_0}(y - Y_t),$$

where $G(\cdot)$ is the distribution function of $K(\cdot)$ and $G_{h_0}(u) = G(u/h_0)$. Clearly, $\hat{F}_u(y|x)$ is continuous and differentiable with respect to y with $\hat{F}_u(-\infty|x) = 0$ and $\hat{F}_u(\infty|x) = 1$. Note that the differentiability of the estimated distribution function can make the asymptotic analysis much easier for the nonparametric estimators of CVaR and CES (see later).

Although Yu and Jones (1998) showed that the double kernel local linear estimator has some attractive properties such as no boundary effects, design adaptation, and mathematical efficiency (see, e.g., Fan and Gijbels, 1996), it has the disadvantage of producing conditional distribution function estimators that are not constrained either to lie between zero and one or to be monotone increasing, which is not good for estimating CVaR if the inverting method is used. In both these respects, the NW method is superior, despite its rather large bias and boundary effects. The properties of positivity and monotonicity are particularly advantageous if the method of inverting conditional distribution estimator is applied to produce the estimator of a conditional quantile or CVaR. To overcome these difficulties, Hall, Wolff, and Yao (1999) and Cai (2002) proposed the WNW estimator based on an empirical likelihood principle, which is designed to possess the superior properties of local linear methods such as bias reduction and no boundary effects, and to preserve the property that the NW estimator is always a distribution function, although it might require more computational efforts since it requires estimating and optimizing additional weights aimed at the bias correction. Cai (2002) discussed the asymptotic properties of the WNW estimator at both interior and boundary points for the mixing time series under some regularity assumptions and showed that the WNW estimator has a better performance than other competitors. See Cai (2002) for details. Recently, Cosma, Scaillet and von Sachs (2006) proposed a shape preserving estimation method to estimate cumulative distribution functions and probability

density functions using the wavelet methodology for multivariate dependent data and then to estimate a conditional quantile or CVaR.

The WNW estimator of the conditional distribution $F(y|x)$ of Y_t given $X_t = x$ is defined by

$$\widehat{F}_{c1}(y|x) = \sum_{t=1}^n W_{c,t}(x, h) I(Y_t \leq y), \quad (6.5)$$

where the weights $\{W_{c,t}(x, h)\}$ are given by

$$W_{c,t}(x, h) = \frac{p_t(x) W_h(x - X_t)}{\sum_{t=1}^n p_t(x) W_h(x - X_t)}, \quad (6.6)$$

and $\{p_t(\mathbf{x})\}$ is chosen to be $p_t(x) = n^{-1} \{1 + \lambda (X_t - x) W_h(x - X_t)\}^{-1} \geq 0$ with λ , a function of data and x , uniquely defined by maximizing the logarithm of the empirical likelihood

$$L_n(\lambda) = - \sum_{t=1}^n \log \{1 + \lambda (X_t - x) W_h(x - X_t)\}$$

subject to the constraints $\sum_{t=1}^n p_t(x) = 1$ and the discrete moments conditions in (6.4); that is,

$$\sum_{t=1}^n W_{c,t}(x, h) (X_t - x)^j = \delta_{0,j} \quad (6.7)$$

for $0 \leq j \leq 1$. Also, see Cai (2002) for details on this aspect. In implementation, Cai (2002) recommended using the Newton-Raphson scheme to find the root of equation $L'_n(\lambda) = 0$. Note that $0 \leq \widehat{F}_{c1}(y|x) \leq 1$ and it is monotone in y . But $\widehat{F}_{c1}(y|x)$ is not continuous in y and of course, not differentiable in y either. Note that under regression setting, Cai (2001) provided a comparison of the local linear estimator and the WNW estimator and discussed the asymptotic minimax efficiency of the WNW estimator.

To accommodate all nice properties (monotonicity, continuity, differentiability, and lying between zero and one) and the attractive asymptotic properties (design adaption, avoiding boundary effects, and mathematical efficiency, see Cai (2002) for detailed discussions) of both estimators $\widehat{F}_{ll}(y|x)$ and $\widehat{F}_{c1}(y|x)$ under a unified framework, we propose the following nonparametric estimators for the conditional density function $f(y|x)$ and its conditional distribution function $F(y|x)$, termed as *weighted double kernel local linear estimation*,

$$\widehat{f}_c(y|x) = \sum_{t=1}^n W_{c,t}(x, h) Y_t^*(y),$$

where $W_{c,t}(x, h)$ is given in (6.6), and

$$\widehat{F}_c(y|x) = \int_{-\infty}^y \widehat{f}_c(y|x) dy = \sum_{t=1}^n W_{c,t}(x, h) G_{h_0}(y - Y_t). \quad (6.8)$$

Note that if $p_t(x)$ in (6.6) is a constant for all t , or $\lambda = 0$, then $\widehat{f}_c(y|x)$ becomes the classical NW type double kernel estimator used by Scaillet (2005). However, Scaillet (2005) adopted a single bandwidth for smoothing in both the y and x directions. Clearly, $\widehat{f}_c(y|x)$ is a probability density function so that $\widehat{F}_c(y|x)$ is a cumulative distribution function (monotone, $0 \leq \widehat{F}_c(y|x) \leq 1$, $\widehat{F}_c(-\infty|x) = 0$, and $\widehat{F}_c(\infty|x) = 1$). Also, $\widehat{F}_c(y|x)$ is continuous and differentiable in y . Further, as expected, it will be shown that like $\widehat{F}_{c1}(y|x)$, $\widehat{F}_c(y|x)$ has the attractive properties such as no boundary effects, design adaptation, and mathematical efficiency.

6.3.2 Estimation of Conditional VaR and ES

We now are ready to formulate the nonparametric estimators for $\nu_p(x)$ and $\mu_p(x)$. To this end, from (6.8), $\nu_p(x)$ is estimated by inverting the estimated conditional survival distribution $\widehat{S}_c(y|x) = 1 - \widehat{F}_c(y|x)$, denoted by $\widehat{\nu}_p(x)$ and defined as $\widehat{\nu}_p(x) = \widehat{S}_c^{-1}(p|x)$. Note that $\widehat{\nu}_p(x)$ always exists since $\widehat{S}_c(p|x)$ is a survival function itself. Plugging-in $\widehat{\nu}_p(x)$ and $\widehat{f}_c(y|x)$ into (6.1), we obtain the nonparametric estimation of $\mu_p(x)$,

$$\begin{aligned} \widehat{\mu}_p(x) &= p^{-1} \int_{\widehat{\nu}_p(x)}^{\infty} y \widehat{f}_c(y|x) dy = p^{-1} \sum_{t=1}^n W_{c,t}(x, h) \int_{\widehat{\nu}_p(x)}^{\infty} y K_{h_0}(y - Y_t) dy \\ &= p^{-1} \sum_{t=1}^n W_{c,t}(x, h) [Y_t \bar{G}_{h_0}(\widehat{\nu}_p(x) - Y_t) + h_0 G_{1,h_0}(\widehat{\nu}_p(x) - Y_t)], \end{aligned} \quad (6.9)$$

where $\bar{G}(u) = 1 - G(u)$, $G_{1,h_0}(u) = G_1(u/h_0)$, and $G_1(u) = \int_u^{\infty} v K(v) dv$. Note that as mentioned earlier, $\widehat{\nu}_p(x)$ in (6.9) can be an any consistent estimator.

6.4 Distribution Theory

6.4.1 Assumptions

Before we proceed with the asymptotic properties of the proposed nonparametric estimators, we first list all assumptions needed for the asymptotic theory, although some of them might not be the weakest possible. Note that proofs of the asymptotic results presented in this

section may be found in Section 6.6 with some lemmas and their detailed proofs relegated to Section 6.7. First, we introduce some notation. Let $\alpha(K) = \int_{-\infty}^{\infty} u K(u) \bar{G}(u) du$ and $\mu_j(W) = \int_{-\infty}^{\infty} u^j W(u) du$. Also, for any $j \geq 0$, write

$$l_j(u|v) = E[Y_t^j I(Y_t \geq u) | X_t = v] = \int_u^{\infty} y^j f(y|v) dy, \quad l_j^{a,b}(u|v) = \frac{\partial^{ab}}{\partial u^a \partial v^b} l_j(u|v),$$

and $l_j^{a,b}(\nu_p(x)|x) = l_j^{a,b}(u|v) \Big|_{u=\nu_p(x), v=x}$. Clearly, $l_0(u|v) = S(u|v)$ and $l_1(\nu_p(x)|x) = p\mu_p(x)$. Finally, $l_j^{1,0}(u|v) = -u^j f(u|v)$ and $l_j^{2,0}(u|v) = -[u^j f^{1,0}(u|v) + j u^{j-1} f(u|v)]$. We now list the following regularity conditions.

Assumption A:

- A1. For fixed y and x , $0 < F(y|x) < 1$, $g(x) > 0$, the marginal density of X_t , and is continuous at x , and $F(y|x)$ has continuous second order derivative with respect to both x and y .
- A2. The kernels $K(\cdot)$ and $W(\cdot)$ are symmetric, bounded, and compactly supported density.
- A3. $h \rightarrow 0$ and $nh \rightarrow \infty$, and $h_0 \rightarrow 0$ and $nh_0 \rightarrow \infty$, as $n \rightarrow \infty$.
- A4. Let $g_{1,t}(\cdot, \cdot)$ be the joint density of X_1 and X_t for $t \geq 2$. Assume that $|g_{1,t}(u, v) - g(u)g(v)| \leq M < \infty$ for all u and v .
- A5. The process $\{(X_t, Y_t)\}$ is a stationary α -mixing with the mixing coefficient satisfying $\alpha(t) = O(t^{-(2+\delta)})$ for some $\delta > 0$.
- A6. $nh^{1+2/\delta} \rightarrow \infty$.
- A7. $h_0 = o(h)$.

Assumption B:

- B1. Assume that $E(|Y_t|^\delta | X_t = u) \leq M_3 < \infty$ for some $\delta > 2$, in a neighborhood of x .
- B2. Assume that $|g_{1,t}(y_1, y_2 | x_1, x_2)| \leq M_1 < \infty$ for all $t \geq 2$, where $g_{1,t}(y_1, y_2 | x_1, x_2)$ be the conditional density of Y_1 and Y_t given $X_1 = x_1$ and $X_t = x_2$.
- B3. The mixing coefficient of the α -mixing process $\{(X_t, Y_t)\}_{t=-\infty}^{\infty}$ satisfies $\sum_{t \geq 1} t^a \alpha^{1-2/\delta}(t) < \infty$ for some $a > 1 - 2/\delta$, where δ is given in Assumption B1.

- B4. Assume that there exists a sequence of integers $s_n > 0$ such that $s_n \rightarrow \infty$, $s_n = o((nh)^{1/2})$, and $(n/h)^{1/2}\alpha(s_n) \rightarrow 0$, as $n \rightarrow \infty$.
- B5. There exists $\delta^* > \delta$ such that $E(|Y_t|^{\delta^*} | X_t = u) \leq M_4 < \infty$ in a neighborhood of x , $\alpha(t) = O(t^{-\theta^*})$, where δ is given in Assumption B1, $\theta^* \geq \delta^* \delta / \{2(\delta^* - \delta)\}$, and $n^{1/2-\delta/4} h^{\delta/\delta^*-1/2-\delta/4} = O(1)$.

Remark 1. Note that Assumptions A1 - A5 and B1 - B5 are used commonly in the literature of time series data (see, e.g., Masry and Fan, 1997, Cai, 2001). Note that α -mixing imposed in Assumption A5 is weaker than β -mixing in Hall, Wolff, and Yao (1999) and ρ -mixing in Fan, Yao, and Tong (1996). Because A6 is satisfied by the bandwidths of optimal size (i.e., $h \approx n^{-1/5}$) if $\delta > 1/2$, we do not concern ourselves with such refinements. Indeed, Assumptions A1 - A6 are also required in Cai (2002). Assumption A7 means that the initial step bandwidth should be chosen as small as possible so that the bias from the initial step can be ignored. Since the common technique – truncation approach for time series data is not applicable to our setting (see, e.g., Masry and Fan, 1997), the purpose of Assumption B5 is to use the moment inequality. If $\alpha(t)$ decays geometrically, then Assumptions B4 and B5 are satisfied automatically. Note that Assumptions B3, B4, and B5 are stronger than Assumptions A5 and A6. This is not surprising because the higher moments involved, the faster decaying rate of $\alpha(\cdot)$ is required. Finally, Assumptions B1 - B5 are also imposed in Cai (2001).

6.4.2 Asymptotic Properties for Conditional PDF and CDF

First, we investigate the asymptotic behaviors of $\hat{f}_c(y | x)$, including the asymptotic normality stated in the following theorem.

Theorem 6.1: *Under Assumptions A1 - A6 with h in A3 and A6 replaced by $h_0 h$, we have*

$$\sqrt{n h_0 h} \left[\hat{f}_c(y | x) - f(y | x) - B_f(y | x) \right] \rightarrow N \{0, \sigma_f^2(y | x)\},$$

where the asymptotic bias is

$$B_f(y | x) = \frac{h^2}{2} \mu_2(W) f^{0,2}(y | x) + \frac{h_0^2}{2} \mu_2(K) f^{2,0}(y | x),$$

and the asymptotic variance is $\sigma_f^2(y | x) = \mu_0(K^2) \mu_0(W^2) f(y | x) / g(x)$.

Remark 2: The asymptotic results for $\widehat{f}_c(y|x)$ in Theorem 6.1 are similar to those for $\widehat{f}_u(y|x)$ in Fan, Yao, and Tong (1996) for the ρ -mixing sequence, which is stronger than α -mixing, but as mentioned earlier, $\widehat{f}_u(y|x)$ is not always a probability density function. The asymptotic bias and variance are intuitively expected. The bias comes from the approximations in both x and y directions and the variance is from the local conditional variance in the density estimation setting, which is $f(y|x)$.

Next, we study the asymptotic behaviors for $\widehat{S}_c(y|x)$ at both interior and boundary points. Similar to Theorem 6.1 for $\widehat{f}_c(y|x)$, we have the following asymptotic normality for $\widehat{S}_c(y|x)$.

Theorem 6.2: *Under Assumptions A1 - A6, we have*

$$\sqrt{n h} \left[\widehat{S}_c(y|x) - S(y|x) - B_S(y|x) \right] \rightarrow N \{0, \sigma_S^2(y|x)\},$$

where the asymptotic bias is given by

$$B_S(y|x) = \frac{h^2}{2} \mu_2(W) S^{0,2}(y|x) - \frac{h_0^2}{2} \mu_2(K) f^{1,0}(y|x),$$

and the asymptotic variance is $\sigma_S^2(y|x) = \mu_0(W^2) S(y|x) [1 - S(y|x)]/g(x)$. In particular, if Assumption A7 holds true, then,

$$\sqrt{n h} \left[\widehat{S}_c(y|x) - S(y|x) - \frac{h^2}{2} \mu_2(W) S^{0,2}(y|x) \right] \rightarrow N \{0, \sigma_S^2(y|x)\}.$$

Remark 3: Note that the asymptotic results for $\widehat{S}_c(y|x)$ in Theorem 6.2 are analogous to those for $\widehat{S}_u(y|x) = 1 - \widehat{F}_u(y|x)$ in Yu and Jones (1998) for iid data, but as mentioned previously, $\widehat{F}_u(y|x)$ is not always a distribution function. A comparison of $B_s(y|x)$ with the asymptotic bias for $\widehat{S}_{cl}(y|x)$ (see Theorem 1 in Cai (2002)), it reveals that there is an extra term $\frac{h_0^2}{2} f^{1,0}(y|x) \mu_2(K)$ in the asymptotic bias expression $B_s(y|x)$ due to the vertical smoothing in the y direction. Also, there is an extra term in the asymptotic variance (see (6.20)). These extra terms are carried over from the initial estimate but they can be ignored if the bandwidth at the initial step is taken to be a higher order than the bandwidth at the smoothing step.

Remark 4: It is important to examine the performance of $\widehat{S}_c(y|x)$ by considering the asymptotic mean squared error (AMSE). Theorem 6.2 concludes that the AMSE of $\widehat{S}_c(y|x)$

is

$$\begin{aligned} \text{AMSE} \left(\widehat{S}_c(y|x) \right) &= \frac{\{h^2 \mu_2(W) S^{0,2}(y|x) - h_0^2 \mu_2(K) f^{1,0}(y|x)\}^2}{4} \\ &+ \frac{1}{n h} \frac{\mu_0(W^2) S(y|x) [1 - S(y|x)]}{g(x)}. \end{aligned} \quad (6.10)$$

By minimizing AMSE in (6.10) and taking $h_0 = o(h)$, therefore, we obtain the optimal bandwidth given by

$$h_{opt,S}(y|x) = \left[\frac{\mu_0(W^2) S(y|x) [1 - S(y|x)]}{\{\mu_2(W) S^{0,2}(y|x)\}^2 g(x)} \right]^{1/5} n^{-1/5}.$$

Therefore, the optimal rate of the AMSE of $\widehat{S}_c(y|x)$ is $n^{-4/5}$.

As for the boundary behavior of the WDKLL estimator, we can follow Cai (2002) to establish a similar result for $\widehat{S}_c(y|x)$ like Theorem 2 in Cai (2002). Without loss of generality, we consider the left boundary point $x = ch$, $0 < c < 1$. From Fan, Hu, and Troung (1994), we take $W(\cdot)$ to have support $[-1, 1]$ and $g(\cdot)$ to have support $[0, 1]$. Then, under Assumptions A1 - A7, by following the same proof as that for Theorem 6.2 and using the second assertion in Lemma 6.1, although not straightforward, we can show that

$$\sqrt{n h} \left[\widehat{S}_c(y|ch) - S_c(y|ch) - B_{S,c}(y) \right] \rightarrow N(0, \sigma_{S,c}^2(y)), \quad (6.11)$$

where the asymptotic bias term is given by $B_{S,c}(y) = h^2 \beta_0(c) S^{0,2}(y|0+)/[2 \beta_1(c)]$ and the asymptotic variance is $\sigma_{S,c}^2(y) = \beta_2(0) S(y|0+)[1 - S(y|0+)]/[\beta_1^2(c) g(0+)]$ with $g(0+) = \lim_{z \downarrow 0} g(z)$,

$$\beta_0(c) = \int_{-1}^c \frac{u^2 W(u)}{1 - \lambda_c u W(u)} du, \quad \beta_j(c) = \int_{-1}^c \frac{W^j(u)}{\{1 - \lambda_c u W(u)\}^j} du, \quad 1 \leq j \leq 2,$$

and λ_c being the root of equation $L_c(\lambda) = 0$

$$L_c(\lambda) = \int_{-1}^c \frac{u W(u)}{1 - \lambda u W(u)} du.$$

Note that the proof of (6.11) is similar to that for Theorem 2 in Cai (2002) and omitted. Theorem 6.2 and (6.11) reflect two of the major advantages of the WKDLL estimator: (a) the asymptotic bias does not depend on the design density $g(x)$, and indeed it is dependent only on the simple conditional distribution curvature $S^{0,2}(y|x)$ and conditional density curvature $f^{1,0}(y|x)$; and (b) it has an automatic good behavior at boundaries. See Cai (2002) for the detailed discussions.

Finally, we remark that if the point 0 were an interior point, then, (6.11) would hold with $c = 1$, which becomes Theorem 6.2. Therefore, Theorem 6.2 shows that the WKDLL estimation has the automatic good behavior at boundaries without the need of the boundary correction.

6.4.3 Asymptotic Theory for CVaR and CES

By the differentiability of $\widehat{S}_c(\widehat{\nu}_p(x) | x)$, we use the Taylor expansion and ignore the higher terms to obtain

$$\widehat{S}_c(\widehat{\nu}_p(x) | x) = p \approx \widehat{S}_c(\nu_p(x) | x) - \widehat{f}_c(\nu_p(x) | x) (\widehat{\nu}_p(x) - \nu_p(x)), \quad (6.12)$$

then, by Theorem 6.1,

$$\widehat{\nu}_p(x) - \nu_p(x) \approx [\widehat{S}_c(\nu_p(x) | x) - p] / \widehat{f}_c(\nu_p(x) | x) \approx [\widehat{S}_c(\nu_p(x) | x) - p] / f(\nu_p(x) | x).$$

As an application of Theorem 6.2, we can establish the following theorem for the asymptotic normality of $\widehat{\nu}_p(x)$ but the proof is omitted since it is similar to that for Theorem 6.2.

Theorem 6.3: *Under Assumptions A1 - A6, we have*

$$\sqrt{n h} [\widehat{\nu}_p(x) - \nu_p(x) - B_\nu(x)] \rightarrow N \{0, \sigma_\nu^2(x)\},$$

where the asymptotic bias is $B_\nu(x) = B_S(\nu_p(x) | x) / f(\nu_p(x) | x)$ and the asymptotic variance is $\sigma_\nu^2(x) = \mu_0(W^2) p(1 - p) / [g(x) f^2(\nu_p(x) | x)]$. In particular, if Assumption A7 holds, then,

$$\sqrt{n h} \left[\widehat{\nu}_p(x) - \nu_p(x) - \frac{h^2}{2} \frac{S^{0,2}(\nu_p(x) | x)}{f(\nu_p(x) | x)} \mu_2(W) \right] \rightarrow N \{0, \sigma_\nu^2(x)\}.$$

Remark 5: First, as a consequence of Theorem 6.3, $\widehat{\nu}_p(x) - \nu_p(x) = O_p(h^2 + h_0^2 + (n h)^{-1/2})$ so that $\widehat{\nu}_p(x)$ is a consistent estimator of $\nu_p(x)$ with a convergence rate. Also, note that the asymptotic results for $\widehat{\nu}_p(x)$ in Theorem 6.3 are akin to those for $\widehat{\nu}_{u,p}(x) = \widehat{S}_u^{-1}(p | x)$ in Yu and Jones (1998) for iid data. But in the bias term of Theorem 6.3, the quantity $S^{0,2}(\nu_p(x) | x) / f(\nu_p(x) | x)$, involving the second derivative of the conditional distribution function with respect to x , replaces $\nu_p''(x)$, the second derivative of the conditional VaR function itself, which is in the bias term of the “check” function type local linear estimator in Yu and Jones (1998) for iid data and Cai and Xu (2005) for time series. See Cai and Xu

(2005) for details. This is not surprising since the bias comes only from the approximation. The former utilizes the approximation of the conditional distribution function but the later uses the approximation of the conditional VaR function. Finally, Theorems 6.2 and 6.3 imply that if the initial bandwidth h_0 is chosen small as possible such as $h_0 = o(h)$, the final estimates of $S(y|x)$ and $\nu_p(x)$ are not sensitive to the choice of h_0 as long as it satisfies Assumption A7. This makes the selection of bandwidths much easier in practice, which will be elaborated later (see Section 6.5.1).

Remark 6: Similar to Remark 5, we can derive the asymptotic mean squared error for $\hat{\nu}_p(x)$. By following Yu and Jones (1998), Theorem 6.3 and (6.20) (given in Section 6.6) imply that the AMSE of $\hat{\nu}_p(x)$ is given by

$$\begin{aligned} \text{AMSE}(\hat{\nu}_p(x)) &= \frac{\{h^2 S^{0,2}(\nu_p(x)|x) \mu_2(W) - h_0^2 f^{1,0}(\nu_p(x)|x) \mu_2(K)\}^2}{4 f^2(\nu_p(x)|x)} \\ &+ \frac{1}{n h} \frac{\mu_0(W^2) [p(1-p) + 2 h_0 f(\nu_p(x)|x) \alpha(K)]}{f^2(\nu_p(x)|x) g(x)}. \end{aligned} \quad (6.13)$$

Note that the above result is similar to that in Theorem 1 in Yu and Jones (1998) for the double kernel local linear conditional quantile estimator. But, a comparison of (6.13) with Theorem 3 in Cai (2002) for the WNW estimator reveals that (6.13) has two extra terms (negligible if Assumption A7 is satisfied) due to the vertical smoothing in the y direction, as mentioned previously. By minimizing AMSE in (6.13) and taking $h_0 = o(h)$, therefore, we obtain the optimal bandwidth given by

$$h_{opt,\nu}(x) = \left[\frac{\mu_0(W^2) p(1-p)}{\{\mu_2(W) S^{0,2}(\nu_p(x)|x)\}^2 g(x)} \right]^{1/5} n^{-1/5}.$$

Therefore, the optimal rate of the AMSE of $\hat{\nu}_p(x)$ is $n^{-4/5}$. By comparing $h_{opt,\nu}(x)$ with $h_{opt,S}(y|x)$, it turns out that $h_{opt,\nu}(x)$ is $h_{opt,\nu}(y|x)$ evaluated at $y = \nu_p(x)$. Therefore, the best choice of the bandwidth for estimating $S_c(y|x)$ can be used for estimating $\nu_p(x)$.

Remark 7: Similar to (6.11), one can establish the asymptotic result at boundaries for $\nu_p(x)$ as follows, one can show that under Assumption A7,

$$\sqrt{n h} [\hat{\nu}_p(c h) - \nu_p(c h) - B_{\nu,c}] \rightarrow N(0, \sigma_{\nu,c}^2),$$

where the asymptotic bias is $B_{\nu,c} = h^2 \beta_2(c) S^{0,2}(\nu_p(0+)|0+)/[2 \beta_1(c) f(\nu_p(0+)|0+)]$ and the asymptotic variance is $\sigma_{\nu,c}^2 = \beta_0(0) p [1-p]/[\beta_1^2(c) f^2(\nu_p(0+)|0+) g(0+)]$. Clearly, $\hat{\nu}_p(x)$

inherits all good properties from the WDKLL estimator of $S_c(y|x)$. Note that the above result can be established by using the second assertion in Lemma 6.1 and following the same lines along with those used in the proof of Theorem 6.2 and omitted.

Finally, we examine the asymptotic behavior for $\hat{\mu}_p(x)$ at both interior and boundary points. First, we establish the following theorem for the asymptotic normality for $\hat{\mu}_p(x)$ when x is an interior point.

Theorem 6.4: *Under Assumptions A1 - A4 and B2 - B5, we have*

$$\sqrt{n h} [\hat{\mu}_p(x) - \mu_p(x) - B_\mu(x)] \rightarrow N \{0, \sigma_\mu^2(x)\},$$

where the asymptotic bias is $B_\mu(x) = B_{\mu,0}(x) + \frac{h_0^2}{2} \mu_2(K) p^{-1} f(\nu_p(x)|x)$ with

$$B_{\mu,0}(x) = \frac{h^2}{2} \mu_2(W) p^{-1} [l_1^{0,2}(\nu_p(x)|x) - \nu_p(x) S^{0,2}(\nu_p(x)|x)],$$

and the asymptotic variance is

$$\sigma_\mu^2(x) = \frac{\mu_0(W^2)}{p g(x)} [p^{-1} l_2(\nu_p(x)|x) - p \mu_p^2(x) + (1-p) \nu_p(x) \{\nu_p(x) - 2 \mu_p(x)\}].$$

In particular, if Assumption A7 holds true, then,

$$\sqrt{n h} [\hat{\mu}_p(x) - \mu_p(x) - B_{\mu,0}(x)] \rightarrow N \{0, \sigma_\mu^2(x)\}.$$

Remark 8: First, Theorem 6.4 concludes that $\hat{\mu}_p(x) - \mu_p(x) = O_p(h^2 + h_0^2 + (n h)^{-1/2})$ so that $\hat{\mu}_p(x)$ is a consistent estimator of $\mu_p(x)$ with a convergence rate. Also, note that the asymptotic results in Theorem 6.4 imply that $\hat{\mu}_p(x)$ is a consistent estimator for $\mu_p(x)$ with a convergence rate $(n h)^{-1/2}$. Further, note that although the asymptotic variance $\sigma_\mu^2(x)$ is the same as that in Scaillet (2005) for $\tilde{\mu}_p(x)$, Scaillet (2005) did not provide an expression for the asymptotic bias term like $B_\mu(x)$ in the first result or $B_{\mu,0}(x)$ in the second conclusion in Theorem 6.4. Clearly, the second term in the asymptotic bias expression is carried over from the y direction smoothing at the initial step and it is negligible if Assumption A7 is satisfied. Clearly, Assumption A7 implies that $B_\mu(x)$ becomes $B_{\mu,0}(x)$.

Remark 9: Like Remark 5, the AMSE for $\hat{\mu}_p(x)$ can be derived in the same manner. It follows from Theorem 6.4 that the AMSE of $\hat{\mu}_p(x)$ is given by

$$\text{AMSE}(\hat{\mu}_p(x)) = \frac{1}{n h} \sigma_\mu^2(x) + \left\{ B_{\mu,0}(x) + \frac{h_0^2}{2} \mu_2(K) p^{-1} f(\nu_p(x)|x) \right\}^2. \quad (6.14)$$

Under Assumption A7, minimizing AMSE in (6.14) with respect to h yields the optimal bandwidth given by

$$h_{opt,\mu}(x) = \left[\frac{\sigma_\mu(x)}{\mu_2(W) p^{-1} \{l_1^{0,2}(\nu_p(x) | x) - \nu_p(x) S^{0,2}(\nu_p(x) | x)\}} \right]^{2/5} n^{-1/5}.$$

Therefore, as expected, the optimal rate of the AMSE of $\hat{\mu}_p(x)$ is $n^{-4/5}$.

Finally, we offer the asymptotic results for $\hat{\mu}_p(x)$ at the left boundary point $x = ch$. By the same fashion, one can show that under Assumption A7,

$$\sqrt{nh} [\hat{\mu}_p(ch) - \mu_p(ch) - B_{\mu,c}] \rightarrow N(0, \sigma_{\mu,c}^2),$$

where the asymptotic bias is

$$B_{\mu,c} = h^2 \beta_2(c) p^{-1} [l_1^{0,2}(\nu_p(0+) | 0+) - \nu_p(0+) S^{0,2}(\nu_p(0+) | 0+)] / [2\beta_1(c)],$$

and the asymptotic variance is

$$\sigma_{\mu,c}^2 = \frac{\beta_0(0)}{p \beta_1^2(c) g(0+)} [p^{-1} l_2(\nu_p(0+) | 0+) - p \mu_p^2(0+) + (1-p) \nu_p(0+) \{\nu_p(0+) - 2\mu_p(0+)\}].$$

Note that the proof of the above result can be carried over by using the second assertion in Lemma 6.1 and following the same lines along with those used in the proof of Theorem 6.4 and omitted. Next, we consider the comparison of the performance of the WDKLL estimation $\hat{\mu}_p(x)$ with the NW type kernel estimator $\tilde{\mu}_p(x)$ as in Scaillet (2005). To this effect, it is not very difficult to derive the asymptotic results for the NW type kernel estimator but the proof is omitted since it is along the same line with the proof of Theorem 6.2. See Scaillet (2005) for the results at the interior point. Under some regularity conditions, it can be shown although tediously (see Cai (2002) for details) that at the left boundary $x = ch$, the asymptotic bias term for the NW type kernel estimator $\tilde{\mu}_p(x)$ is of the order h by comparing to the order h^2 for the WDKLL estimate (see $B_{\mu,c}$ above). This shows that the WDKLL estimate does not suffer from boundary effects but the NW type kernel estimator estimate does. This is another advantage of the WDKLL estimator over the WW type kernel estimator $\tilde{\mu}_p(x)$.

6.5 Empirical Examples

To illustrate the proposed methods, we consider two simulated examples and two real data examples on stock index returns and security returns. Throughout this section, the Epanechnikov kernel $K(u) = 0.75(1 - u^2)_+$ is used and bandwidths are selected as described in the next section.

6.5.1 Bandwidth Selection

With the basic model at hand, one must address the important bandwidth selection issue, as the quality of the curve estimates depends sensitively on the choice of the bandwidth. For practitioners, it is desirable to have a convenient and effective data-driven rule. However, almost nothing has been done so far about this problem in the context of estimating $\nu_p(x)$ and $\mu_p(x)$ although there are some results available in the literature in other contexts for some specific purposes.

As indicated earlier, the choice of the initial bandwidth h_0 is not very sensitive to the final estimation but it needs to be specified. First, we use a very simple idea to choose h_0 . As mentioned previously, the WNW method involves only one bandwidth in estimating the conditional distribution and VaR. Because the WNW estimate is a linear smoother (see (6.5)), we recommend using the optimal bandwidth selector, the so-called nonparametric AIC proposed by Cai and Tiwari (2000), to select the bandwidth, called \tilde{h} . Then we take $0.1 \times \tilde{h}$ or smaller as the initial bandwidth h_0 . For the given h_0 , we can select h as follows. According to (6.8), $\hat{F}_c(\cdot|\cdot)$ is a linear estimator so that the nonparametric AIC selector of Cai and Tiwari (2000) can be applied here to select the optimal bandwidth for $\hat{F}_c(\cdot|\cdot)$, denoted by h_S . As mentioned at the end of Remark 6, the bandwidth for $\hat{\nu}_p(x)$ is the same as that for $\hat{F}_c(\cdot|\cdot)$ so that it is simply to take h_S as h_ν . From (6.9), $\hat{\mu}_p(x)$ is a linear estimator too for given $\hat{\nu}_p(x)$. Therefore, by the same token, the nonparametric AIC selector is applied to selecting h_μ for $\hat{\mu}_p(x)$. This simple approach is used in our implementation in the next sections.

6.5.2 Simulated Examples

In the simulated examples, we demonstrate the finite sample performance of the estimators in terms of the mean absolute deviation error (MADE). For example, the MADE for $\hat{\mu}_p(x)$ is defined as

$$\mathcal{E}_{\mu_p} = \frac{1}{n_0} \sum_{k=1}^{n_0} |\hat{\mu}_p(x_k) - \mu_p(x_k)|,$$

where $\{x_k\}_{k=1}^{n_0}$ are the pre-determined regular grid points. Similarly, we can define the MADE for $\hat{\nu}_p(x)$, denoted by \mathcal{E}_{ν_p} .

Example 6.1. We consider an ARCH type model with $X_t = Y_{t-1}$,

$$Y_t = 0.9 \sin(2.5X_t) + \sigma(X_t)\varepsilon_t,$$

where $\sigma^2(x) = 0.8\sqrt{1.2 + x^2}$ and $\{\varepsilon_t\}$ are iid standard normal random variables. We consider three sample sizes: $n = 250, 500$, and 1000 and the experiment is repeated 500 times for each sample size. The mean absolute deviation errors are computed for each sample size and each replication.

The 5% WDKLL and NW estimations are summarized in Figure 6.1 for CVaR and in Figure 6.2 for CES. For each n , the boxplots of 500 \mathcal{E}_{ν_p} -values of the WDKLL and NW

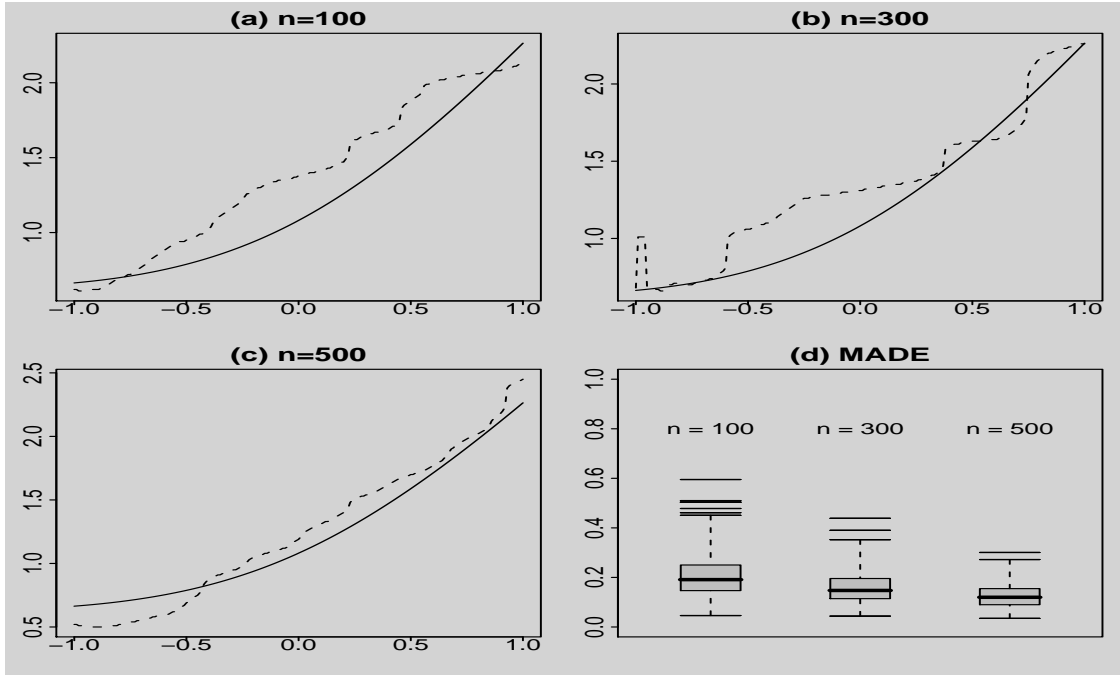


Figure 6.1: Simulation results for Example 1 when $p = 0.05$. Displayed in (a) - (c) are the true CVaR functions (solid lines), the estimated WDKLL CVaR functions (dashed lines), and the estimated NW CVaR functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both the WDKLL and NW estimations of CVaR are plotted in (d).

estimations are plotted in Figure 6.1(d) for CVaR and in Figure 6.2(d) for CES.

From Figures 6.1(d) and 6.2(d), we can observe that the estimation becomes stable as the sample size increases for both the WDKLL and NW estimators. This is in line with our asymptotic theory that the proposed estimators are consistent. Further, it is obvious that the MADEs of the WDKLL estimator are smaller than those for the NW estimator. This

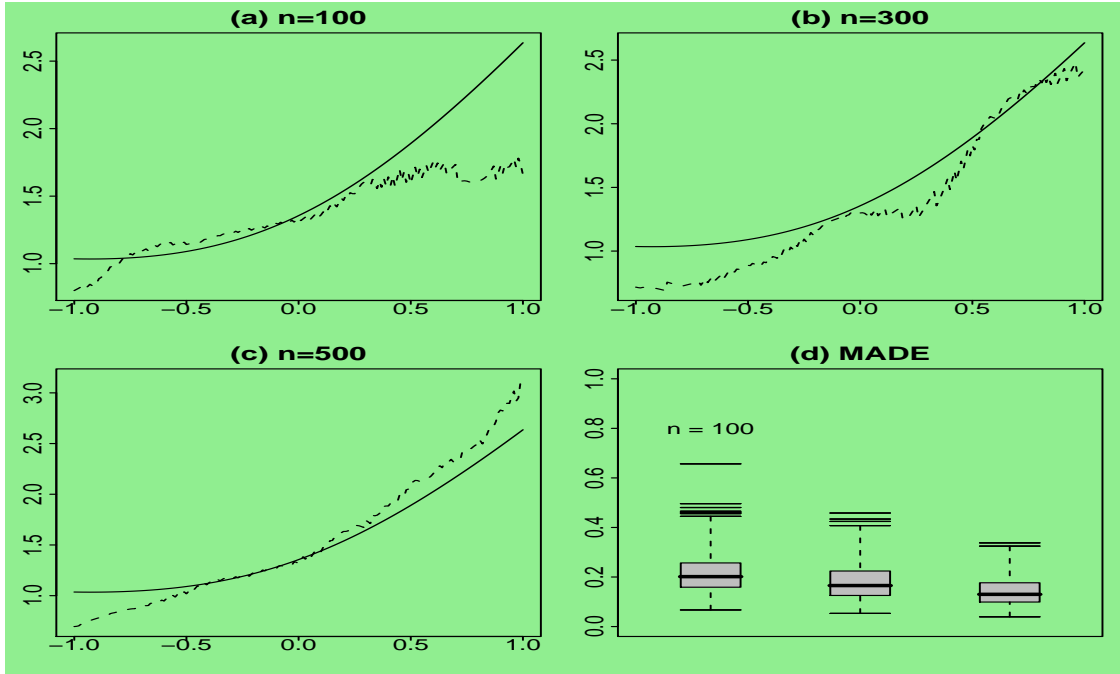


Figure 6.2: Simulation results for Example 1 when $p = 0.05$. Displayed in (a) - (c) are the true CES functions (solid lines), the estimated WDKLL CES functions (dashed lines), and the estimated NW CES functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both the WDKLL and NW estimations of CES are plotted in (d).

indicates that our WDKLL estimator has smaller bias than that for the NW estimator. This implies that the overall performance of the WDKLL estimator should be better than that for the NW estimator.

Figures 6.1(a) – (c) for $n = 250, 500$ and 1000 , respectively, display the true CVaR function (solid line) $\nu_p(x) = 0.9 \sin(2.5x) + \sigma(x)\Phi^{-1}(1 - p)$, where $\Phi(\cdot)$ is the standard normal distribution function, together with the dashed and dotted lines representing the proposed WDKLL (dashed) and NW (dotted) estimates of CVaR, respectively, which are computed based on a typical sample. The typical sample is selected in such a way that its \mathcal{E}_{ν_p} value is equal to the median in the 500 replications. From Figures 6.1(a) – (c), we can observe that both the estimated curves are closer to the true curve as n increases and the performance of the WDKLL estimator is better than that for the NW estimator, especially at boundaries.

In Figures 6.2(a) – (c), the true CES function $\mu_p(x) = 0.9 \sin(2.5x)p + \sigma(x)\mu_1(\Phi^{-1}(1 - p))$ is displayed by the solid line, where $\mu_1(t) = \int_t^\infty u\phi(u)du$ and $\phi(\cdot)$ is the standard normal

distribution density function, and the dashed and dotted lines present the proposed WDKLL (dashed) and NW (dotted) estimates of CES, respectively, from a typical sample. The typical sample is selected in such a way that its \mathcal{E}_{μ_p} -value is equal to the median in the 500

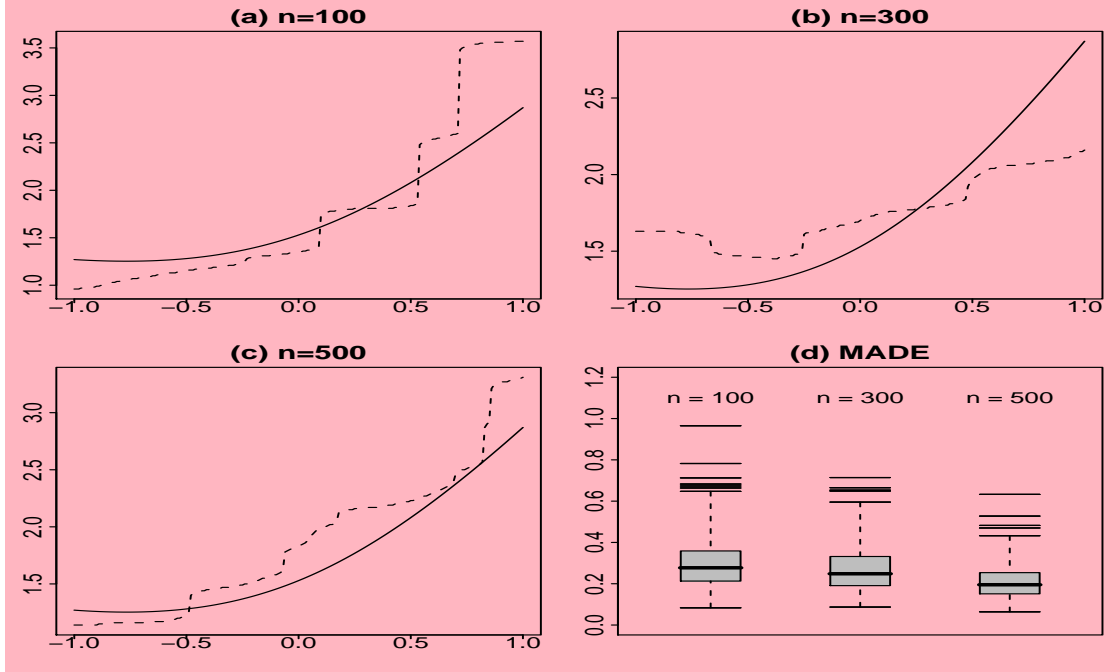


Figure 6.3: Simulation results for Example 1 when $p = 0.01$. Displayed in (a) - (c) are the true CVaR functions (solid lines), the estimated WDKLL CVaR functions (dashed lines), and the estimated NW CVaR functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both WDKLL and NW estimation of the conditional VaR are plotted in (d).

replications. We can conclude from Figures 6.2(a) – (c) that the CES estimator has a similar performance as that for the CVaR estimator.

The 1% WDKLL and NW estimates of CVaR and CES are computed under the same setting and they are displayed in Figures 6.3 and 6.4, respectively. Similar conclusions to those for the 5% estimates can be observed. But it is not surprising to see that the performance of the 1% CVaR and CES estimates is not good as that for the 5% estimates due to the sparsity of data.

Example 6.2. In the above example, we consider only the case when X_t is one-dimensional. In this example, we consider the multivariate situation, i.e. X_t consists of two lagged variables: $X_{t1} = Y_{t-1}$ and $X_{t2} = Y_{t-2}$. The data generating model is given below:

$$Y_t = m(X_t) + \sigma(X_t)\varepsilon_t,$$

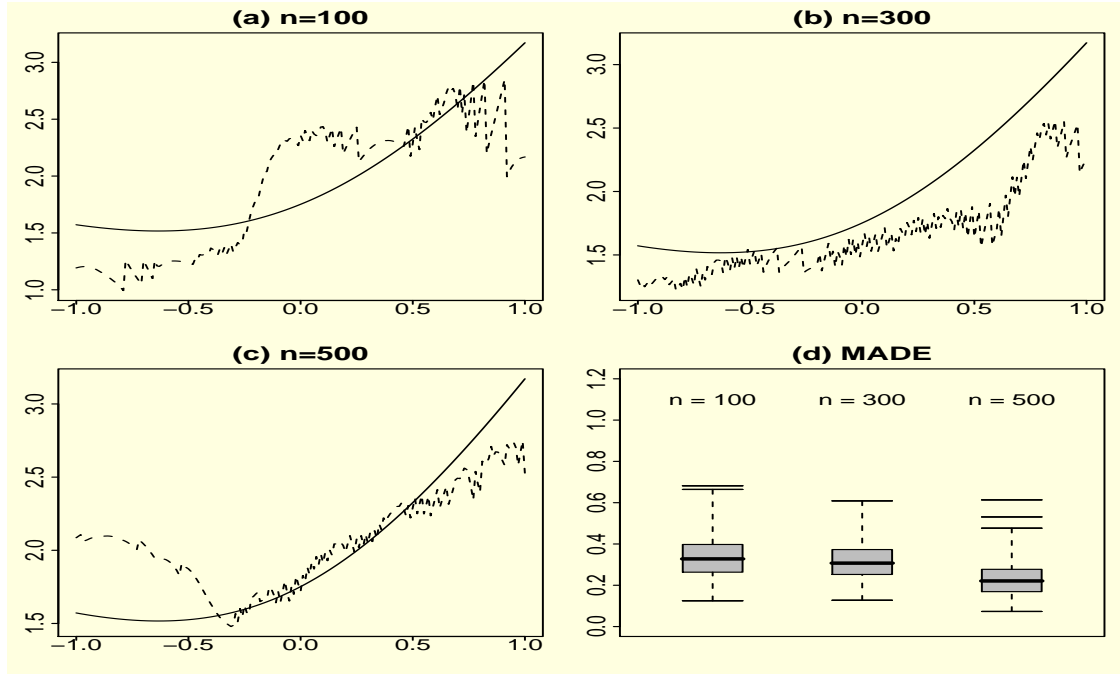


Figure 6.4: Simulation results for Example 1 when $p = 0.01$. Displayed in (a) - (c) are the true CES functions (solid lines), the estimated WDKL CES functions (dashed lines), and the estimated NW CES functions (dotted lines) for $n = 250, 500$ and 1000 , respectively. Boxplots of the 500 MADE values for both the WDKL and NW estimations of CVaR are plotted in (d).

where $m(x) = 0.63x_1 - 0.47x_2$, $\sigma^2(x) = 0.5 + 0.23x_1^2 + 0.3x_2^2$, and $\{\varepsilon_t\}$ are iid generated from $N(0, 1)$. Three sample sizes: $n = 200, 400$, and 600 , are considered here. For each sample size, we replicate the design 500 times. Here we present only the boxplots of the 500 MADEs for the CVaR and CES estimates in Figure 6.5. Figure 6.5(a) displays the boxplots of the 500 \mathcal{E}_{ν_p} -values of the WDKL and NW estimates of CVaR and the boxplots of the 500 \mathcal{E}_{μ_p} -values of the WDKL and NW estimates of CES are given in Figure 6.5(b). From Figures 6.5(a) and (b), it is visually verified that both WDKL and NW estimations become stable as the sample size increases and the performance of the WDKL estimator is better than that for the NW estimator.

6.5.3 Real Examples

Example 6.3. Now we illustrate our proposed methodology by considering a real data set on Dow Jones Industrials (DJI) index returns. We took a sample of 1801 daily prices from DJI index, from November 3, 1998 to January 3, 2006, and computed the daily returns as

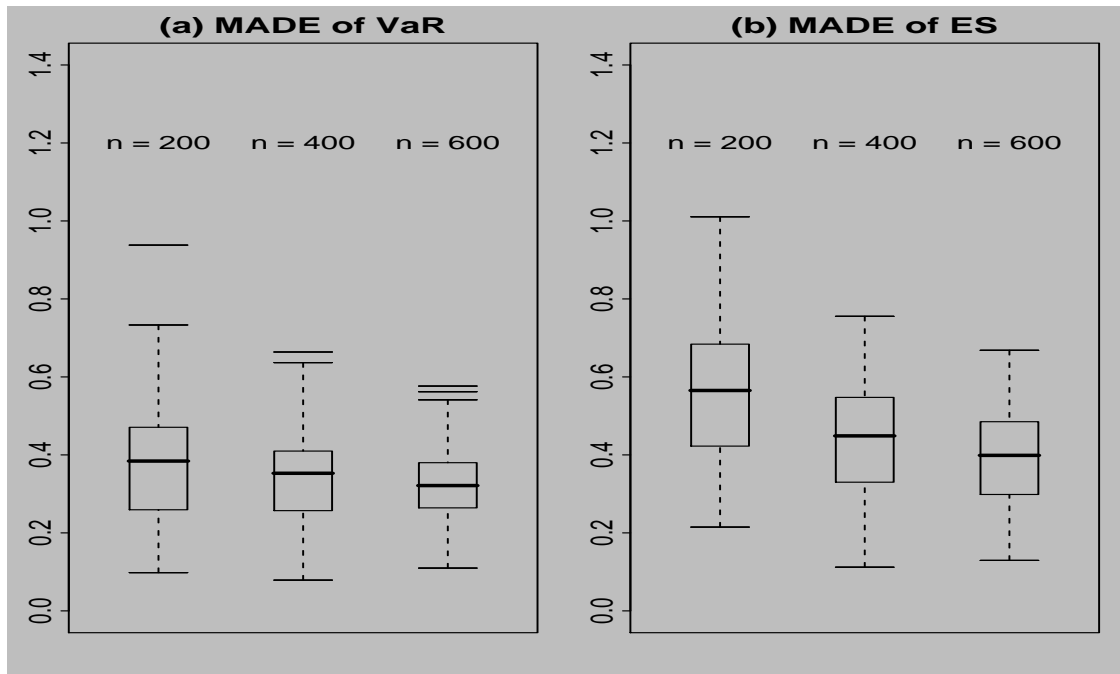


Figure 6.5: Simulation results for Example 2 when $p = 0.05$. (a) Boxplots of MADEs for both the WDKLL and NW CVaR estimates. (b) Boxplots of MADEs for Both the WDKLL and NW CES estimates.

100 times the difference of the log of prices. Let Y_t be the daily negative log return (log loss) of DJI and X_t be the first lagged variable of Y_t . The estimators proposed in this chapter are used to estimate the 5% CVaR and CES functions. The estimation results are shown in Figure 6.6 for the 5% CVaR estimate in Figure 6.6(a) and the 5% CES estimate in Figure 6.6(b). Both CVaR and CES estimates exhibit a U-shape, which corresponds to the so-called “volatility smile”. Therefore, the risk tends to be lower when the lagged log loss of DJI is close to the empirical average and larger otherwise. We can also observe that the curves are asymmetric. This may indicate that the DJI is more likely to fall down if there was a loss within the last day than there was a same amount positive return.

Example 6.4. We apply the proposed methods to estimate the conditional value-at-risk and expected shortfall of the International Business Machine Co. (NYSE: IBM) security returns. The data are daily prices recorded from March 1, 1996 to April 6, 2005. We use the same method to calculate the daily returns as in Example 3. In order to estimate the value-at-risk of a stock return, generally, the information set X_t may contain a market index of corresponding capitalization and type, the industry index, and the lagged values of stock return. For this example, Y_t is the log loss of IBM stock returns and only two variables are

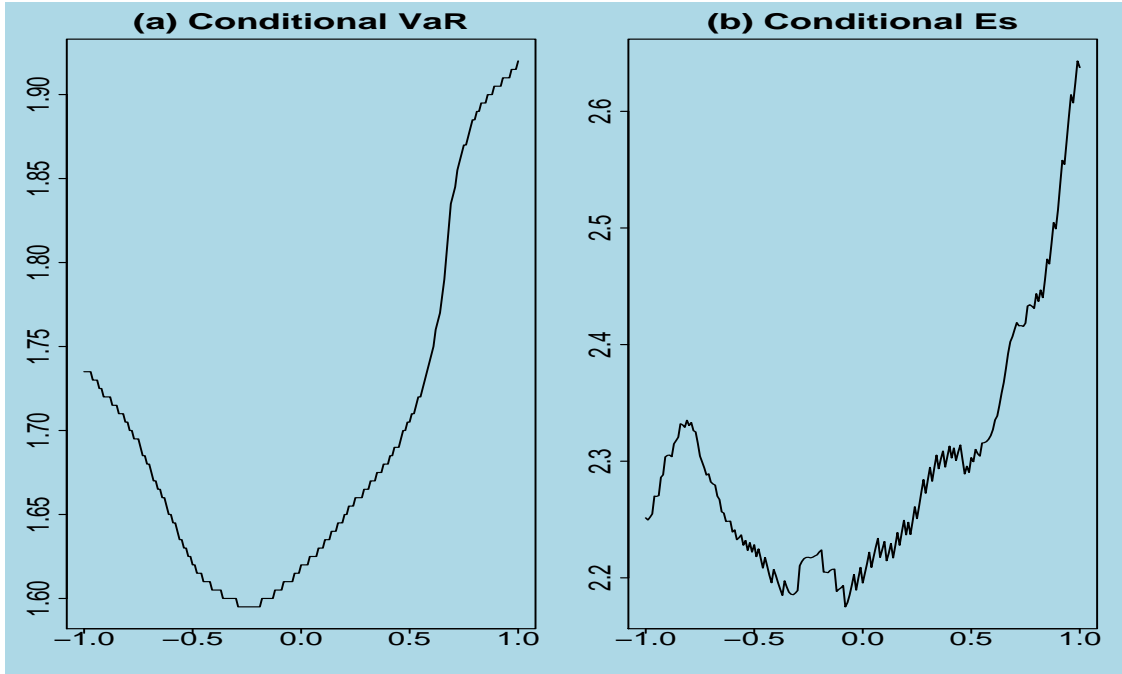


Figure 6.6: (a) 5% CVaR estimate for DJI index. (b) 5% CES estimate for DJI index.

chosen as information set for the sake of simplicity. Let X_{t1} be the first lagged variable of Y_t and X_{t2} denote the first lagged daily log loss of Dow Jones Industrials (DJI) index. Our main results from the estimation of the model are summarized in Figure 6.7. The surfaces of the estimators of IBM returns are given in Figure 6.7(a) for CVaR and in Figure 6.7(b) for CES. For visual convenience, Figures 6.7(c) and (e) depict the estimated CVaR and CES curves (as function of X_{t2}) for three different values of $X_{t1} = (-0.275, -0.025, 0.325)$ and Figures 6.7(d) and (f) display the estimated CVaR and CES curves (as function of X_{t1}) for three different values of $X_{t2} = (-0.225, 0.025, 0.425)$.

From Figures 6.7(c) - (f), we can observe that most of these curves are U-shaped. This is consistent with the results observed in Example 3. Also, we can see that these three curves in each figure are not parallel. This implies that the effects of lagged IBM and lagged DJI variables on the risk of IBM are different and complex. To be concrete, let us examine Figure 6.7(d). Three curves are close to each other when the lagged IBM log loss is around -0.2 and far away otherwise. This implies that DJI has fewer effects (less information) on CVaR around this value. Otherwise, DJI has more effects when the lagged IBM log loss is far from this value.

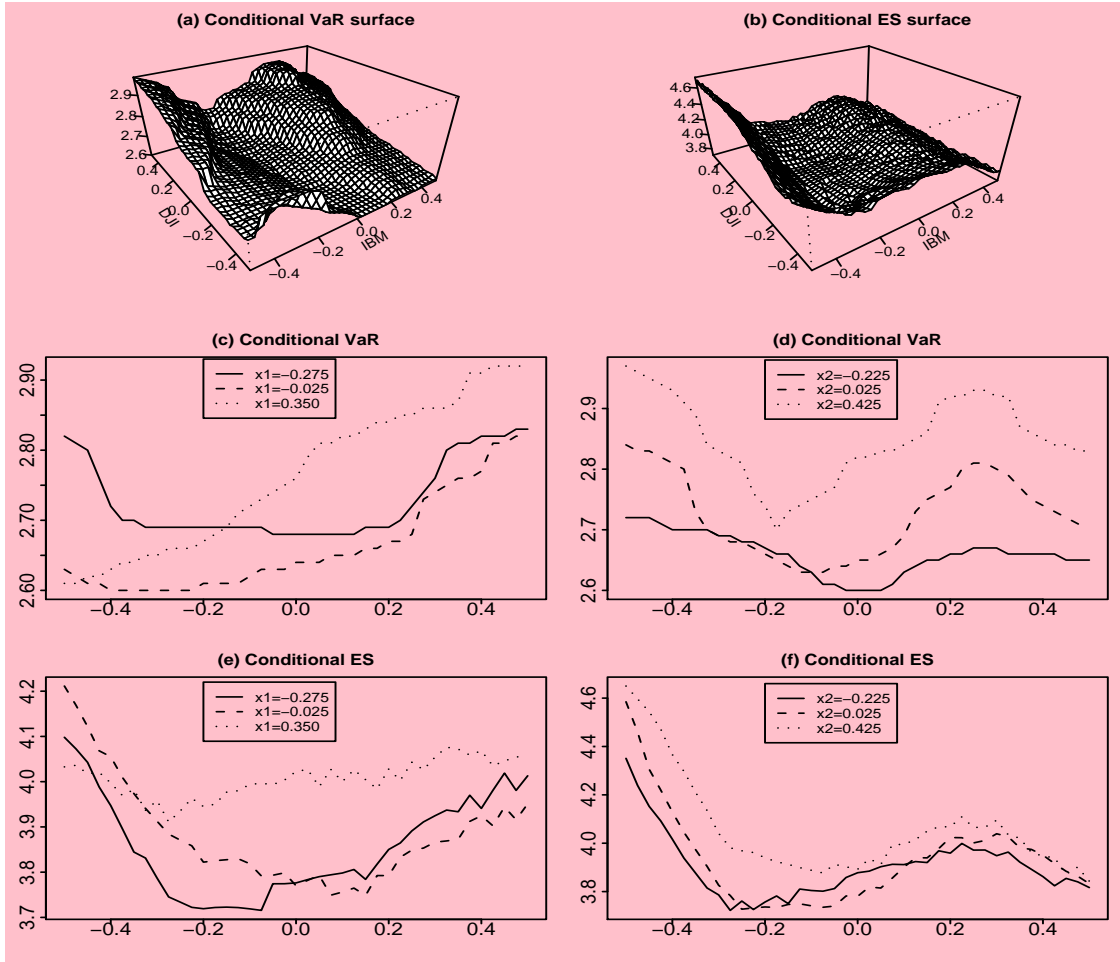


Figure 6.7: (a) 5% CVaR estimates for IBM stock returns. (b) 5% CES estimates for IBM stock returns index. (c) 5% CVaR estimates for three different values of lagged negative IBM returns ($-0.275, -0.025, 0.325$). (d) 5% CVaR estimates for three different values of lagged negative DJI returns ($-0.225, 0.025, 0.425$). (e) 5% CES estimates for three different values of lagged negative IBM returns ($-0.275, -0.025, 0.325$). (f) 5% CES estimates for three different values of lagged negative DJI returns ($-0.225, 0.025, 0.425$).

6.6 Proofs of Theorems

In this section, we present the proofs of Theorems 6.1 - 6.4. First, we list two lemmas. The proof of Lemma 6.1 can be found in Cai (2002) and the proof of Lemma 6.2 is relegated to Section 6.7.

Lemma 6.1: *Under Assumptions A1 - A5, we have*

$$\lambda = -h \lambda_0 \{1 + o_p(1)\} \quad \text{and} \quad p_t(x) = n^{-1} b_t(x) \{1 + o_p(1)\},$$

where $\lambda_0 = \mu_2(W) g'(x) / [2 \mu_2(W^2) g(x)]$ and $b_t(x) = [1 - h \lambda_0 (X_t - x) W_h(x - X_t)]^{-1}$. Fur-

ther, we have

$$p_t(ch) = n^{-1} b_t^c(ch) \{1 + o_p(1)\},$$

where $b_t^c(x) = [1 + \lambda_c(X_t - x) K_h(x - X_t)]^{-1}$.

Lemma 6.2: Under Assumptions A1 - A5, we have, for any $j \geq 0$,

$$J_j = n^{-1} \sum_{t=1}^n c_t(x) \left(\frac{X_t - x}{h} \right)^j = g(x) \mu_j(W) + O_p(h^2),$$

where $c_t(x) = b_t(x) W_h(x - X_t)$.

Before we start to provide the main steps for proofs of theorems. First, it follows from Lemmas 6.1 and 6.2 that

$$W_{c,t}(x, h) \approx \frac{b_t(x) W_h(x - X_t)}{\sum_{t=1}^n b_t(x) W_h(x - X_t)} \approx n^{-1} g^{-1}(x) b_t(x) W_h(x - X_t) = \frac{c_t(x)}{n g(x)}. \quad (6.15)$$

Now we embark on the proofs of the theorems.

Proof of Theorem 6.1: By (6.7), we decompose $\hat{f}_c(y|x) - f(y|x)$ into three parts as follows

$$\hat{f}_c(y|x) - f(y|x) \equiv I_1 + I_2 + I_3, \quad (6.16)$$

where with $\varepsilon_{t,1} = Y_t^*(y) - E(Y_t^*(y)|X_t)$,

$$I_1 = \sum_{t=1}^n \varepsilon_{t,1} W_{c,t}(x, h), \quad I_2 = \sum_{t=1}^n [E(Y_t^*(y)|X_t) - f(y|X_t)] W_{c,t}(x, h),$$

and

$$I_3 = \sum_{t=1}^n [f(y|X_t) - f(y|x)] W_{c,t}(x, h).$$

An application of the Taylor expansion, (6.7), (6.15), and Lemmas 6.1 and 6.2 gives

$$\begin{aligned} I_3 &= \sum_{t=1}^n \frac{1}{2} f^{0,2}(y|x) W_{c,t}(x, h) (X_t - x)^2 + o_p(h^2) \\ &= \frac{1}{2} g^{-1}(x) f^{0,2}(y|x) n^{-1} \sum_{t=1}^n c_t(x) (X_t - x)^2 + o_p(h^2) \\ &= \frac{h^2}{2} \mu_2(W) f^{0,2}(y|x) + o_p(h^2). \end{aligned}$$

By (6.2) and following the same steps as in the proof of Lemma 6.2, we have

$$I_2 = \frac{h_0^2 \mu_2(K)}{2 g(x)} n^{-1} \sum_{t=1}^n f^{2,0}(y|X_t) c_t(x) + o_p(h_0^2 + h^2) = \frac{h_0^2}{2} \mu_2(K) f^{2,0}(y|x) + o_p(h_0^2 + h^2).$$

Therefore,

$$I_2 + I_3 = \frac{h^2}{2} \mu_2(W) f^{0,2}(y|x) + \frac{h_0^2}{2} \mu_2(K) f^{2,0}(y|x) + o_p(h^2 + h_0^2) = B_f(y|x) + o_p(h^2 + h_0^2).$$

Thus, (6.16) becomes

$$\begin{aligned} & \sqrt{nh_0h} \left[\widehat{f}_c(y|x) - f(y|x) - B_f(y|x) + o_p(h^2 + h_0^2) \right] = \sqrt{nh_0h} I_1 \\ & = g^{-1}(x) I_4 \{1 + o_p(1)\} \rightarrow N\{0, \sigma_f^2(y|x)\}, \end{aligned}$$

where $I_4 = \sqrt{h_0h/n} \sum_{t=1}^n \varepsilon_{t,1} c_t(x)$. This, together with Lemma 6.3 in Section 6.7, therefore, proves the theorem. \square

Proof of Theorem 6.2: Similar to (6.16), we have

$$\widehat{S}_c(y|x) - S(y|x) \equiv I_5 + I_6 + I_7, \quad (6.17)$$

where with $\varepsilon_{t,2} = \bar{G}_{h_0}(y - Y_t) - E(\bar{G}_{h_0}(y - Y_t)|X_t)$,

$$I_5 = \sum_{t=1}^n \varepsilon_{t,2} W_{c,t}(x, h), \quad I_6 = \sum_{t=1}^n [E\{\bar{G}_{h_0}(y - Y_t) | X_t\} - S(y|X_t)] W_{c,t}(x, h),$$

and

$$I_7 = \sum_{t=1}^n [S(y|X_t) - S(y|x)] W_{c,t}(x, h).$$

Similar to the analysis of I_2 , by the Taylor expansion, (6.7), and Lemmas 6.1 and 6.2, we have

$$\begin{aligned} I_7 &= \sum_{t=1}^n \frac{1}{2} S^{0,2}(y|x) W_{c,t}(x, h) (X_t - x)^2 + o_p(h^2) \\ &= \frac{1}{2} S^{0,2}(y|x) g^{-1}(x) n^{-1} \sum_{t=1}^n c_t(x) (X_t - x)^2 + o_p(h^2) \\ &= \frac{h^2}{2} \mu_2(W) S^{0,2}(y|x) + o_p(h^2). \end{aligned}$$

To evaluate I_6 , first, we consider the following

$$\begin{aligned} E[\bar{G}_{h_0}(y - Y_t) | X_t = x] &= \int_{-\infty}^{\infty} K(u) S(y - h_0 u | x) du \\ &= S(y|x) + \frac{h_0^2}{2} \mu_2(K) S^{2,0}(y|x) + o(h_0^2) \\ &= S(y|x) - \frac{h_0^2}{2} \mu_2(K) f^{1,0}(y|x) + o(h_0^2). \end{aligned} \quad (6.18)$$

By (6.18) and following the same arguments as in the proof of Lemma 6.2, we have

$$I_6 = -\frac{h_0^2 \mu_2(K)}{2g(x)} n^{-1} \sum_{t=1}^n f^{1,0}(y | X_t) c_t(x) + o_p(h_0^2 + h^2) = -\frac{h_0^2}{2} \mu_2(K) f^{1,0}(y | x) + o_p(h_0^2 + h^2).$$

Therefore,

$$I_6 + I_7 = \frac{h^2}{2} \mu_2(W) S^{0,2}(y | x) - \frac{h_0^2}{2} \mu_2(K) f^{1,0}(y | x) + o_p(h^2 + h_0^2) = B_S(y | x) + o_p(h^2 + h_0^2),$$

so that by (6.17),

$$\sqrt{n h} \left[\widehat{S}_c(y | x) - S(y | x) - B_S(y | x) + o_p(h^2 + h_0^2) \right] = \sqrt{n h} I_5.$$

Clearly, to accomplish the proof of theorem, it suffices to establish the asymptotic normality of $\sqrt{n h} I_5$. To this end, first, we compute $\text{Var}(\varepsilon_{t,2} | X_t = x)$. Note that

$$\begin{aligned} E[\bar{G}_{h_0}^2(y - Y_t) | X_t = x] &= \int_{-\infty}^{\infty} \bar{G}_{h_0}^2(y - u) f(u | x) du \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} K(u_1) K(u_2) S(\max(y - h_0 u_1, y - h_0 u_2) | x) du_1 du_2 \\ &= S(y | x) + 2 h_0 \alpha(K) f(y | x) + O(h_0^2), \end{aligned} \tag{6.19}$$

which, in conjunction with (6.18), implies that

$$\text{Var}(\varepsilon_{t,2} | X_t = x) = S(y | x) [1 - S(y | x)] + 2 h_0 \alpha(K) f(y | x) + o(h_0).$$

This, together with the fact that

$$\text{Var}(\varepsilon_{t,2} c_t(x)) = E [c_t^2(x) E\{\varepsilon_{t,2}^2 | X_t\}] = E [c_t^2(x) \text{Var}(\varepsilon_{t,2} | X_t)],$$

leads to

$$h \text{Var}\{\varepsilon_{t,2} c_t(x)\} = \mu_0(W^2) g(x) [S(y | x)\{1 - S(y | x)\} + 2 h_0 \alpha(K) f(y | x)] + o(h_0).$$

Now, since $|\varepsilon_{t,2}| \leq 1$, by following the same arguments as those used in the proofs of Lemmas 6.2 and 6.3 in Section 6.7 (or Lemma 1 and Theorem 1 in Cai (2002)), we can show although tediously that

$$\text{Var}(I_8) = \sigma_S^2(y | x) g^2(x) + 2 \mu_0(W^2) h_0 \alpha(K) f(y | x) g(x) + o(h_0), \tag{6.20}$$

where $I_8 = \sqrt{h/n} \sum_{t=1}^n \varepsilon_{t,2} c_t(x)$, and

$$\sqrt{n h} I_5 = g^{-1}(x) I_8 \{1 + o_p(1)\} \rightarrow N \{0, \sigma_S^2(y | x)\}.$$

This completes the proof of Theorem 6.2. \square

Proof of Theorem 6.4: Similar to (6.12), we use the Taylor expansion and ignore the higher terms to obtain

$$\begin{aligned} \int_{\hat{\nu}_p(x)}^{\infty} y K_{h_0}(y - Y_t) dy &\approx \int_{\nu_p(x)}^{\infty} y K_{h_0}(y - Y_t) dy - \nu_p(x) K_{h_0}(\nu_p(x) - Y_t) [\hat{\nu}_p(x) - \nu_p(x)] \\ &= Y_t \bar{G}_{h_0}(\nu_p(x) - Y_t) - \nu_p(x) K_{h_0}(\nu_p(x) - Y_t) [\hat{\nu}_p(x) - \nu_p(x)] + h_0 G_{1,h_0}(\nu_p(x) - Y_t). \end{aligned}$$

Plugging the above into (6.9) leads to

$$p \hat{\mu}_p(x) \approx \hat{\mu}_{p,1}(x) + I_9, \quad (6.21)$$

where

$$\hat{\mu}_{p,1}(x) = \sum_{t=1}^n W_{c,t}(x, h) Y_t \bar{G}_{h_0}(\nu_p(x) - Y_t) - \nu_p(x) \hat{f}_c(\nu_p(x) | x) [\hat{\nu}_p(x) - \nu_p(x)],$$

which will be shown later to be the source of both the asymptotic bias and variance, and

$$I_9 = h_0 \sum_{t=1}^n W_{c,t}(x, h) G_{1,h_0}(\nu_p(x) - Y_t),$$

which will be shown to contribute only the asymptotic bias (see Lemma 6.4 in Section 6.7).

From (6.12) and (6.8),

$$\hat{f}_c(\nu_p(x) | x) [\hat{\nu}_p(x) - \nu_p(x)] \approx \sum_{t=1}^n W_{c,t}(x, h) \{ \bar{G}_{h_0}(\nu_p(x) - Y_t) - p \}.$$

Therefore, by (6.15),

$$\begin{aligned} \hat{\mu}_{p,1}(x) &= \sum_{t=1}^n W_{c,t}(x, h) [\{Y_t - \nu_p(x)\} \bar{G}_{h_0}(\nu_p(x) - Y_t) - p \nu_p(x)] \\ &= \sum_{t=1}^n W_{c,t}(x, h) \varepsilon_{t,3} + \sum_{t=1}^n W_{c,t}(x, h) E\{\zeta_t(x) | X_t\} \\ &\approx g^{-1}(x) n^{-1} \sum_{t=1}^n \varepsilon_{t,3} c_t(x) + \sum_{t=1}^n W_{c,t}(x, h) E\{\zeta_t(x) | X_t\} \\ &\equiv \hat{\mu}_{p,2}(x) + \hat{\mu}_{p,3}(x), \end{aligned}$$

where $\zeta_t(x) = [Y_t - \nu_p(x)] \bar{G}_{h_0}(\nu_p(x) - Y_t) + p \nu_p(x)$ and $\varepsilon_{t,3} = \zeta_t(x) - E\{\zeta_t(x) | X_t\}$. Next, we derive the asymptotic bias and variance for $\hat{\mu}_{p,1}(x)$. Indeed, we will show that asymptotic

bias of $\hat{\mu}_p(x)$ comes from both $\hat{\mu}_{p,3}(x)$ and I_9 , and the asymptotic variance for $\hat{\mu}_{p,1}(x)$ is only from $\hat{\mu}_{p,2}(x)$. First, we consider $\hat{\mu}_{p,3}(x)$. Now, it is easy to see by the Taylor expansion that

$$\begin{aligned} E[Y_t \bar{G}_{h_0}(\nu_p(x) - Y_t) | X_t = v] &= \int_{-\infty}^{\infty} K(u) du \int_{\nu_p(x) - h_0 u}^{\infty} y f(y | v) dy \\ &= \int_{-\infty}^{\infty} l_1(\nu_p(x) - h_0 u | v) K(u) du = l_1(\nu_p(x) | v) + \frac{h_0^2}{2} \mu_2(K) l_1^{2,0}(\nu_p(x) | v) + o(h_0^2) \\ &= l_1(\nu_p(x) | v) - \frac{h_0^2}{2} \mu_2(K) [\nu_p(x) f^{1,0}(\nu_p(x) | v) + f(\nu_p(x) | x)] + o(h_0^2), \end{aligned}$$

which, in conjunction with (6.18), leads to

$$\zeta(v) = E[\zeta_t(x) | X_t = v] = A(\nu_p(x) | v) - \frac{h_0^2}{2} \mu_2(K) f(\nu_p(x) | v) + o(h_0^2), \quad (6.22)$$

where $A(\nu_p(x) | v) = l_1(\nu_p(x) | v) - \nu_p(x) [S(\nu_p(x) | v) - p]$. It is easy to verify that $A(\nu_p(x) | v) = E[\{Y_t - \nu_p(x)\} I(Y_t \geq \nu_p(x)) | X_t = v] + p \nu_p(x)$, $A(\nu_p(x) | x) = p \mu_p(x)$, and $A^{0,2}(\nu_p(x) | x) = l_1^{0,2}(\nu_p(x) | x) - \nu_p(x) S^{0,2}(\nu_p(x) | x)$. Therefore, by (6.22), the Taylor expansion, and (6.7), $\hat{\mu}_{p,3}(x)$ becomes

$$\hat{\mu}_{p,3}(x) = \sum_{t=1}^n W_{c,t}(x, h) \zeta(X_t) = \zeta(x) + \frac{1}{2} \zeta''(x) \sum_{t=1}^n W_{c,t}(x, h) (X_t - x)^2 + o_p(h^2).$$

Further, by Lemmas 6.1 and 6.2,

$$\begin{aligned} \hat{\mu}_{p,3}(x) &= \zeta(x) + \frac{h^2}{2} \mu_2(W) \zeta''(x) + o_p(h^2) \\ &= p \mu_p(x) + \frac{h^2}{2} \mu_2(W) A^{0,2}(\nu_p(x) | x) - \frac{h_0^2}{2} \mu_2(K) f(\nu_p(x) | x) + o_p(h_0^2). \end{aligned}$$

This, in conjunction with Lemma 6.4 in Section 6.7, concludes that

$$\hat{\mu}_{p,3}(x) + I_9 = p [\mu_p(x) + B_\mu(x)] + o_p(h^2 + h_0^2),$$

so that by (6.21),

$$\hat{\mu}_{p,1}(x) - p [\mu_p(x) + B_\mu(x)] = \hat{\mu}_{p,2}(x) + o_p(h^2 + h_0^2),$$

and

$$\hat{\mu}_p(x) - \mu_p(x) - B_\mu(x) = p^{-1} \hat{\mu}_{p,2}(x) + o_p(h^2 + h_0^2).$$

Finally, by Lemma 6.5 in Section 6.7, we have

$$\sqrt{nh} [\hat{\mu}_p(x) - \mu_p(x) - B_\mu(x) + o_p(h^2 + h_0^2)] = \frac{1}{pg(x)} I_{10} \{1 + o_p(1)\} \rightarrow N\{0, \sigma_\mu^2(x)\},$$

where $I_{10} = \sqrt{h/n} \sum_{t=1}^n \varepsilon_{t,3} c_t(x)$. Thus, we prove the theorem. \square

6.7 Proofs of Lemmas

In this section, we present the proofs of Lemmas 6.2, 6.3, 6.4, and 6.5. Note that we use the same notation as in Sections 6.2 - 6.6. Also, throughout this section, we denote a generic constant by C , which may take different values at different appearances.

Proof of Lemma 6.2: Let $\xi_t = c_t(x)(X_t - x)^j/h^j$. It is easy to verify by the Taylor expansion that

$$E(J_j) = E(\xi_t) = \int \frac{v^j W(v) g(x - h v)}{1 + h \lambda_0 v W(v)} dv = g(x) \mu_j(W) + O(h^2), \quad (6.23)$$

and

$$E(\xi_t^2) = h^{-1} \int \frac{v^{2j} W^2(v) g(x - h v)}{[1 + h \lambda_0 v W(v)]^2} dv = O(h^{-1}).$$

Also, by the stationarity, a straightforward manipulation yields

$$n \text{Var}(J_j) = \text{Var}(\xi_1) + \sum_{t=2}^n l_{n,t} \text{Cov}(\xi_1, \xi_t), \quad (6.24)$$

where $l_{n,t} = 2(n-t+1)/n$. Now decompose the second term on the right hand side of (6.24) into two terms as follows

$$\sum_{t=2}^n |\text{Cov}(\xi_1, \xi_t)| = \sum_{t=2}^{d_n} (\dots) + \sum_{t=d_n+1}^n (\dots) \equiv J_{j1} + J_{j2}, \quad (6.25)$$

where $d_n = O(h^{-1/(1+\delta/2)})$. For J_{j1} , it follows by Assumption A4 that $|\text{Cov}(\xi_1, \xi_t)| \leq C$, so that $J_{j1} = O(d_n) = o(h^{-1})$. For J_{j2} , Assumption A2 implies that $|(X_t - x)^j W_h(x - X_t)| \leq C h^{j-1}$, so that $|\xi_t| \leq C h^{-1}$. Then, it follows from the Davydov's inequality (see, e.g., Theorem 17.2.1 of Ibragimov and Linnik (1971)) that $|\text{Cov}(\xi_1, \xi_{t+1})| \leq C h^{-2} \alpha(t)$, which, together with Assumption A5, implies that

$$J_{j2} \leq C h^{-2} \sum_{t \geq d_n} \alpha(t) \leq C h^{-2} d_n^{-(1+\delta)} = o(h^{-1}).$$

This, together with (6.24) and (6.25), therefore implies that $\text{Var}(J_j) = O((n h)^{-1}) = o(1)$. This completes the proof of the lemma. \square

Lemma 6.3: Under Assumptions A1 - A6 with h in A3 and A6 replaced by $h h_0$, we have

$$I_4 = \sqrt{\frac{h_0 h}{n}} \sum_{t=1}^n \varepsilon_{t,1} c_t(x) \rightarrow N \{0, \sigma_f^2(y|x) g^2(x)\}.$$

Proof: It follows by using the same lines as those used in the proof of Lemma 6.2 and Theorem 1 in Cai (2002), omitted. The outline is described as follows. First, similar to the proof of Lemma 6.2, it is easy to see that

$$\text{Var}(I_4) = h_0 h \text{Var}(\varepsilon_{t,1} c_t(x)) + h_0 h \sum_{t=2}^n l_{n,t} \text{Cov}(\varepsilon_{1,1} c_1(x), \varepsilon_{t,1} c_t(x)). \quad (6.26)$$

Next, we compute $\text{Var}(\varepsilon_{t,1} | X_t = x)$. Note that

$$h_0 E[Y_t^*(y)^2 | X_t = x] = \int_{-\infty}^{\infty} K^2(u) f(y - h_0 u | x) du = \mu_0(K^2) f(y | x) + O(h_0^2),$$

which, together with the fact that

$$\text{Var}(\varepsilon_{t,1} c_t(x)) = E [c_t^2(x) E\{\varepsilon_{t,1}^2 | X_t\}] = E [c_t^2(x) \text{Var}(\varepsilon_{t,1} | X_t)]$$

and (6.2), implies that

$$h h_0 \text{Var}(\varepsilon_{t,1} c_t(x)) = \mu_0(K^2) \mu_0(W^2) f(y | x) g(x) + O(h_0^2) = \sigma_f^2(y | x) g^2(x) + O(h_0^2).$$

As for the second term on the right hand side of (6.26), similar to (6.25), it is decomposed into two summons. By using Assumption A4 for the first summon and using the Davydov's inequality and Assumption A5 to the second summon, we can show that the second term on the right hand side of (6.26) goes to zero as n goes to infinity. Thus, $\text{Var}(I_4) \rightarrow \sigma_f^2(y | x) g^2(x)$ by (6.26). To show the normality, we employ Doob's small-block and large-block technique (see, e.g., Ibragimov and Linnik, 1971, p. 316). Namely, partition $\{1, \dots, n\}$ into $2q_n + 1$ subsets with large-block of size $r_n = \lfloor (nh h_0)^{1/2} \rfloor$ and small-block of size $s_n = \lfloor (nh h_0)^{1/2} / \log n \rfloor$, where $q_n = \lfloor n / (r_n + s_n) \rfloor$ with $\lfloor x \rfloor$ denoting the integer part of x . By following the same steps as in the proof of Theorem 1 in Cai (2002), we can accomplish the rest of proofs: the summands for the large-blocks are asymptotically independent, two summands for the small-blocks are asymptotically negligible in probability, and the standard Lindeberg-Feller conditions hold for the summands for the large-blocks. See Cai (2002) for details. So, the proof of the lemma is complete. \square

Lemma 6.4: *Under Assumptions A1 - A6, we have*

$$I_9 = h_0 \sum_{t=1}^n W_{c,t}(x, h) G_{1,h_0}(\nu_p(x) - Y_t) = h_0^2 \mu_2(K) f(\nu_p(x) | x) + o_p(h_0^2).$$

Proof: Define $\xi_{t,1} = c_t(x) G_{1,h_0}(\nu_p(x) - Y_t)$. Then, by Lemma 6.1, $I_9 = I_{10} \{1 + o_p(1)\}$, where $I_{10} = g^{-1}(x) h_0 \sum_{t=1}^n \xi_{t,1}/n$. Similar to (6.23),

$$\begin{aligned} E(\xi_{t,1}) &= E[c_t(x) E\{G_{1,h_0}(\nu_p(x) - Y_t) | X_t\}] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{K(u) W(v) u S(\nu_p(x) - h_0 u) | x) g(x - h v)}{1 + h \lambda_0 v W(v)} dudv \\ &= h_0 \mu_2(K) f(\nu_p(x) | x) g(x) + O(h_0 h^2), \end{aligned}$$

and

$$E(\xi_{t,1}^2) = E[b_t^2(x) W_h^2(x - X_t) E\{G_{1,h_0}^2(\nu_p(x) - Y_t) | X_t\}] = O(h_0/h),$$

so that $\text{Var}(\xi_{t,1}) = O(h_0/h)$. By following the same arguments in the derivation of $\text{Var}(J_j)$ in Lemma 6.2, one can show that $\text{Var}(I_{10}) = O((nh)^{-1}) = o(1)$. This proves the lemma. \square

Lemma 6.5: *Under Assumptions A1 - A4 and B2 - B5, we have*

$$I_{10} = \sqrt{\frac{h}{n}} \sum_{t=1}^n \varepsilon_{t,3} c_t(x) \rightarrow N\{0, p^2 g^2(x) \sigma_\mu^2(x)\}.$$

Proof: It follows by using the same lines as those used in the proof of Lemma A.1 and Theorem 1 in Cai (2001), omitted. The main idea is as follows. First, similar to the proof of Lemmas 6.2 and 6.3, we will show by Assumptions B1 - B3 that

$$\text{Var}(I_{10}) \rightarrow p^2 \sigma_\mu^2(x) g^2(x). \quad (6.27)$$

Finally, we need to compute $\text{Var}(\varepsilon_{t,3} c_t(x))$. Since

$$\text{Var}(\varepsilon_{t,3} c_t(x)) = E[c_t^2(x) E\{\varepsilon_{t,3}^2 | X_t\}] = E[c_t^2(x) \text{Var}(\zeta_t(x) | X_t)],$$

then, we first need to calculate $\text{Var}(\zeta_t(x) | X_t)$. To this effect, by (6.22),

$$\begin{aligned} \text{Var}(\zeta_t(x) | X_t = v) &= \text{Var}[(Y_t - \nu_p(x)) \bar{G}_{h_0}(\nu_p(x) - Y_t) | X_t = v] \\ &= E[(Y_t - \nu_p(x))^2 \bar{G}_{h_0}^2(\nu_p(x) - Y_t) | X_t = v] - [l_1(\nu_p(x) | v) - \nu_p(x) S(\nu_p(x) | v)]^2 + O(h_0^2). \end{aligned}$$

Similar to (6.19),

$$\begin{aligned} E[(Y_t - \nu_p(x))^2 \bar{G}_{h_0}^2(\nu_p(x) - Y_t) | X_t = v] &= \int_{-\infty}^{\infty} G_{h_0}^2(\nu_p(x) - y) (y - \nu_p(x))^2 f(y | v) dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} K(u_1) K(u_2) \tau(\max(\nu_p(x) - h_0 u_1, \nu_p(x) - h_0 u_2) | v) du_1 du_2 \\ &= \tau(\nu_p(x) | v) - 2 h_0 \tau^{1,0}(\nu_p(x) | v) \alpha(K) + O(h_0^2) = \tau(\nu_p(x) | v) + O(h_0^2) \end{aligned}$$

since $\tau^{1,0}(\nu_p(x) | v) = 0$, where $\tau(u | v) = l_2(u | v) - 2\nu_p(x)l_1(u | v) + \nu_p^2(x)S(u | v)$. Therefore,

$$\text{Var}(\zeta_t(x) | X_t = v) = \text{Var}[(Y_t - \nu_p(x))I(Y_t \geq \nu_p(x)) | X_t = v] + O(h_0^2),$$

and

$$h \text{Var}(\varepsilon_{t,3} c_t(x)) = \mu_0(W^2) \text{Var}[(Y_t - \nu_p(x))I(Y_t \geq \nu_p(x)) | X_t = x] g(x) + o(1).$$

Similar to Lemmas 6.2 and 6.3, clearly, we have,

$$\text{Var}(I_{10}) = h \text{Var}(\varepsilon_{t,3} c_t(x)) + h \sum_{t=2}^n l_{n,t} \text{Cov}(\varepsilon_{1,3} c_1(x), \varepsilon_{t,3} c_t(x)),$$

and the first term on right hand side of the above equation converges to $p^2 \sigma_\mu^2(x) g^2(x)$. As for the second term on the right hand side of the above equation, similar to (6.25), it is decomposed into two summons. By using Assumptions A4 and B2 for the first summon and using the Davydov's inequality and Assumptions A5 and B3 to the second summon, we can show that the second term on the right hand side of the above equation goes to zero as n goes to infinity. Thus, (6.27) holds. To show the normality, we employ Doob's small-block and large-block technique (see, e.g., Ibragimov and Linnik, 1971, p. 316). Namely, partition $\{1, \dots, n\}$ into $2q_n + 1$ subsets with large-block of size r_n and small-block of size s_n , where s_n is given in Assumption B4, $q_n = \lfloor n/(r_n + s_n) \rfloor$, and $r_n = \lfloor (nh)^{1/2}/\gamma_n \rfloor$ with γ_n satisfying followings: γ_n is a sequence of positive numbers $\gamma_n \rightarrow \infty$ such that $\gamma_n s_n / \sqrt{nh} \rightarrow 0$ and $\gamma_n (n/h)^{1/2} \alpha(s_n) \rightarrow 0$ by Assumption B4. By following the same steps as in the proof of Theorem 1 in Cai (2001) and using Assumption B5, we can accomplish the rest of proofs: the summands for the large-blocks are asymptotically independent, two summands for the small-blocks are asymptotically negligible in probability, and the standard Lindeberg-Feller conditions hold for the summands for the large-blocks. See Cai (2001) for details. Therefore, the lemma is proved. \square

6.8 Computer Codes

Please see the files [chapter6-1.r](#), [chapter6-2.r](#), [chapter6-3.r](#), and [chapter6-4.r](#) for making figures. If you want to learn the codes for computation, they are available upon request.

6.9 References

- Acerbi, C. and D. Tasche (2002). On the coherence of expected shortfall. *Journal of Banking and Finance*, **26**, 1487-1503.

- Artzner, P., F. Delbaen, J.M. Eber, and D. Heath (1999). Coherent measures of risk. *Mathematical Finance*, **9**, 203-228.
- Boente, G. and R. Fraiman (1995). Asymptotic distribution of smoothers based on local means and local medians under dependence. *Journal of Multivariate Analysis*, **54**, 77-90.
- Cai, Z. (2001). Weighted Nadaraya-Watson regression estimation. *Statistics and Probability Letters*, **51**, 307-318.
- Cai, Z. (2002). Regression quantiles for time series data. *Econometric Theory*, **18**, 169-192.
- Cai, Z. (2007). Trending time varying coefficient time series models with serially correlated errors. *Journal of Econometrics*, **137**, 163-188.
- Cai, Z. and L. Qian (2000). Local estimation of a biometric function with covariate effects. In *Asymptotics in Statistics and Probability* (M. Puri, ed), 47-70.
- Cai, Z. and R.C. Tiwari (2000). Application of a local linear autoregressive model to BOD time series. *Environmetrics*, **11**, 341-350.
- Cai, Z. and X. Wang (2006). Nonparametric methods for estimating conditional value-at-risk and expected shortfall. Forthcoming in *Journal of Econometrics*.
- Cai, Z. and X. Xu (2005). Nonparametric quantile estimations for dynamic smooth coefficient models. Forthcoming in *Journal of the American Statistical Association*.
- Chaudhuri, P. (1991). Nonparametric estimates of regression quantiles and their local Bahadur representation. *The Annals of statistics*, **19**, 760-777.
- Chen, S.X. (2006). Nonparametric estimation of expected shortfall. *Working paper*, Department of Statistics, Iowa State University.
- Chen, S.X. and C.Y. Tang (2005). Nonparametric inference of value at risk for dependent financial returns. *Journal of Financial Econometrics*, **3**, 227-255.
- Chernozhukov, V. and L. Umanstev (2001). Conditional value-at-risk: Aspects of modeling and estimation. *Empirical Economics*, **26**, 271-292.
- Cosma, A., O. Scaillet and R. von Sachs (2006). Multivariate wavelet-based shape preserving estimation for dependent observations. *Bernoulli*, in press.
- Duffie, D. and J. Pan (1997). An overview of value at risk. *Journal of Derivatives*, **4**, 7-49.
- Duffie, D. and K.J. Singleton (2003). *Credit Risk: Pricing, Measurement, and Management*. Princeton: Princeton University Press.
- Embrechts, P., C. Kluppelberg, and T. Mikosch (1997). *Modeling Extremal Events For Finance and Insurance*. New York: Springer-Verlag.

- Engle, R.F. and S. Manganelli (2004). CAViaR: conditional autoregressive value at risk by regression quantile. *Journal of Business and Economics Statistics*, **22**, 367-381.
- Fan, J. and I. Gijbels (1996). *Local Polynomial Modeling and Its Applications*. London: Chapman and Hall.
- Fan, J. and J. Gu (2003). Semiparametric estimation of value-at-risk. *Econometrics Journal*, **6**, 261-290.
- Fan, J., T.-C. Hu and Y.K. Troung (1994). Robust nonparametric function estimation. *Scandinavian Journal of Statistics*, **21**, 433-446.
- Fan, J., Q. Yao, and H. Tong (1996). Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika*, **83**, 189-206.
- Frey, R. and A.J. McNeil (2002). VaR and expected shortfall in portfolios of dependent credit risks: conceptual and practical insights. *Journal of Banking and Finance*, **26**, 1317-1334.
- Hall, P., R.C.L. Wolff, and Q. Yao (1999). Methods for estimating a conditional distribution function. *Journal of the American Statistical Association*, **94**, 154-163.
- Hürlimann, W. (2003). A Gaussian exponential approximation to some compound Poisson distributions. *ASTIN Bulletin*, **33**, 41-55.
- Ibragimov, I.A. and Yu. V. Linnik (1971). *Independent and Stationary Sequences of Random Variables*. Groningen, the Netherlands: Walters-Noordhoff.
- Jorion, P. (2001). *Value at Risk*, 2nd Edition. New York: McGraw-Hill.
- Jorion, P. (2003). *Financial Risk Manager Handbook*, 2nd Edition. New York: John Wiley.
- Koenker, R. and G.W. Bassett (1978). Regression quantiles. *Econometrica*, **46**, 33-50.
- Lejeune, M.G. and P. Sarda (1988). Quantile regression: a nonparametric approach. *Computational Statistics and Data Analysis*, **6**, 229-281.
- Masry, E. and J. Fan (1997). Local polynomial estimation of regression functions for mixing processes. *The Scandinavian Journal of Statistics*, **24**, 165-179.
- McNeil, A. (1997). Estimating the tails of loss severity distributions using extreme value theory. *ASTIN Bulletin*, **27**, 117-137.
- Morgan, J.P. (1996). *Risk Metrics - Technical Documents*, 4th Edition, New York.
- Oakes, D. and T. Dasu (1990). A note on residual life. *Biometrika*, **77**, 409-410.
- Rockafellar, R. and S. Uryasev (2000). Optimization of conditional value-at-risk. *Journal of Risk*, **2**, 21-41.

- Roussas, G.G. (1969). Nonparametric estimation of the transition distribution function of a Markov process. *The Annals of Mathematical Statistics*, **40**, 1386-1400.
- Roussas, G.G. (1991). Estimation of transition distribution function and its quantiles in Markov processes: Strong consistency and asymptotic normality. In G.G. Roussas (ed.), *Nonparametric Functional Estimation and related Topics*, pp. 443-462. Amsterdam: Kluwer Academic.
- Samanta, M. (1989). Nonparametric estimation of conditional quantiles. *Statistics and Probability Letters*, **7**, 407-412.
- Scaillet, O. (2004). Nonparametric estimation and sensitivity analysis of expected shortfall. *Mathematical Finance*, **14**, 115-129.
- Scaillet, O. (2005). Nonparametric estimation of conditional expected shortfall. *Revue Assurances et Gestion des Risques/Insurance and Risk Management Journal*, **74**, 639-660.
- Troung, Y.K. (1989). Asymptotic properties of kernel estimators based on local median. *The Annals of Statistics*, **17**, 606-617.
- Troung, Y.K. and C.J. Stone (1992) Nonparametric function estimation involving time series. *The Annals of Statistics* 20, 77-97.
- Yu, K. and M.C. Jones (1997). A comparison of local constant and local linear regression quantile estimation. *Computational Statistics and Data Analysis*, **25**, 159-166.
- Yu, K. and M.C. Jones (1998). Local linear quantile regression. *Journal of the American Statistical Association*, **93**, 228-237.