

## Series 4

1. The dataset `bmw` is a time series of log returns of the BMW stock (business-daily, between June 1986 and March 1990). The log return is defined as follows:

$$X_t = \log \left( \frac{P_t}{P_{t-1}} \right),$$

where  $P_t$  is the stock price at time  $t$ . Log returns can be modelled by

$$X_t = \sigma_t \epsilon_t, \text{ where } \mathbf{E}[\epsilon_t] = 0, \text{Var}(\epsilon_t) = 1, \quad (1)$$

$\epsilon_t$  independent of  $\{X_s; s < t\}$ ,  $\sigma_t^2 = v(X_{t-1})$ , where  $v: \mathbb{R} \mapsto \mathbb{R}^+$  is the so-called “volatility function”. Thus,  $X_t$  depends on  $\{X_s; s < t\}$  only through  $X_{t-1}$  (Markov-property).

The model can be fitted by nonparametric regression of the function  $v$  in

$$Y_t = X_t^2 = v(X_{t-1}) + \eta_t, \text{ where } \eta_t = \sigma_t^2(\epsilon_t^2 - 1)$$

is treated as error term. In task e) you will prove that  $\mathbf{E}[\eta_t] = 0$ .

**Note:** Other usual model assumptions on errors, such as independence, are not fulfilled by  $\eta_t$ , but with some effort (don’t try!) it can be shown that  $v$  can be optimally estimated by the same estimation methods as if the  $\eta_t$  were independent errors.

- a) Model (1) is often chosen for this kind of data because it leads to observations that are not autocorrelated<sup>1</sup>, but dependent. Dependency can be verified by showing that under the model,  $\text{Cov}(X_t^2, X_{t-h}^2) \neq 0$ ,  $h > 0$  (complicated). Plot and interpret the autocorrelation functions of  $X_t$  and  $X_t^2$  for the BMW-dataset.

The data can be read into R by

```
> bmwlr <- scan("http://stat.ethz.ch/Teaching/Datasets/bmw.dat")
```

`bmwlr` should be a vector of 1000 observations.

**R-hint:** Function `acf`. For example, “autocorrelation of lag 1” (in the plot, with 1000 observations, indicated as lag 1 out of 999) means correlation between  $X_t$  and  $X_{t-1}$ . The plot shows also an acceptance region (at 5%-significance level) for testing the null hypothesis of uncorrelated observations.

- b) Fit the data using the nonparametric regression methods Nadaraya-Watson, local polynomial and smoothing splines for the regression function  $v$ .

Comment on the results and compare the fits obtained using the mentioned nonparametric estimators.

**R-hint:** Use `loess` for local polynomial, `smooth.spline` for smoothing splines and `ksmooth` for Nadaraya-Watson kernel regression.

First, sort the `x` values in ascending order (and reorder `y` values according to the new ordering of `x` values), in order not to get problems with `ksmooth`, which orders the values internally and gives back results corresponding to the ordered values.

Start with fitting local polynomial regression to the data (`bmwloess <- loess(...)`). Extract the estimated degree of freedom by `bmwloess$trace.hat` and use this degree of freedom for the `smooth.spline(x,y,df=...)` function. For the Nadaraya-Watson estimator (`ksmooth` function), it is possible to adaptively choose the bandwidth so that the trace of the hat matrix (`=df`) equals the degrees of freedom calculated above (see Series 3). For this purpose, you may use the following code, where ??? should be replaced with the correct information. Use **R** documentation to look up functions `hatMat` (package `sfsmisc`) and `uniroot`. Note that the code takes several minutes to execute.

---

<sup>1</sup>“Autocorrelated” refers to correlation over time, i.e., correlation between  $X_t$  and  $X_{t-h}$ ,  $h > 0$ .

```
# Function delta.dgf returns the degrees of freedom for the Nadaraya-Watson kernel
# regression with bandwidth 'h' minus the degrees of freedom 'dgf' to be matched
delta.dgf <- function (h, x, dgf){
  hatMat(x, trace = TRUE,
    pred.sm = function(x,y,...) ksmooth(sort(x), y, "normal", x.points=x, ...) $y,
    bandwidth = h) - dgf
}
bandwidth.nw <- uniroot(???, c(3,4), x = x, dgf = ???, tol = 0.01)$root
```

Check model assumptions, but don't spend too much time on this since the structure of the data is pretty unclear. Note that for computing residuals it is necessary to know the fitted values at the data points. For `ksmooth` they are provided via argument `x.points` and for `loess` and `smooth.spline` via `fitted()`.

- c) Fit the data using the functions `glkerns` (kernel regression with global optimal bandwidth) and `lokerns` (kernel regression with local optimal bandwidth) of the R package `lokern`. Compare the fits. Plot the local bandwidths from `lokerns` and compare them to the global bandwidth of the function `glkerns`.
- d) (optional) Compute  $\mathbf{E}[X_t | X_{t-1}, X_{t-2}, \dots]$ ,  $\text{Var}(X_t | X_{t-1}, X_{t-2}, \dots)$ ,  $\text{Cov}(X_t, X_{t-h})$ ,  $h > 0$ .

**Background about conditional expectations:**

For two (possibly multi-dimensional) random variables  $X$  and  $Y$ , the conditional distribution  $P_{Y|X=x}$  can be uniquely defined for  $P$ -almost all values of  $X$ . Define  $h(x) = \mathbf{E}[Y | X = x]$  as the expectation of  $Y$  under the conditional distribution  $P_{Y|X=x}$ . For the random variable  $X$ ,  $h(X) = \mathbf{E}[Y | X]$  is a random variable. Here is a very useful equation (the so-called tower property), which will be needed for parts a and b:

$$\mathbf{E}[Y] = \mathbf{E}[\mathbf{E}[Y | X]], \quad (2)$$

the outer expectation taken over the distribution of  $X$ . Conditional variances and covariances are defined analogously.

- e) (optional) Show  $\mathbf{E}[\eta_t] = 0$ .

**Preliminary discussion:** Friday, March 20.

**Deadline:** Friday, March 27.