

Solutions to Problem Set 0  
Economics 551, Yale University  
Woocheol Kim and John Rust

**Question 1** We are asked to compute maximum likelihood estimates of the parameter vector  $(\beta, \sigma^2, \mu, \delta^2)$ , in the model given by:

$$y_i = \exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^4) + \varepsilon_i \quad (1)$$

where we are instructed to believe that  $(y_i, x_i)$  are *i.i.d.* draws from  $x_i \sim N(\mu, \delta^2)$  and  $\varepsilon_i \sim N(0, \sigma^2)$ . Actually you were misled: the errors  $\varepsilon_i$  are heteroscedastic, so conditional on  $x_i$  we have  $\varepsilon_i \sim N(0, \sigma^2(x_i))$  where  $\sigma^2(x) = \exp\{-4 - .2x^2\}$ . So you will be estimating a misspecified model, and later in Econ 551 we will discuss test statistics which are capable of detecting this misspecification. In the meantime your job is to calculate the MLE's of the parameters,  $(\beta, \sigma^2, \mu, \delta^2)$ . The first step is to write down the likelihood function  $L_N(\theta)$  for the data,  $\{(y_i, x_i)\}_{i=1}^N$ . In general “brute force” maximization of  $L_N(\theta)$  may not be a good idea: it might be better to try a “divide and conquer” strategy. Note that the joint density of  $(y, x)$ ,  $f(y, x|\theta)$ , is a product of a conditional likelihood of  $y$  given  $x$ ,  $g(y|x, \beta, \sigma^2)$ , times the marginal density of  $x$ ,  $h(x|\mu, \delta^2)$ . It is easy to see that this factorization or *separability* in the joint likelihood enables us to compute the MLEs for the  $(\beta, \sigma^2)$  parameters and the  $(\mu, \delta^2)$  parameters independently. It also implies a *block diagonality* property which enable us to show that the asymptotic distributions of these parameters are independent.

$$\begin{aligned} L_N(\theta) &= \frac{1}{N} \sum_{i=1}^N \ln f(y_i, x_i | (\beta, \sigma^2, \mu, \delta^2)) \\ &= \frac{1}{N} \sum_{i=1}^N \ln g(y_i | x_i, \beta, \sigma^2) + \frac{1}{N} \sum_{i=1}^N \ln h(x_i | \mu, \delta^2) \\ &= \left\{ -\frac{1}{2} \ln \sigma^2 - \frac{1}{2N\sigma^2} \sum_{i=1}^n [y_i - \exp(X_i \beta)]^2 \right\} \\ &\quad + \left\{ -\frac{1}{2} \ln \delta^2 - \frac{1}{2N\delta^2} \sum_{i=1}^n [x_i - \mu]^2 \right\} - \frac{1}{2} \ln(2\pi) \\ &= L_N^c(\beta, \sigma^2) + L_N^m(\mu, \delta^2), \end{aligned}$$

where  $L_N^c$  denotes the conditional likelihood of the  $\{y_i\}$ 's given the  $\{x_i\}$ 's,  $L_N^m$  denotes the marginal likelihood of the  $\{x_i\}$ 's,  $X_i = (1, x_i, x_i^2, x_i^3)'$  and  $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)'$ . The separability of parameters in the third and fourth equation allows us to break the estimation problem into two subproblems, which ordinarily makes the programming considerably easier and computations considerably faster:

- (1) calculate the MLE's for  $(\beta, \sigma^2)$  from the conditional likelihood  $L_N^c(\beta, \sigma^2)$ :

$$\max_{(\beta, \sigma^2)} L_N^c(\beta, \sigma^2) = -\frac{1}{2} \ln \sigma^2 - \frac{1}{N} \sum_{i=1}^N \frac{[y_i - \exp(X_i \beta)]^2}{2\sigma^2}, \quad (2)$$

with FOC's

$$\begin{aligned} \frac{\partial L_N^c}{\partial \beta}(\beta, \sigma^2) &= \frac{1}{N\sigma^2} \sum_{i=1}^N [y_i - \exp(X_i \beta)] \exp(X_i \beta) X_i' = 0 \\ \frac{\partial L_N^c}{\partial \sigma^2}(\beta, \sigma^2) &= \frac{-1}{2\sigma^2} + \frac{1}{2N\sigma^4} \sum_{i=1}^N [y_i - \exp(X_i \beta)]^2 = 0, \end{aligned}$$

Note that there is further separability in this first subproblem: the FOC for  $\beta$  is the same as the FOC for nonlinear least squares (NLS) estimation of  $\beta$  in equation (2) above, ignoring  $\sigma^2$  since it doesn't affect the solution for  $\hat{\beta}$ . Once we have computed the NLS estimate of  $\hat{\beta}$ , we use the second equation to compute  $\hat{\sigma}^2$  as the sample variance of the estimated NLS residuals.

- (2) calculate the MLE's for  $(\mu, \delta^2)$  from the marginal likelihood

$$\max_{(\mu, \delta^2)} L_N^m(\mu, \delta^2) = \frac{1}{2} \ln \delta^2 - \frac{1}{2N\delta^2} \sum_{i=1}^N [x_i - \mu]^2, \quad (3)$$

with FOC's:

$$\begin{aligned} \frac{\partial}{\partial \mu} L_N^m(\mu, \delta^2) &= \frac{1}{N\delta^2} \sum_{i=1}^n [x_i - \mu] = 0 \\ \frac{\partial}{\partial \delta} L_N^m(\mu, \delta^2) &= -\frac{1}{2\delta^2} + \frac{1}{2N\delta^4} \sum_{i=1}^n [x_i - \mu]^2 = 0 \end{aligned}$$

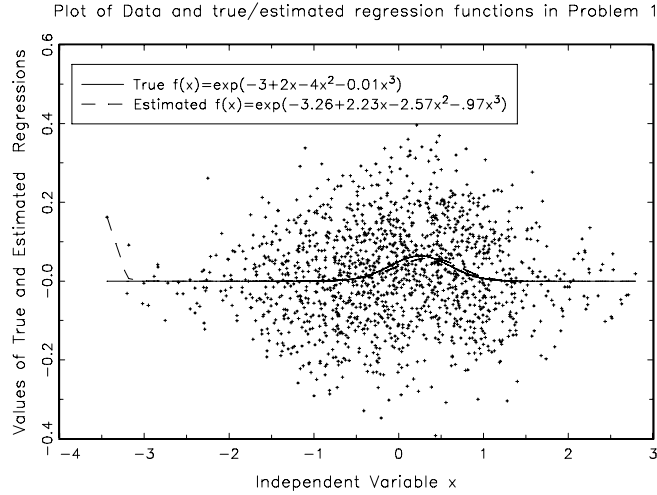
respectively. Using attached Gauss code `nlreg.gpr` for computing NLS estimates (the sum of squared errors, derivatives and hessian are coded in the `eval.g` procedure) we are able to numerically solve for a vector  $\hat{\beta}$  that sets the FOC for  $\beta$  given in equation (2) above to zero. There are closed-form expressions for the MLEs for the remaining parameters:

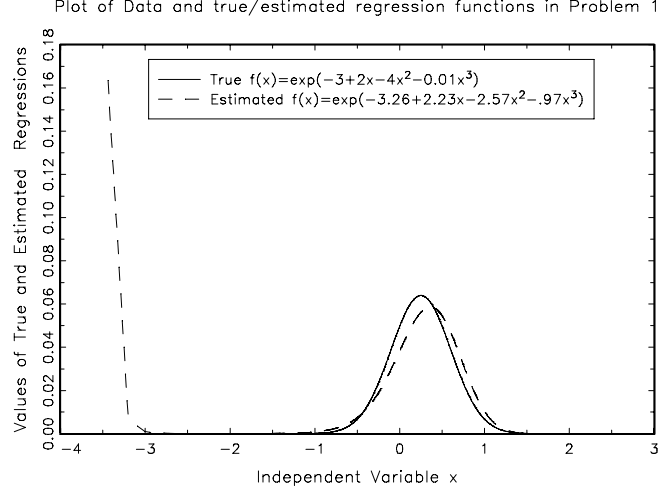
$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{N} \sum_{i=1}^N \left[ y_i - \exp(X_i \hat{\beta}) \right]^2, \\ \hat{\mu} &= \frac{1}{N} \sum_{i=1}^N x_i \\ \hat{\delta}^2 &= \frac{1}{N} \sum_{i=1}^N [x_i - \hat{\mu}]^2. \end{aligned}$$

**Table 1:** Maximum Likelihood Estimates of  $\theta$  using `data1.asc`

Parameter	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\sigma^2$	$\mu$	$\delta^2$
True Value	-3.000	2.000	-4.000	-0.010	0.015	0.000	1.000
MLE	-3.259	2.228	-2.571	-0.971	0.015	-0.047	0.974
Std. Dev.	0.194	0.809	0.836	0.304	.00055	0.025	0.036

Figures 1 and 2 below plot the true and estimated regression functions for this problem. Figure 1 plots the data points also: we see that both the true and estimated regression functions are generally quite close to each other and both go through the middle of the “data cloud”. However we see that the MLE gives a substantially downward biased estimate of  $\beta_3^*$  and this causes the estimated regression function to make big divergences from the true regression function at extreme high and low values of  $x$ , say  $|x| > 4$ . However since there are very few high or low values of  $x$  in the sample, the NLS and MLE are not able to “penalize” this divergence: the MLE sets  $\hat{\beta}_3 = -.97$  since it helps fit the data around  $x = 0$  where most of the data points are. Figure 2 provides a blow up of Figure 1 without the data points to show you how the estimated regression function diverges from the truth near  $x = 0$ . Overall, despite the misspecification of the heteroscedasticity, the MLE and NLS estimators seem to do a pretty good job of uncovering the true regression function, at least for those  $x$ ’s where we have sufficient data. Note that the data plot indicates heteroscedasticity, since the variance of the data around the regression function is bigger in the middle of the graph (near  $x = 0$ ) than at large positive or negative values of  $x$ .

**Figure 1:** Plot of `data1.asc` and true and estimated regression functions



**Figure 2:** Blow up of true and estimated regression functions

Since the true model used to generate `data1.asc` has heteroscedastic and not homoscedastic error terms  $\{\epsilon_i\}$  as assumed here, it is easy to show that the MLE  $\hat{\sigma}^2$  of the misspecified model converges to the expectation of  $\sigma^2(x) = \exp\{-4 - .2x^2\}$ . We can calculate  $E\{\sigma^2(\tilde{x})\}$  as follows:

$$\begin{aligned}
 E\{\sigma^2(x)\} &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \exp\{-4 - .2x^2\} \exp\{-\frac{x^2}{2}\} dx \\
 &= \frac{\exp\{-4\}}{\sqrt{1.4}} \frac{\sqrt{1.4}}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \exp\{-\frac{\sqrt{1.4}x^2}{2}\} dx \\
 &= \frac{\exp\{-4\}}{\sqrt{1.4}} = .015479540.
 \end{aligned}$$

We leave it to you to show that even though the model is misspecified, the MLE  $\hat{\sigma}^2$  converges to  $E\{\sigma^2(\tilde{x})\}$  as  $N \rightarrow \infty$ . Indeed in this case we find that  $\hat{\sigma}^2 = 0.015$ , which happens to be almost exactly equal to the “true” value.

Recall that the asymptotic distribution of the standardized MLE estimator is given by:

$$\sqrt{N}(\hat{\theta} - \theta^*) \xRightarrow{d} N\left(0, H(\theta^*)^{-1} \mathcal{I}(\theta^*) H(\theta^*)^{-1}\right), \quad (4)$$

where  $H(\theta^*)$  is the Hessian and  $\mathcal{I}(\theta^*)$  is the information matrix (both evaluated at  $\theta^*$ ). However the likelihood is misspecified in this case (due to heteroscedasticity), and it is easy to verify that the equality of  $\mathcal{I}(\theta^*) = -H(\theta^*)$  does not hold, so the correct asymptotic covariance matrix is given by the White “misspecification consistent” formula in equation (4) rather than by the inverse of the information matrix  $\mathcal{I}(\theta^*)$ . For example, consider the  $(\beta, \beta)$  block of  $\mathcal{I}(\theta^*)$ , or  $\mathcal{I}_{\beta\beta}(\theta^*)$ :

$$\mathcal{I}_{\beta\beta}(\theta^*) = E\left\{\frac{\sigma^2(x) \exp\{X\beta^*\} XX'}{\sigma^4}\right\} \neq -E\left\{-\frac{\exp\{2X\beta^*\} XX'}{\sigma^2}\right\} = -H_{\beta\beta}(\theta^*)$$

We see that a sufficient condition for the two expressions above to equal each other is  $\sigma^2(x) = \sigma^2$  for all  $x$ , i.e. for the model to be homoscedastic as you were asked to assume. The failure of this equality can be a basis for a specification test statistic that can detect model misspecification which we will discuss in more detail below. Note that even despite the misspecification, the separability property implies that the Hessian is a block diagonal matrix:

$$E \begin{bmatrix} \frac{\partial^2 \ln g(y|x; (\beta, \sigma^2))}{\partial(\beta, \sigma^2) \partial(\beta, \sigma^2)'} & 0 \\ 0 & \frac{\partial^2 \ln h(x; (\mu, \delta^2))}{\partial(\mu, \delta^2) \partial(\mu, \delta^2)'} \end{bmatrix} \\ = \begin{bmatrix} \frac{1}{\sigma^2} E(Z_i X_i) & 0 & 0 & 0 \\ 0 & -\frac{1}{2\sigma^4} & 0 & 0 \\ 0 & 0 & -\frac{1}{\delta^2} & 0 \\ 0 & 0 & 0 & -\frac{1}{2\delta^4} \end{bmatrix},$$

where  $Z_i \equiv [y_i - 2 \exp(X_i \beta)] \exp(X_i \beta) X_i'$ . It is also easy to verify that despite the misspecification of the heteroscedasticity, the information matrix  $\mathcal{I}(\theta^*)$  is still block diagonal. Block diagonality of  $\mathcal{I}(\theta^*)$  and  $H(\theta^*)$  implies that the covariance matrix of  $\sqrt{N}(\hat{\theta} - \theta^*)$  is block diagonal. One can see further block diagonality between  $\beta$  and  $\sigma^2$  and between  $\mu$  and  $\delta^2$ . This block diagonality is not just a consequence of the separability in both the marginal and conditional likelihood in the parameters describing the mean or conditional mean ( $\mu$  and  $\beta$ , respectively) and the variance parameters ( $\delta^2$  and  $\sigma^2$ , respectively). You should verify through direct calculation that this block diagonality is a result of the symmetry of the normal distribution, which implies that  $E\{\epsilon^3\} = 0$ , where the conditional distribution of  $\epsilon = (y - \exp\{X\beta\})$  given  $X$  is  $N(0, \sigma^2(x))$ , and similarly, the distribution of  $\eta = (x - \mu^*) \sim N(0, \delta^2)$ , which implies that  $E\{\eta^k\} = 0$  for any positive odd integer  $k$ . If the normal distribution were not symmetric the block diagonality property wouldn't hold.

Using the block diagonality property, it is easy to compute the standard errors of the full parameter vector,  $\hat{\theta}$ . The covariance matrix of  $\hat{\beta}$  is given by  $1/N$  times the upper  $(\beta, \beta)$  block of  $H(\theta^*)^{-1} \mathcal{I}(\theta^*) H(\theta^*)^{-1}$  and it is easy to verify that this is the same as the covariance matrix for the nonlinear least squares estimator for  $\beta^*$  (which is also numerically identical to the MLE) that is output from the `nlreg.gpr` program. By working with equation (4), you can show that the estimated variance of  $\sigma^2$  is given by  $\text{var}(\hat{\sigma}^2) = [\hat{\mu}_4 - \hat{\sigma}^4]/N$  where  $\hat{\mu}_4$  is the sample analog of the fourth central moment of  $\tilde{X}$ :

$$\hat{\mu}_r = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^r.$$

There is a similar formula for the estimated variance of  $\delta^2$ . We have  $\text{var}(\hat{\delta}^2) = 1.948/1500$ , so the estimated standard error of  $\hat{\delta}^2$  is 0.036. The estimated standard error of  $\hat{\sigma}^2$  is  $\text{std}(\hat{\sigma}^2) = 0.0005499$ .

**Question 2.** If the model is correctly specified, we know that the information equality will hold which implies that:

$$H(\theta^*)^{-1} \mathcal{I}(\theta^*) H(\theta^*)^{-1} = \mathcal{I}^{-1}(\theta^*). \quad (5)$$

We compare White's misspecification consistent estimate to the inverse of information:

$$\hat{H}(\hat{\theta}_{ML})^{-1} \hat{\mathcal{I}}(\hat{\theta}_{ML}) \hat{H}(\hat{\theta}_{ML})^{-1} = \begin{bmatrix} 0.0374 & -0.0872 & 0.0255 & 0.0157 \\ -0.0872 & 0.6543 & -0.5643 & -0.2215 \\ 0.0255 & -0.5643 & 0.6998 & 0.2518 \\ 0.0157 & -0.2215 & 0.2518 & 0.0923 \end{bmatrix}$$

$$\hat{\mathcal{I}}^{-1}(\hat{\theta}_{ML}) = \begin{bmatrix} 114.9 & -254.7 & 9.0 & 95.8 \\ -254.7 & 2155.5 & -2045.8 & -810.5 \\ 9.0 & -2045.8 & 3765.1 & 6.2 \\ 95.8 & -810.5 & 6.2 & 1721.5 \end{bmatrix},$$

where all the estimates are given from the results of eval.g. Following White ("Maximum Likelihood Estimation of Misspecified Models" *Econometrica* 1982) we can construct a formal hypothesis test statistic using the difference between the estimates of the upper diagonal elements of  $H(\theta^*)$  and the corresponding elements of  $\mathcal{I}(\theta^*)$ . This statistic should be small if the null hypothesis of correct specification is true (since  $-H(\theta^*) = \mathcal{I}(\theta^*)$  in that case), and large if the model is misspecified (since it is not necessarily true that  $-H(\theta^*) = \mathcal{I}(\theta^*)$  if the model is misspecified). The large difference in the two difference estimates of the covariance matrix for  $\hat{\beta}$  suggests that  $-H(\theta^*)$  and  $\mathcal{I}(\theta^*)$  are different, and hence that the model is misspecified. However we did not actually compute the actual test statistic to see at what level of significance null would actually be rejected (i.e. to compute the marginal significance level of the test statistic) since we didn't expect you to know about this particular specification test statistic at this stage of the course.

**Question 3.** If the true model were *log-linear*, i.e. if

$$y_i = \exp\{X_i\beta\} \exp\{\epsilon_i\}$$

where  $\epsilon_i \sim N(0, \sigma^2)$ , then it is easy to see that the  $\{y_i\}$  are conditionally lognormally distributed so it is valid to take log transformation and estimate  $(\beta, \sigma^2)$  by OLS:

$$\ln(y_i) = X_i\beta + \epsilon_i$$

. It would not difficult to show that the OLS estimates of the log-linearized model are actually the maximum likelihood estimates of the original lognormal specification (make sure you understand this by writing down the lognormal likelihood function and verifying what we just said is true)! However the error term for the specification of Model I in question 1 is additive and not multiplicative, and so the log transformation is generally not valid. Indeed, there is

a positive probability of observing negative realizations of  $y_i$ , something that has zero probability under the lognormal specification. Thus, to even do OLS one must screen out all  $(y_i, x_i)$  pairs where  $y_i$  is negative, something that is generally not a good idea. When one does OLS on this nonrandomly selected subsample, it is not surprising that the estimates are highly biased:

**Table 2:** Comparison of MLE and OLS estimates of  $\beta$

Parameter	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$
True values	-3.000	2.000	-4.000	-0.010
Model I (MLE)	-3.259	2.228	-2.570	-0.971
Model II (OLS)	-2.543	0.118	-0.152	-0.033

We can show analytically why the OLS estimates for Model II will be inconsistent when the true model is Model I with additive normal errors rather than multiplicative lognormal errors. After screening out negative  $y$ 's, the OLS estimator solves:

$$\hat{\beta}_{OLS} = \underset{\beta \in R^4}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N [\ln(y_i) - X_i \beta]^2 I\{y_i > 0\}. \quad (6)$$

As  $N \rightarrow \infty$ , we can show that the right hand side of the above equation converges uniformly to

$$E \left\{ [\ln(y) - X\beta]^2 I\{y > 0\} \right\}.$$

In general the  $\beta$  that minimizes the expression above is not the same as the  $\beta$  that minimizes

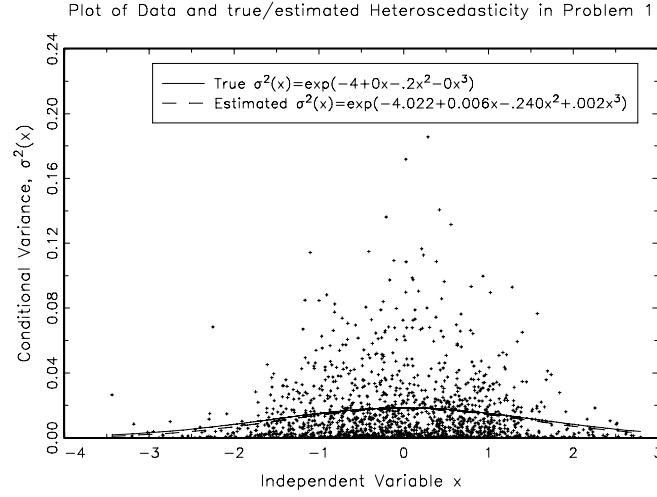
$$E\{[y - \exp(X\beta)]^2\},$$

which is the true  $\beta^*$  when the conditional mean function  $\exp\{X\beta\}$  is correctly specified. Therefore, we conclude that the probability limits for  $\hat{\beta}$ 's from the two different models (Model I and Model II) are not the same.

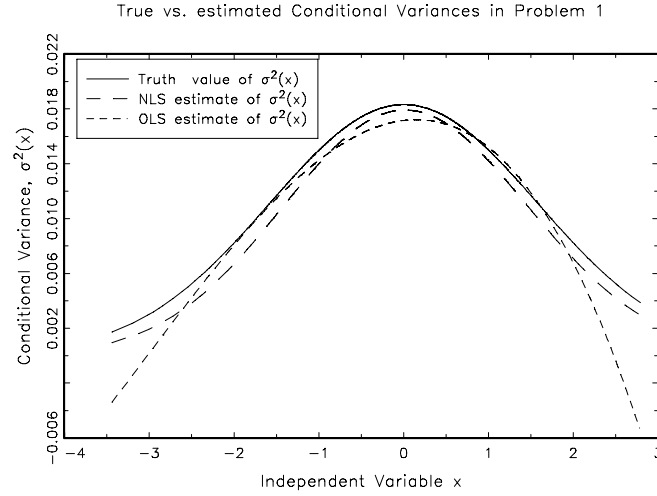
**Question 4** Figure 3 below plots the squared residuals  $\hat{\epsilon}_i^2 = (y_i - \exp\{X_i \hat{\beta}\})^2$  from the MLE/NLS estimation results in Question 1. The figure also plots the results of the following nonlinear regression:

$$\hat{\epsilon}_i^2 = \exp\{X_i \gamma\} + \eta_i, \quad (7)$$

where as before  $X_i = (1, x_i, x_i^2, x_i^3)$  and  $\gamma$  is a conformable  $4 \times 1$  parameter vector. We see substantial evidence of heteroscedasticity, confirming our earlier visual impression in looking at the data in figure 1. The estimated conditional variance function looks concave and symmetric w.r.t.  $y$ -axis, almost like a normal density. Table 3 summarizes the estimation results. According to table 3, the regression coefficient for the quadratic term is significant and seems to dominate the form of the heteroscedasticity plotted in figure 3. Figure 4 does a blow up, plotting the estimated and true conditional variance functions.



**Figure 3:** Estimated Squared Residuals and True vs. Estimated  $\sigma^2(x)$  Functions



**Figure 4:** Comparison of True, NLS, and OLS Estimates of  $\sigma^2(x)$  Functions

**Table 3:** NLS Estimates of  $\gamma$  using squared residuals from first stage as data

Parameter	$\gamma_0$	$\gamma_1$	$\gamma_2$	$\gamma_3$
True Value	-4.0	0.0	-0.2	0.0
NLS estimate	-4.022	0.006	-0.240	0.002
Standard Deviation	0.050	0.069	0.039	0.025

Some students may have used simple OLS, estimating a specification like

$$\hat{\epsilon}_i^2 = X_i \gamma + \eta_i \quad (8)$$



rather than the exponential specification in equation (7). This is also OK since we weren't specific about what type of tool to use to check for heteroscedasticity. The only disadvantage of the OLS specification in (8) over the exponential specification in (7) is that the latter doesn't guarantee that  $\hat{\sigma}^2(x) \geq 0$  for all  $x$ . However we find that even in the linear specification most of the predicted  $\hat{\sigma}^2(x_i)$  values are indeed positive. Figure 4 compares the predicted values of  $\sigma^2(x)$  using both specifications, and we can see that the negative predicted values of  $\sigma^2(x)$  occur at the extreme high and low values of the observed  $x$ 's.

**Question 5.** Now we consider full information maximum likelihood (FIML) estimation of Model III, which is the same as model I, but with an exponential specification for conditional heteroscedasticity. Thus the joint density for  $(y, x)$  is given by:

$$y_i = \exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3) + \varepsilon_i \quad (9)$$

where  $\varepsilon_i \sim N(0, \exp(\gamma_0 + \gamma_1 x_i + \gamma_2 x_i^2))$ , and  $x_i \sim N(\mu, \delta^2)$ . We want to consider simultaneous or "full information maximum likelihood" (FIML) estimation of the parameter vector  $\theta = (\beta, \gamma, \mu, \delta^2)$ . The log-likelihood function  $L_N(\theta)$  is given by:

$$\begin{aligned} L_N(\theta) &= \frac{1}{N} \sum_{i=1}^N \ln f(y_i, x_i; (\beta, \gamma)) \\ &= \frac{1}{N} \sum_{i=1}^N \ln g(y|x; (\beta, \gamma)) + \frac{1}{N} \sum_{i=1}^N \ln h(x|(\mu, \delta^2)) \\ &= \left\{ \frac{-1}{2N} \sum_{i=1}^N X_i \gamma - \frac{1}{2N} \sum_{i=1}^N \frac{[y_i - \exp(X_i \beta)]^2}{\exp(X_i \gamma)} \right\} \\ &\quad + \left\{ -\frac{N}{2} \ln \delta^2 - \frac{1}{2\delta^2} \sum_{i=1}^N [x_i - \mu]^2 \right\} - \ln(2\pi) \\ &\equiv L_N^c(\beta, \gamma) + L_N^m(\mu, \delta^2). \end{aligned}$$

where  $X_i = (1, x_i, x_i^2, x_i^3)$  and  $\gamma = (\gamma_0, \gamma_1, \gamma_2, \gamma_3)'$ . The gradients of  $L_N(\theta)$  with respect to  $(\beta, \gamma)$  are:

$$\begin{aligned} \frac{\partial}{\partial \beta} L_N^c(\beta, \gamma) &= -\frac{1}{N} \sum_{i=1}^N \frac{[y_i - \exp(X_i \beta)] \exp(X_i \beta) X_i'}{\exp(X_i \gamma)} \\ \frac{\partial}{\partial \gamma} L_N^c(\beta, \gamma) &= -\frac{1}{2N} \sum_{i=1}^N X_i' + \frac{1}{2N} \sum_{i=1}^N \frac{[y_i - \exp(X_i \beta)]^2}{\exp(X_i \gamma)} X_i'. \end{aligned}$$

The gradients for  $L_N$  with respect to  $(\mu, \delta^2)$  are the same as in Question 1. The hessian matrix for  $L_N^c$  with respect to  $(\beta, \gamma)$  is given by:

$$\frac{\partial^2}{\partial \beta \partial \beta'} L_N^c(\beta, \gamma) = \frac{1}{N} \sum_{i=1}^N [y_i - 2 \exp(X_i \beta)] \exp\{X_i(\beta - \gamma)\} X_i' X_i$$

$$\begin{aligned}\frac{\partial^2}{\partial \beta \partial \gamma'} L_N^c(\beta, \gamma) &= -\frac{1}{N} \sum_{i=1}^N [y_i - \exp(X_i \beta)] \exp\{X_i(\beta - \gamma)\} X_i' X_i \\ \frac{\partial^2}{\partial \gamma \partial \gamma'} L_N^c(\beta, \gamma) &= -\frac{1}{N} \sum_{i=1}^N \frac{[y_i - \exp(X_i \beta)]^2}{2 \exp(X_i \gamma)} X_i' X_i.\end{aligned}$$

It is easy to verify (using the law of iterated expectations) that when  $(\beta, \gamma) = (\beta^*, \gamma^*)$ , the expectation of  $\partial^2 L_N^c(\beta, \gamma) / \partial \beta \partial \gamma' = 0$ , i.e. we have block diagonality between the  $\beta$  and  $\gamma$  parameters (assuming the model is correctly specified). Similarly one can verify that the  $(\beta, \gamma)$  block of the information matrix  $\mathcal{I}$  is zero. This implies that the asymptotic covariance between the maximum likelihood estimates  $\hat{\beta}$  and  $\hat{\gamma}$  is zero, so they are asymptotically independently distributed. This independence suggests the following 2-step procedure to obtain initial consistent estimates of  $(\beta^*, \gamma^*)$ : 1) estimate  $\beta$  by NLS (see attached Gauss code `eval_nls.g` and shell program `nlreg.gpr`), 2) use the estimated squared residuals  $\epsilon_i^2$  to estimate the  $\gamma$  parameters by NLS using the exponential specification in equation (7). We did this using the same `eval_nls.g` procedure we used for step 1, with a slight modification of `nlreg.gpr` to substitute  $\{\hat{\epsilon}_i^2\}$  instead of  $\{y_i\}$  as the dependent variable in the regression.

However we can do even better than this. We can do a 3rd step, weighted NLS or feasible generalized least squares (FGLS) estimation of  $\beta$  using the estimated conditional variance  $\hat{\sigma}^2(x_i) = \exp\{X_i \hat{\gamma}\}$  from step 2 as weights. The procedure `eval_fgls.g` provides the code to do the FGLS estimation. Due to the block diagonality property, it is not hard to show that the FGLS estimates of  $\beta^*$  have the same asymptotic distribution as the MLE: i.e. FGLS is asymptotically efficient in this case. To see this, note that the gradient and hessian of  $L_N^c(\beta, \gamma)$  with respect to  $\beta$  is the same as the gradient and hessian for the following FGLS criterion function:

$$\hat{\beta}_{\text{fgls}} = \underset{\beta \in R^4}{\operatorname{argmax}} -\frac{1}{N} \sum_{i=1}^N \frac{[y_i - \exp\{X_i \beta\}]^2}{\exp\{X_i \hat{\gamma}\}} \quad (10)$$

We know that the block diagonality property implies that as long as  $\hat{\gamma}$  is any consistent estimator of  $\gamma^*$  that a solution  $\hat{\beta}$  to the FOC  $\partial L_N^c(\beta, \hat{\gamma}) / \partial \beta = 0$  is asymptotically efficient. But since this is also the FOC for the FGLS estimator (10), it follows that  $\hat{\beta}_{\text{fgls}}$  is also an asymptotically efficient estimator, i.e. it attains not only the Chamberlain efficiency bound for condition moment restrictions, but the Cramer-Rao lower bound as well. It is not hard to show that these two bounds coincide in this case: make sure you understand this by verifying the equality yourself.

It is not apparent that the  $\hat{\gamma}$  obtained from the 4th step of our estimation procedure which regresses the squared residuals from the FGLS estimation in step 3 on  $\exp\{X_i \gamma\}$  will be asymptotically efficient since the first order condition for  $\gamma$  from maximizing  $L_N^c(\hat{\beta}, \gamma)$  with respect to  $\gamma$  does not appear to

be the same as the FOC for  $\gamma$  from the nonlinear regression in step 4 of our suggested estimation procedure. So to get fully efficient estimates, we can use  $(\hat{\beta}_{\text{fgls}}, \hat{\gamma}_{\text{fgls}})$  as starting values for direct FIML estimation of the full parameter vector  $(\beta^*, \gamma^*)$  using  $L_N^c(\beta, \gamma)$ . The procedure `eval_fiml.g` and the shell program `mle.gpr` implement full maximum likelihood estimation of Model III (note we have also provided the procedure `hesschk.g` to allow you to compare numerical and analytically calculated values of the hessian matrix, verifying that the analytic formulas for the hessian matrix given above are correct). This full model is rather delicate and we were unable to get it to converge starting from  $(\beta, \gamma) = 0$ . However we had no problems with convergence starting from  $(\hat{\beta}_{\text{fgls}}, \hat{\gamma}_{\text{fgls}})$ . Table 4 below compares the FGLS and MLE estimates.

**Table 4:** Comparison of MLE and FGLS Estimates of Model III

Parameter	Truth	MLE	Std. Error	FGLS	Std. Error
$\beta_0$	-3.000	-3.256	0.190	-3.256	0.190
$\beta_1$	2.000	2.209	0.777	2.200	0.789
$\beta_2$	-4.000	-2.558	0.819	-2.548	0.833
$\beta_3$	-0.010	-0.966	0.295	-0.962	0.300
$\gamma_0$	-4.000	-4.012	0.044	-4.022	0.050
$\gamma_1$	0.000	0.030	0.053	0.006	0.069
$\gamma_2$	-0.200	-0.254	0.026	-0.239	0.039
$\gamma_3$	0.000	-0.010	0.014	0.002	0.025

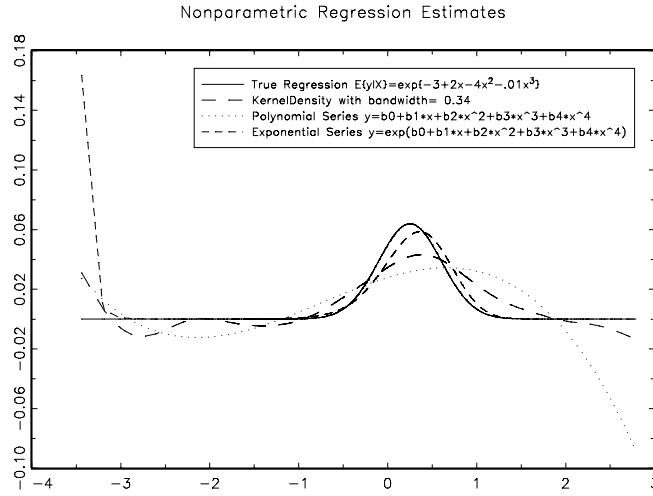
We see that the FIML and FGLS estimates of  $\beta$  are very close to each other and the standard errors are nearly identical, as we would expect from the theoretical result shown above that the FGLS estimator of  $\beta$  is asymptotically efficient. There are more significant differences in the FIML and FGLS estimates of  $\gamma$ . In particular the standard errors of the FGLS estimates are significantly larger than the FIML estimates of  $\gamma$ , which suggests that the FGLS estimates are not asymptotically efficient. Students should be able to verify that this is the case by deriving analytic formulas for the asymptotic covariance matrix for the MLE and FGLS estimators of  $\gamma$ .

**Question 6** The Nadaraya-Watson estimator is used to provide a nonparametric estimation of  $y = f(x) + \varepsilon$ :

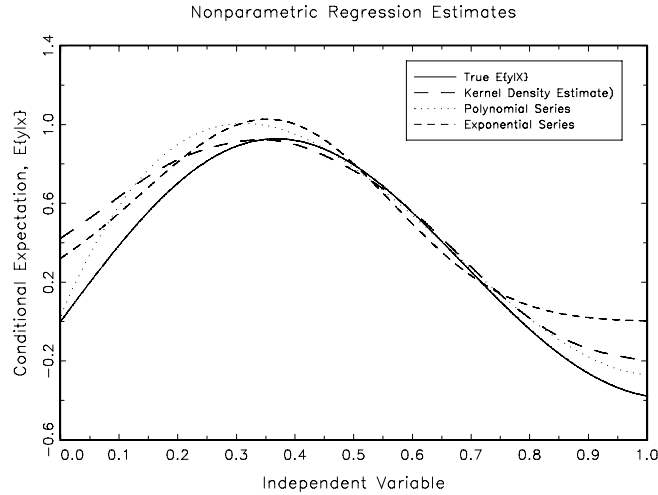
$$\hat{f}(x) = \frac{\sum_{i=1}^N K_h(x_i - x) y_i}{\sum_{i=1}^N K_h(x_i - x)}, \quad (11)$$

where  $K_h(x_i - x) = \frac{1}{h} K\left(\frac{x_i - x}{h}\right)$  and  $K(\cdot)$  is defined to be a Gaussian density function. For the choice of a bandwidth parameter,  $h$ , a rule of thumb is used:  $\hat{h} = c * \text{std}(x) * n^{-\frac{1}{5}}$ , with  $c \in [0.5, 2.5]$ . The Gauss program that computed these estimates is `kernel.gpr`. You should experiment with different bandwidths, showing that when  $h$  is much less than the automatically chosen value of  $h = .34$  the fitted regression tends to be too wiggly, and too smooth

when  $h$  is substantially greater than  $h = 0.34$ , so in this case the automatically chosen bandwidth seems like the best compromise. There are more sophisticated ways to choose bandwidths such as “least squares cross validation” described in class. The simple rule given above does almost as well and is much simpler to compute. The `kernel.gpr` program also includes two types of polynomial series estimation, ordinary polynomials and nonlinear least squares using the exponential specification considered in problem 4. Figures 5 and 6 compare the different estimators for the data in `data1.asc` and `data2.asc`, respectively. We can see that all of the estimators give similar answers in the region of the  $x$ -space where most of the data points lie, i.e. around  $x = 0$ .

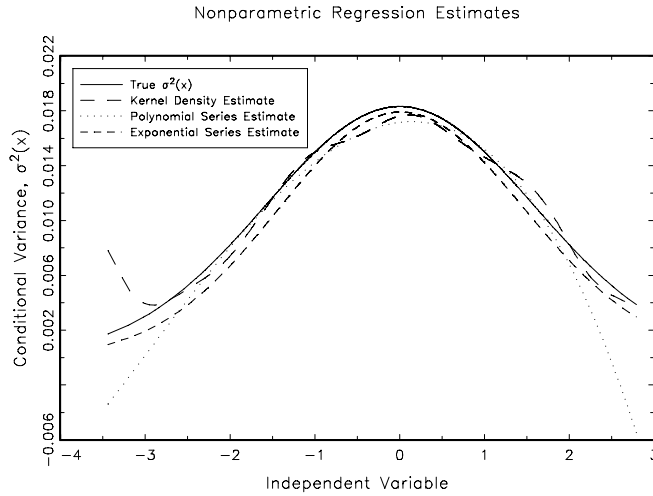


**Figure 5:** Nonparametric Estimates of Conditional Expectation Using `data1.asc`

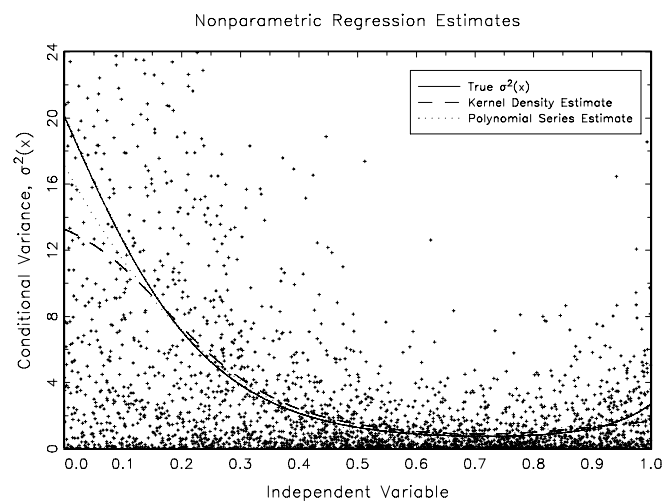


**Figure 6:** Nonparametric Estimates of Conditional Mean Expectation Using `data2.asc`

**Questions 7 and 8** To explore the form of heteroscedasticity, we repeat the procedures described above, but using the squared residuals from Question 6 as dependent variable. Figure 7, which plots true and estimated conditional variance functions  $\hat{\sigma}^2(x)$ , shows that similar to the results on estimation of the conditional expectation in part 6, all of the different parametric and nonparametric estimation methods give similar answers in the region of the  $x$ -space where most of the data points lie, i.e. around  $x = 0$ . Note particularly that the parametric exponential specification of  $\sigma^2(x)$  and the nonparametric kernel density estimate of  $\sigma^2(x)$  are quite close to each other except at extreme values of  $x$ . The same story also emerges in Figure 8, which plots the estimated squared residuals from the first stage NLS estimates of the exponential specification of the conditional expectation function. We didn't plot estimates of  $\hat{\sigma}^2(x)$  from the exponential specification,  $\sigma^2(x) = \exp(x\gamma)$  since we were unable to get the `nlreg.gpr` program to converge, even when we started the estimation from the true values of  $\gamma^*$ . This problem is probably a result of the large values chosen for the  $\gamma^*$  components,  $\gamma^* = (3.0, -4.0, -8.0, 10)$ , especially under the polynomial specification at extreme values of  $x$  where the computer runs into underflow or overflow problems. With some “playing around” one might be able to coax `nlreg.gpr` into convergence, but here is a case where the parameter estimates from the exponential specification seems rather fragile and non-robust.



**Figure 7:** Nonparametric Estimates of Conditional Variance Using `data1.asc`



**Figure 8:** Nonparametric Estimates of Conditional Variance Using `data2.asc`