

# Advanced Quantitative Methods in Political Science: Maximum Likelihood Estimation and Heteroskedastic Regression

---

Thomas Gschwend | Lisa-Marie Müller | Domantas Undzėnas

Week 5 - 11 March 2026

# Intro

---

## What should you take home from this class today?

- Log-likelihoods can be approximated around the maximum by a matrix of second derivatives (aka the *Hessian*) that measures the curvature in the neighborhood of the MLE.
- We get standard errors as square roots of diagonal terms of the VarCov matrix.
- We will implement our first MLE estimator in R and also estimate a (heteroskedastic) regression model.

# Three Steps to come up with a suitable ML Estimator for your Research Question

1. Formulate a suitable probability model of the data-generating process including assumptions of how  $Y$  is distributed (i.e., stochastic component) and a parametrization of stuff that gets estimated (i.e., systematic component).
2. Write down the (log-)likelihood function based on your parametrization and assumptions.
3. Maximize the log-likelihood, analytically (often hard, even impossible) or numerically (use functions in R).

## MLE and Standard Errors

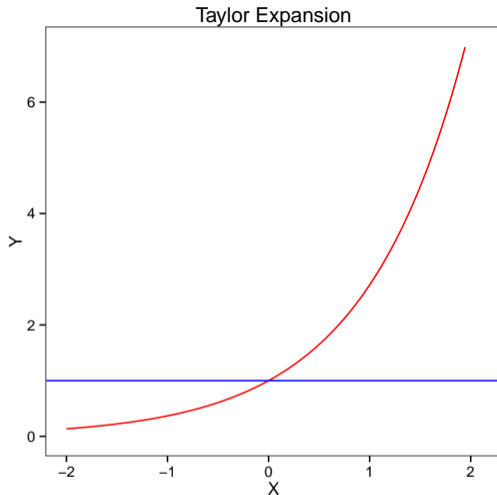
---

## Justifying Standard Errors

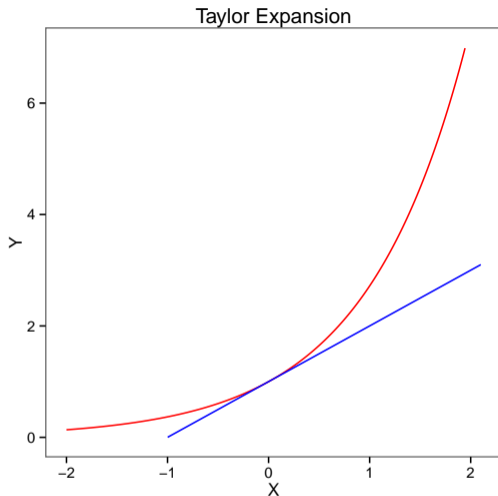
- The degree of *curvature* of the LL of the normal depends on the second derivative, because (remember from last week?) the LL of the normal is a quadratic polynomial around the MLE.
- This is generally not the case, but every (i.e. non-normal) LL can be **approximated** by a quadratic polynomial around the maximum.
- We take the *second order Taylor series expansion* of the log-likelihood with respect to  $\theta$  around the maximum  $\hat{\theta}$ :

$$f(\theta) = \ln L(\theta|y) \approx \ln L(\hat{\theta}|y) + \left( \frac{\partial \ln L(\hat{\theta}|y)}{\partial \theta} \right)' (\theta - \hat{\theta}) + \frac{1}{2} (\theta - \hat{\theta})' \frac{\partial^2 \ln L(\hat{\theta}|y)}{\partial \theta \partial \theta'} (\theta - \hat{\theta})$$

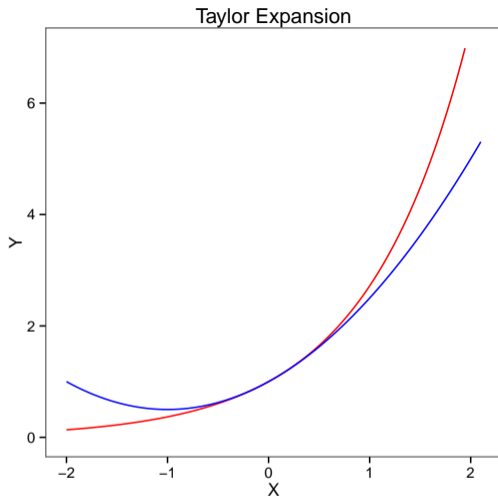
- This is fairly general. In fact, any function can be approximated by a quadratic function around a particular point (see Taylor series demonstration!)



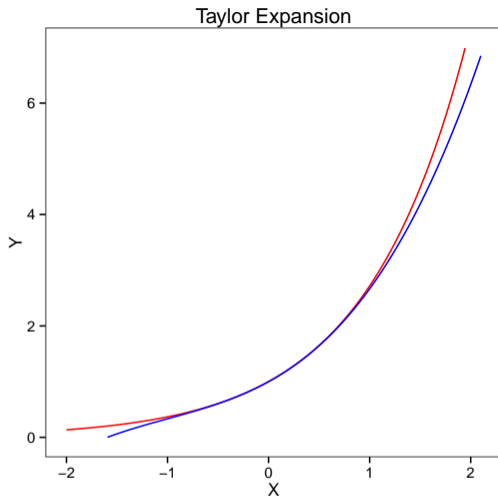
**Figure 1:** The exponential function,  $f(x) = e^x$ , and the Taylor series approximation:  $x_0 = 0$ ,  $f_0(x_1) = e^0 = 1$  (from Wikipedia)



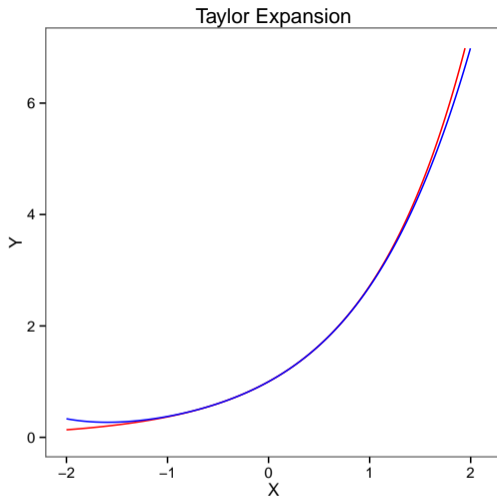
**Figure 2:** The exponential function,  $f(x) = e^x$ , and the Taylor series approximation:  $x_0 = 0$ ,  $f_1(x_1) = e^0 + e^0(x_1 - 0) = 1 + x_1$  (from Wikipedia)



**Figure 3:** The exponential function,  $f(x) = e^x$ , and the Taylor series approximation:  $x_0 = 0$ ,  $f_2(x_1) = e^0 + e^0(x_1 - 0) + \frac{1}{2}(x_1 - 0)^2 = 1 + x_1 + \frac{x_1^2}{2}$  (from Wikipedia)



**Figure 4:** The exponential function,  $f(x) = e^x$ , and the Taylor series approximation:  $x_0 = 0$ ,  $f_3(x_1) = 1 + x_1 + \frac{x_1^2}{2} + \frac{x_1^3}{6}$  (from Wikipedia)



**Figure 5:** The exponential function,  $f(x) = e^x$ , and the Taylor series approximation:  $x_0 = 0$ ,  $f_4(x_1) = 1 + x_1 + \frac{x_1^2}{2} + \frac{x_1^3}{6} + \frac{x_1^4}{24}$  (from Wikipedia)

- Instead of plotting the entire likelihood function, we can summarize the curvature in the neighborhood of the maximum with the *Fisher Information Matrix* denoted by  $\mathcal{I}(\theta|y)$ .
- The *information* in the data (i.e., degree of curvature) is defined as negative expectation in terms of the second derivatives (the so-called *Hessian Matrix*) of the log-likelihood with respect to  $\theta$

$$\mathcal{I}(\theta|y) = -E\left(\frac{\partial^2 \ln L(\theta|y)}{\partial \theta \partial \theta'}\right) = -E(H(\theta))$$

- Thus, the larger  $\mathcal{I}(\theta|y)$ , the more curved the log-likelihood and, thus, the more information in the data to estimate  $\hat{\theta}$ . Hence, we expect to get more precise estimates (i.e., smaller standard errors).
- More curvature  $\rightarrow$  data contain more information  $\rightarrow$  estimator is more precise.

The *Hessian Matrix*  $H(\theta)$  reflects the degree of curvature of the second-order approximation of the log-likelihood, i.e.,

$$H(\theta) = \left( \frac{\partial^2 \ln L(\theta|y)}{\partial \theta \partial \theta'} \right) = \begin{pmatrix} \frac{\partial^2 \ln L(\theta)}{\partial \theta_0 \partial \theta_0'} & \frac{\partial^2 \ln L(\theta)}{\partial \theta_0 \partial \theta_1'} & \cdots \\ \frac{\partial^2 \ln L(\theta)}{\partial \theta_1 \partial \theta_0'} & \frac{\partial^2 \ln L(\theta)}{\partial \theta_1 \partial \theta_1'} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

This is a square, symmetric matrix.

Given that we think about having more precision when having more information (about the curvature around the maximum), the variance-covariance matrix  $\text{Var}(\hat{\theta})$  is inversely related to the Information Matrix  $\mathcal{I}(\theta|y)$ , which is the negative of the expected value of the Hessian.

## Standard Errors

Thus,

$$\text{Var}(\hat{\theta}) = [\mathcal{I}(\theta|y)]^{-1} = [-E(H(\theta))]^{-1} = \left[ -E\left(\frac{\partial^2 \ln L(\theta|y)}{\partial \theta \partial \theta'}\right) \right]^{-1}$$

Because the above expression is evaluated at  $\theta$  to calculate an expected value but we only have  $\hat{\theta}_{ML}$ , we plug-in our estimate  $\hat{\theta}_{ML}$  (and use the “observed” information matrix) to estimate

$$\widehat{\text{Var}}(\hat{\theta}) = \left( \begin{array}{ccc} -\frac{\partial^2 \ln L(\hat{\theta})}{\partial \theta_0 \partial \theta_0'} & -\frac{\partial^2 \ln L(\hat{\theta})}{\partial \theta_0 \partial \theta_1'} & \cdots \\ -\frac{\partial^2 \ln L(\hat{\theta})}{\partial \theta_1 \partial \theta_0'} & -\frac{\partial^2 \ln L(\hat{\theta})}{\partial \theta_1 \partial \theta_1'} & \cdots \\ \vdots & \vdots & \ddots \end{array} \right)^{-1} = \left( -H(\hat{\theta}) \right)^{-1}$$

As you can see, we can read off the standard errors of  $\hat{\theta}_{ML}$  from the square roots of the diagonal elements of this matrix. Thus,  $\widehat{\text{Var}}(\hat{\theta})$  can be estimated as a function of the matrix of second derivatives. Remember (last week), standard errors are only correct asymptotically (as  $\hat{\theta}$  is asymptotically normal distributed)!

## Variance-Covariance Matrix of a Linear Regression Model

We start with the normal equations (see last week!) and look at the a *gradient* vector

$$\frac{\partial \ln L}{\partial \theta} = \begin{pmatrix} \frac{\partial \ln L}{\partial \beta} \\ \frac{\partial \ln L}{\partial \sigma^2} \end{pmatrix} = \begin{pmatrix} \frac{X'(y-X\beta)}{\sigma^2} \\ -\frac{N}{2\sigma^2} + \frac{(y-X\beta)'(y-X\beta)}{2\sigma^4} \end{pmatrix}$$

Next we take the derivative of each element of the gradient vector wrt  $\beta$  and  $\sigma^2$ .

$$\frac{\partial^2 \ln L}{\partial \beta \partial \beta'} = \frac{\partial \left( \frac{X'(y-X\beta)}{\sigma^2} \right)}{\partial \beta} = -\frac{X'X}{\sigma^2}$$

$$\frac{\partial^2 \ln L}{\partial \beta \partial \sigma^2} = \frac{\partial \left( \frac{X'(y-X\beta)}{\sigma^2} \right)}{\partial \sigma^2} = -\frac{X'\epsilon}{\sigma^4}$$

$$\frac{\partial^2 \ln L}{\partial \sigma^2 \partial \sigma^2} = \frac{\partial \left( -\frac{N}{2\sigma^2} + \frac{(y-X\beta)'(y-X\beta)}{2\sigma^4} \right)}{\partial \sigma^2} = \frac{N}{2\sigma^4} - \frac{\epsilon'\epsilon}{\sigma^6}$$

## Variance-Covariance Matrix of a Linear Regression Model

Subsequently, we can calculate the remaining entries of the *Hessian Matrix* as

$$H(\theta) = \left( \frac{\partial^2 \ln L}{\partial \theta \partial \theta'} \right) = \begin{pmatrix} -\frac{X'X}{\sigma^2} & -\frac{X'\epsilon}{\sigma^4} \\ -\frac{X'\epsilon}{\sigma^4} & \frac{N}{2\sigma^4} - \frac{\epsilon'\epsilon}{\sigma^6} \end{pmatrix}$$

Taking the expectation yields ...

$$E(H(\theta)) = E\left(\frac{\partial^2 \ln L}{\partial \theta \partial \theta'}\right) = \begin{pmatrix} E\left(-\frac{X'X}{\sigma^2}\right) & E\left(-\frac{X'\epsilon}{\sigma^4}\right) \\ E\left(-\frac{X'\epsilon}{\sigma^4}\right) & E\left(\frac{N}{2\sigma^4} - \frac{\epsilon'\epsilon}{\sigma^6}\right) \end{pmatrix} = \begin{pmatrix} -\frac{X'X}{\sigma^2} & 0 \\ 0 & -\frac{N}{2\sigma^4} \end{pmatrix}$$

Thus, the variance-covariance matrix is

$$\text{Var}(\hat{\theta}) = [-E(H(\theta))]^{-1} = \begin{pmatrix} \sigma^2(X'X)^{-1} & 0 \\ 0 & \frac{2\sigma^4}{N} \end{pmatrix}$$

and can be estimated through ...

$$\widehat{\text{Var}}(\hat{\theta}) = [-E(H(\hat{\theta}))]^{-1} = \begin{pmatrix} \hat{\sigma}^2(X'X)^{-1} & 0 \\ 0 & \frac{2\hat{\sigma}^4}{N} \end{pmatrix} = \begin{pmatrix} \text{Var}(\hat{\beta}_{OLS}) & 0 \\ 0 & \text{Var}(\hat{\sigma}_{MLE}^2) \end{pmatrix}$$

# Bootstrapping

- OK, I understand that we can use the LL to get standard errors.
- But what should I do in small samples, all standard errors in MLE are only correct asymptotically?
- Use **bootstrapping!**
  - A bootstrap provides a way to perform a statistical inference by re-sampling (i.e. drawing potentially infinite - or at least a really large number of samples) from the data you have.
  - Thus, assuming the data you have is equivalent to the population you wanna draw inferences to, bootstrap produces multiple samples (obviously by replacement) from the current population.
  - Based on every drawn sample calculate an estimate. Thus, you get a *bootstrap sampling distribution* of  $\theta$ .
  - The  $se(\hat{\theta}_{ML})$  would be the standard deviation of this distribution.
  - Use the relevant percentiles to construct confidence intervals.

## Likelihood ratio for nested models

- $L^*$  is the likelihood value of the *unrestricted* model.
- $L_R^*$  is the likelihood value of the (nested) *restricted* model.
- Thus,  $L^* \geq L_R^*$ , i.e.  $\frac{L_R^*}{L^*} \leq 1$ .
- Substantively, the likelihood ratio is a ratio of two probabilities (aka *risk ratio*):

$$\begin{aligned}\frac{L(\theta_1|y)}{L(\theta_2|y)} &= \frac{k(y) P(y|\theta_1)}{k(y) P(y|\theta_2)} \\ &= \frac{P(y|\theta_1)}{P(y|\theta_2)}\end{aligned}$$

- Statistically, let  $R = -2\ln\left(\frac{L_R^*}{L^*}\right) = 2(\ln L^* - \ln L_R^*)$ , then – under  $H_0$  of no difference between the two models –  $R$  is asymptotically  $\chi^2$  distributed, with the degree of freedom equal to the number of restrictions.

- Is  $\hat{\theta}_j$  systematically different from a theoretical  $\theta^*$ ?
- Generalized version of a  $t$ -test.
- Let  $\hat{\theta}_j$  the  $j$ th element of  $\hat{\theta}$ ,  $\hat{\sigma}_j$  its standard error, the square root of the  $j$ th diagonal element of the variance-covariance matrix. Then,

$$\mathcal{W} = \frac{\hat{\theta}_j - \theta^*}{\hat{\sigma}_j}$$

is asymptotically standard normal distributed, assuming  $H_0 : \theta_j = \theta^*$ .

- For the formal Wald test, we can instead also use that  $\mathcal{W}^2 \sim \chi^2(1)$ .

- The *score test* is aka *Lagrange multiplier test*.
- Again, if the null is valid, i.e.  $H_0 : \theta_j = \theta^*$ , then the restricted estimator should be near the point that maximizes the log-likelihood.
- Therefore, the respective slope should be near zero.
- Given that the *score*  $S(\theta^*)$  is the slope of the log-likelihood at  $\theta^*$ , it can be shown that the score statistic  $S$  with

$$S = \frac{S(\theta^*)}{\sqrt{\mathcal{I}(\theta^*)}}$$

is asymptotically standard normal distributed, assuming  $H_0 : \theta_j = \theta^*$ .

# Overview: Statistical Inference using ML

(Figure is taken from Fox, Appendix D)

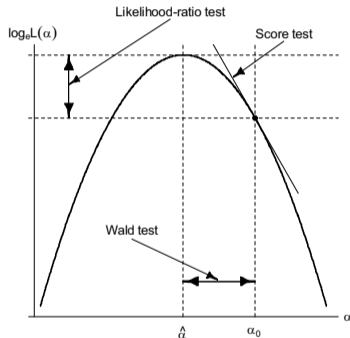


Figure D.19: Tests of the hypothesis  $H_0: \alpha = \alpha_0$ : The likelihood-ratio test compares  $\log_e L(\hat{\alpha})$  to  $\log_e L(\alpha_0)$ ; the Wald test compares  $\hat{\alpha}$  to  $\alpha_0$ ; and the score test examines the slope of  $\log_e L(\alpha)$  at  $\alpha = \alpha_0$ .

## To sum-up

- If the ML model is correct, then  $\hat{\theta}_{ML}$  is a consistent point estimate of  $\theta$ .
- As the number of observations become large, ...
  - ...the sampling distribution of  $\hat{\theta}_{ML}$  becomes normal.
  - ...the log-likelihood becomes quadratic.
  - ...the assumed second-order approximation of the log-likelihood improves.
- There are also several numerical algorithms (e.g. *Newton-Raphson*, *BHHH*, *Method of Scoring*, *L-BFGS-B*, *BFGS*, *simulated-annealing*) to find a maximum and estimate the variance-covariance matrix.

# Numerical Optimization

- Newton-Raphson works well and quickly for simple functions with global maxima
- Method of Scoring, BHHH and simulated-annealing can be better alternatives when likelihood is complex
- Some practical tips
  - The likelihood can have local maxima or saddle points with which numerical algorithms have a hard time (because they “think” its a global maximum).
  - Use different starting values. They should not matter if a global maximum is detected.
  - Use OLS instead to find first reasonable parameter values.
  - Graph LL by fixing all parameters (but 1 or 2) at reasonable values and graph the rest to eyeball maximum in order to find good starting values.
- When encountering convergence problems, ...
  - ...you may delete missing values explicitly and try again.
  - ...rescale the variables so that they are measured (ideally) on the same scale
  - ...try another numerical algorithm

## Implementation in R

---

- Let's estimate a linear regression model via maximum likelihood instead of using ordinary least squares
- *Step 1:* Assume the following model:  
$$Y_i \sim f_N(y_i | \mu_i, \sigma^2) \quad \text{stochastic}$$
$$\mu_i = X\beta (= \beta_0 + \beta_1 x_i) \quad \text{systematic}$$
- The parameters we are going to estimate using the above parameterization are  $\theta = (\mu_i, \sigma^2) = (\beta_0, \beta_1, \sigma^2)$
- We further assume that  $y_i$  is iid.

## Implementation in R

- *Step 2:* Using our assumptions about the model and the chosen parameterization of the systematic component, we can set up the likelihood function as follows:

$$L(\beta, \sigma^2 | y) = (2\pi\sigma^2)^{-N/2} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i)^2 \right]$$

- Then (although this is optional) we can take the log of the likelihood function, because it simplifies the next step (i.e. maximization):

$$\begin{aligned} \log L(\beta, \sigma^2 | y) &= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i)^2 \\ &= -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2} \sum_{i=1}^N \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\sigma^2} \\ &= -\frac{N}{2} \log(\sigma^2) - \frac{1}{2} \sum_{i=1}^N \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\sigma^2} \end{aligned}$$

## Implementation in R

- Now, let's write the log-likelihood as a R-function `lm.lik`:

```
lm.lik <- function(theta, y, x) {  
  beta0 <- theta[1]  
  beta1 <- theta[2]  
  gamma <- theta[3]  
  # Parametrize sigma2 to be non-negative  
  sigma2 <- exp(gamma)  
  # Residual  
  e <- y - beta0 - beta1*x  
  # Log lik function for one observation  
  logl <- -1/2*log(sigma2) - 1/2*(e^2/(sigma2))  
  # Log lik function is sum over N observations  
  logl <- sum(logl)  
  return(logl)  
}
```

## Implementation in R

- Here is a slightly more general code of the same likelihood:

```
lm.lik1 <-function(theta,y,X){
  N<-nrow(X) # number of observations
  k<-ncol(X) # number of parameters
  # Supstring paramters theta
  beta<-theta[1:k]
  gamma<-theta[k+1]
  # Parametrize sigma2 to be non-negative
  sigma2 <- exp(gamma)
  # Residual
  e<- y-(X%*%beta)
  # Log lik function fover N observations
  logl <- - 1/2*N*log(sigma2)-1/2*((t(e)%*%e)/(sigma2))
  return(logl)
}
```

## Implementation in R

- *Step 3:* Maximize the log-likelihood numerically. Of course, we could do it analytically (see last week). Now we let the computer do all the work for us.
- R provides a tool named `optim()` which maximizes arbitrary functions numerically if we specify `control=list(fnscale=-1)` (`optim()` tries to minimize by default).
- To maximize our likelihood function, we need to feed `optim()` with a set of starting values (the `optim(stval, ...)`'s first guesses for the parameters).

```
stval <- c(1,1,1)
```

- Then we simply call `optim()` to maximize a likelihood function (`fn=lm.lik`), with particular starting values (`stval`) and data (`y=y, x=x`)

```
res <- optim(stval, fn=lm.lik, control=list(fnscale=-1),  
            y=y, x=x, hessian=TRUE)
```

```
> res$par
```

```
[1] 49.708304  1.125821 10.378797
```

```
> sqrt(diag(solve(-1 * res$hessian)))
```

```
[1] 1.6249732 0.4578586 3.7924240
```

- Take some data and see how our  $\hat{\theta}_{ML}$  compares to  $\hat{\theta}_{OLS}$ !

# Heteroskedastic Regression

---

# Heteroskedastic Regression

- Now, what if we instead relax the homoskedasticity assumption?
- *Step 1:* Assume the following model:
  - $Y_i \sim f_N(y_i|\mu_i, \sigma_i^2)$  stochastic
  - $\mu_i = X\beta (= \beta_0 + \beta_1 x_i)$  systematic
  - $\sigma_i^2 = \exp(\gamma Z) (= \exp(\gamma_0 + \gamma_1 z_i))$  systematic
- The parameters we are going to estimate using the above parametrization of the model's systematic component are  $\theta = (\beta_0, \beta_1, \gamma_0, \gamma_1)$
- We further assume that the  $y_i$  are independently distributed.
- Thus, we get the following log-likelihood function:

$$\begin{aligned}\log L(\theta|y) &= -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \log(\sigma_i^2) - \frac{1}{2} \sum_{i=1}^N \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\sigma_i^2} \\ &= -\frac{1}{2} \sum_{i=1}^N (\gamma_0 + \gamma_1 z_i) - \frac{1}{2} \sum_{i=1}^N \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\exp(\gamma_0 + \gamma_1 z_i)}\end{aligned}$$

# Heteroskedastic Regression - Implementation in R

- Lets write the LL as a R-function `hetero.lik`, but this time with four arguments ( $\theta, y, x, z$ ):

```
hetero.lik <- function(theta, y, x, z) {  
  beta0 <- theta[1]  
  beta1 <- theta[2]  
  gamma0 <- theta[3]           # This line is new  
  gamma1 <- theta[4]           # This line is new  
  # Residual  
  e <- y - beta0 - beta1*x  
  # Variance parameterization  
  sigma2 = exp(gamma0 + gamma1*z) # This line is new  
  # Log lik function for one observation  
  logl <- -1/2*log(sigma2) - 1/2*(e^2/(sigma2))  
  # Log lik function is sum over N observations  
  logl <- sum(logl)  
  return(logl)  
}
```

- Note, we need to feed `optim()` with four starting values!

## Heteroskedastic Regression - Implementation in R

```
# start values for maximization algorithm - now we need 4 values
stval <- c(0,0,0,0)

# maximize the likelihood function numerically using optim()

res2 <- optim(stval,                # starting values
              fn=hetero.lik,        # the likelihood function
              control=list(fnscale=-1), # maximize rather than minimize funct
              y=y, x=x, z=z,        # the data
              hessian=TRUE)         # return numerical Hessian

cat("MLE Betas\n", res2$par[1:2], "\n\n")
cat("MLE Gammas\n", (res2$par[3:4]), "\n\n")
cat("Hessian\n")
print(res2$hessian)
cat("\n MLE St. Errors\n", sqrt(diag(solve(-1*res2$hessian))), "\n\n")
```

## Quo vadis AQM

---

# Infrastructure of “Advanced Quantitative Methods” Course

Three steps to come up with a suitable ML Estimator for your research question

1. Formulate a suitable probability model of the data-generating process including assumptions of how  $Y$  is distributed (i.e., stochastic component) and a parametrization of stuff that gets estimated (i.e., systematic component).
2. Write down the (log-)likelihood function based on your parametrization and assumptions.
3. Maximize the Log-Likelihood, analytically (often hard, even impossible) or numerically (use functions in R).

There are two more things we need to talk about this semester:

- Interpretation of estimation results through simulating quantities of interest (*you have seen this last semester as well as in the lab*)
- How to check whether the assumed model does fit the data? (*Coming soon!*)

Then, we can apply this infrastructure to any existing model or come-up with our own model.