

# Violating the homoskedasticity assumption of the linear regression model

Econometrics for minor Finance, Lecture 6

Tunga Kantarcı, Fall 2025

# Linear regression model

When we introduced the linear regression model, we assumed  
homoskedasticity:

$$\text{Var}[u_i | x] = \sigma^2$$

# Generalized linear regression model

What defines an econometric model is largely about the assumptions on

$$u_i$$

If the homoskedasticity and no autocorrelation assumptions about  $u_i$  hold, the linear model is the **standard** linear regression model. If one of them does not, so that

$$u_i$$

is heteroskedastic or autocorrelated, the linear model becomes the **generalized** linear regression model.

# Generalized linear regression model

In this lecture we focus on violating the homoskedasticity assumption.

# Generalized linear regression model: Model assumptions

All the assumptions we make for the linear regression model holds except that we now have **heteroskedasticity**

# Generalized linear regression model

In Greek, **hetero** means different, and **skedasis** means dispersion. Different dispersion. That is, non-constant variance.

# Generalized linear regression model

That is, if

$$E[u_i | x] = 0$$

and if

$$u_i$$

is heteroskedastic, for each observation  $i$ , we have

$$\begin{aligned}\text{Var}[u_i | x] &= E[u_i u_i | x] - E[u_i | x] E[u_i | x] \\ &= E[u_i u_i | x] \\ &= \sigma_i^2\end{aligned}$$

# Generalized linear regression model

$$\sigma_{\textcolor{red}{i}}^2$$

is shorthand for

$$\sigma^2 \omega(\textcolor{red}{x}_i)$$

for ease of notation. It says that the variance of  $u_i$  depends on the different values of an explanatory variable in some given functional form. We think of this as

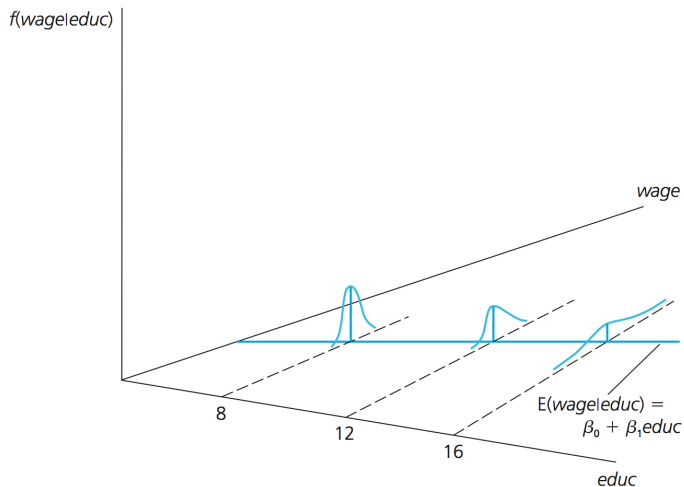
$$u_i$$

being drawn from a **different distribution** at each given

$$x_i$$



# Generalized linear regression model: Example



# Generalized linear regression model

For  $n$  errors, we have

$$\begin{aligned}\text{Var}[u \mid x] &= E[uu' \mid x] \\ &= \sigma^2 \Omega\end{aligned}$$

is the variance-covariance matrix, where

$$\Omega$$

is shorthand for

$$\Omega(x)$$

# Generalized linear regression model

How does

$$E[uu' | x]$$

look like?

# Generalized linear regression model

$$\begin{aligned} E[uu' | x] &= \begin{bmatrix} E[u_1 u_1 | x] & E[u_1 u_2 | x] & \dots & E[u_1 u_n | x] \\ E[u_2 u_1 | x] & E[u_2 u_2 | x] & \dots & E[u_2 u_n | x] \\ \vdots & \vdots & \ddots & \vdots \\ E[u_n u_1 | x] & E[u_n u_2 | x] & \dots & E[u_n u_n | x] \end{bmatrix} \\ &= \sigma^2 \underbrace{\begin{bmatrix} \omega_1 & 0 & \dots & 0 \\ 0 & \omega_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \omega_n \end{bmatrix}}_{\Omega} \end{aligned}$$

# Generalized linear regression model: Can we use the OLS estimator?

Does the OLS estimator

$$\hat{\beta}$$

have desirable statistical properties in the generalized linear regression model where the error

$$u_i$$

is heteroskedastic?

# Generalized linear regression model: Statistical properties of the OLS estimator: Unbiasedness

The answer is easy. Recall that for unbiasedness of

$$\hat{\beta}$$

only the exogeneity assumption

$$E[u_i | x_i] = 0$$

was required. This means that

$$\hat{\beta}$$

is still unbiased in the generalized model:

$$E[\hat{\beta} | \mathbf{x}] = \beta$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

For efficiency and asymptotic efficiency both the exogeneity assumption

$$E[u_i | x_i] = 0$$

and the homoskedasticity assumption

$$\text{Var}[u_i | x_i] = \sigma^2$$

were required.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

Now that homoskedasticity of

$$u_i$$

is violated,

$$\hat{\beta}$$

is not efficient and it is not asymptotically efficient. Let's see why.



# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

From an earlier lecture recall that we have

$$\text{Var} \left[ \hat{\beta}_1 \mid x \right] = \text{Var} \left[ \frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \mid x \right]$$

Expand the numerator variance:

$$\begin{aligned} \text{Var} \left[ \sum_{i=1}^n (x_i - \bar{x}) u_i \mid x \right] &= \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var} [u_i \mid x] \\ &\quad + \sum_{i \neq j} (x_i - \bar{x})(x_j - \bar{x}) \text{Cov} [u_i, u_j \mid x] \end{aligned}$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

Assume that errors are uncorrelated:

$$\text{Cov}[u_i, u_j \mid x] = 0$$

for all  $i \neq j$ .

Assume that errors are heteroskedastic:

$$\text{Var}[u_i \mid x] = \sigma_i^2$$

so they differ across  $i$ .

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

With these assumptions:

$$\text{Var} \left[ \sum_{i=1}^n (x_i - \bar{x}) u_i \middle| x \right] = \sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

For the denominator, as we condition on  $x$ , the fraction acts as constant, and is out of the variance operator as a square using var. property 2 in the first lecture slides:

$$\text{Var} \left[ \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x \right] = \left[ \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^2$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

The variance of

$$\hat{\beta}_1$$

becomes

$$\text{Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

In the linear regression model with

$$\text{Var}[u_i | x_i] = \sigma^2$$

we had

$$\text{Var}[\hat{\beta}_1 | x] = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

Now in the **generalized linear regression model** we have

$$\text{Var}[u_i | x_i] = \sigma_i^2$$

and

$$\text{Var}[\hat{\beta}_1 | x] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

Let's study how the sampling distribution of the OLS estimator

$$\hat{\beta}$$

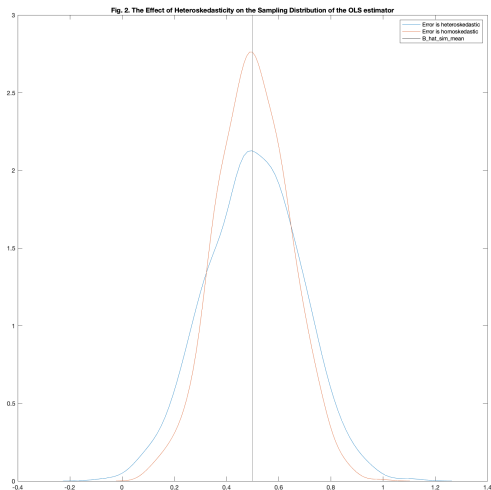
behaves when

$$u_i$$

is homoskedastic and when it is heteroskedastic.



# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency



# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency

In this simulation exercise the sampling distribution of

$$\hat{\beta}$$

has a larger variance in the generalized model. It can also be smaller. It depends on the functional form of heteroskedasticity:

$$\omega(x_i)$$

What matters is that OLS no longer guarantees the smallest possible variance.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

What are the implications of heteroskedasticity?

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

In the linear regression model with

$$\text{Var}[u_i | x_i] = \sigma^2$$

we had

$$\text{Var}[\hat{\beta}_1 | x] = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

and showed that

$$\text{Var}[\hat{\beta}_1 | x] \leq \text{Var}[\hat{\beta}_1^0 | x]$$

where  $\hat{\beta}_0$  was a competitor linear, unbiased estimator. That is, the OLS estimator

$$\hat{\beta}$$

was the most efficient estimator.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

Now in the **generalized linear regression model** we have

$$\text{Var}[u_i | x_i] = \sigma_i^2$$

and

$$\text{Var}[\hat{\beta}_1 | x] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

Is

$$\text{Var}[\hat{\beta}_1 | x] \leq \text{Var}[\hat{\beta}_1^0 | x]$$

still true in the generalized model? That is, is the OLS estimator

$$\hat{\beta}$$

still the most efficient estimator?

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

It is not. It can be shown that

$$\text{Var} \left[ \hat{\beta}^{GLS} \mid x \right] \leq \text{Var} \left[ \hat{\beta} \mid x \right]$$

where

$$\hat{\beta}^{GLS}$$

is the Generalized Least Squares estimator.

$$\hat{\beta}$$

is **not** the most efficient estimator in the generalized linear regression model.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

The Generalized Least Squares estimator

$$\hat{\beta}_1^{GLS} = \frac{\sum_{i=1}^n w_i (x_i - \bar{x}_w)(y_i - \bar{y}_w)}{\sum_{i=1}^n w_i (x_i - \bar{x}_w)^2}$$

where

$$w_i = \frac{1}{\sigma_i^2}$$

and

$$\bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}, \quad \bar{y}_w = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i}$$

is an alternative estimator which is efficient in the generalized linear regression model.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

We do not study this estimator here. We note, however, that using this estimator is another way to deal with heteroskedasticity or autocorrelation in regression errors. In practice we do not use this estimator because it is computationally costly. To deal with heteroskedasticity, we adapt another strategy, which we study next.



# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

Under heteroskedasticity, the OLS estimator remains unbiased but is inefficient, so we must correct the variance of it.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

In the **linear regression model** with homoskedasticity

$$\text{Var}[u_i | x_i] = \sigma^2$$

we had

$$\text{Var}[\hat{\beta}_1 | x] = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

We used this variance to calculate the  $t$  statistic, showing that it has an exact standard normal distribution if errors are assumed to be normal:

$$z = \frac{\hat{\beta} - \beta^0}{\sqrt{\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim N[0, 1]$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

Now in the generalized linear regression model we have

$$\text{Var}[u_i | x_i] = \sigma_i^2$$

and

$$\text{Var}[\hat{\beta}_1 | x] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

When the error is heteroskedastic, we **cannot use**

$$\text{Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

to calculate the **t** statistic. This is not the correct variance to use for this statistic. Should we use

$$\text{Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

instead? Yes, we should.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

The only problem is that  $\sigma$  is unobserved so we need to estimate it. We will derive this estimator later.

# Generalized linear regression model: Statistical properties of the OLS estimator: Efficiency: Implications

The same holds for the  $F$  statistic.

# Generalized linear regression model: Statistical properties of the OLS estimator: Normality

Recall from an earlier lecture that

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Assume that That is,

$$u_i \mid x_i \sim N[0, \sigma_i^2]$$

Does  $\hat{\beta}$  have a normal distribution in the generalized model?

## Generalized linear regression model: Statistical properties of the OLS estimator: Normality

Conditioning on  $x_i$ , the denominator is fixed. The numerator is a linear combination of the error terms  $u_j$ . Linear combinations of normal random variables are normal. Therefore

$$\hat{\beta}_1$$

is normal.



# Generalized linear regression model: Statistical properties of the OLS estimator: Normality

Using the mean and variance of

$$\hat{\beta}_1$$

derived above, we have

$$\hat{\beta}_1 \mid x \sim N \left[ \beta_1, \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

# Generalized linear regression model: Statistical properties of the OLS estimator: Consistency

The answer is easy. Recall that for consistency of

$$\hat{\beta}$$

only the exogeneity assumption

$$E[u_i | x_i] = 0$$

was required. This means that

$$\hat{\beta}$$

is still consistent in the generalized model.

# Generalized linear regression model: Statistical properties of the OLS estimator: Consistency

It can be shown that, in the generalized linear regression model, the variance of the OLS estimator

$$\hat{\beta}$$

approaches 0 when the sample size increases:

$$\text{Var} \left[ \hat{\beta} \mid x \right] \xrightarrow{p} 0$$

so that  $\hat{\beta}$  collapses to  $\beta$ :

$$\hat{\beta} \xrightarrow{p} \beta$$

Hence, the OLS estimator is consistent when

$$u_i$$

is heteroskedastic.

# Generalized linear regression model: Statistical properties of the OLS estimator: Asymptotic efficiency

The OLS estimator

$$\hat{\beta}$$

is asymptotically not efficient when

$$\text{Var}[u | x] = \sigma_i^2$$

We do not prove this. This means that increasing the sample size does not help to achieve efficiency. This is not surprising because when error variances differ, observations with smaller variance carry more precise information about the regression line. OLS treats all observations equally, instead of giving greater weight to those with lower error variance, and therefore fails to achieve efficiency.

Generalized least squares estimator treats all observations individually and achieves a lower asymptotic variance, while OLS cannot.

# Generalized linear regression model: Statistical properties of the OLS estimator: Asymptotic normality

Assume that  $u_i$  are uncorrelated but allow them to be heteroskedastic, that the random sampling assumption hold, and that the exogeneity assumption holds. Do **not** assume normality of  $u_i$ . Using the Central Limit Theorem,  $\hat{\beta}_1$  is asymptotically normal:

$$\hat{\beta}_1 \mid x \stackrel{a}{\sim} N \left[ \beta_1, \frac{1}{n} \frac{E[(x_i - \mu_x)^2 \sigma_i^2]}{E[(x_i - \mu_x)^2] E[(x_i - \mu_x)^2]} \right]$$

where  $\mu_x = E[x_i]$ . This is the asymptotic distribution of  $\hat{\beta}$ . It says that as  $N$  increases, the sampling distribution of  $\hat{\beta}$  approaches normality.

## Generalized linear regression model: Statistical properties of the OLS estimator: Asymptotic normality

Note that this result is for heteroskedasticity. The other type of violation is where  $u_i$  are correlated across  $i$ . Provided that this correlation diminishes with observations further away from each other,  $\hat{\beta}$  is also asymptotically normal when errors are autocorrelated.

## Generalized linear regression model: Heteroskedasticity consistent variance estimator

Consider the asymptotic variance of the OLS estimator shown above:

$$\text{Asy. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{1}{n} \frac{E \left[ (x_i - \mu_x)^2 \sigma_i^2 \right]}{E \left[ (x_i - \mu_x)^2 \right] E \left[ (x_i - \mu_x)^2 \right]}$$

The two expected value terms and  $\sigma$  are population terms. We need consistent estimators for these so that we can use it in practice. We also do not know the functional form of  $\sigma_i^2$ , that is  $\sigma^2(\omega_i)$ .

For what purpose we need this estimator is still to come.

# Generalized linear regression model: Heteroskedasticity consistent variance estimator

Relying on large  $N$ , a **consistent** estimator is

$$\text{Est. Asy. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Dropping the  $\frac{1}{n}$  terms,

$$\text{Est. Asy. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$



## Generalized linear regression model: Heteroskedasticity consistent variance estimator

$$\text{Est. Asy. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

This estimator is called the **heteroskedasticity consistent variance estimator**, HCVE, of

$$\text{Asy. Var} \left[ \hat{\beta} \right]$$

# Generalized linear regression model: Heteroskedasticity consistent variance estimator

We said that the  $t$  and  $F$  statistics are not valid if we use

$$\text{Est. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

But they are valid if we use the HCVE:

$$\text{Est. Asy. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

They are then called the heteroskedasticity-consistent  $t$  and  $F$  statistics.

## Generalized linear regression model: Heteroskedasticity consistent variance estimator

HCVE is powerful.  $\sigma^2\omega(x_i)$  is often unknown. HCVE does not need to figure this out. We can use the HCVE to make inference on  $\beta$ . We only need to keep in mind that the HCVE, and the test statistics that make use of the HCVE, require that  $N$  is large. We need this because otherwise the estimator is not consistent. We also do not need to assume that the errors are normal.

# Generalized linear regression model: Heteroskedasticity consistent variance estimator: Intuition

What is the intuition?

# Generalized linear regression model: Heteroskedasticity consistent variance estimator: Intuition

Under [homoskedasticity](#):

$$\text{Est. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

where

$$\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n \hat{u}_i^2}{n - K}}$$

So we have

$$\text{Est. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\frac{\sum_{i=1}^n \hat{u}_i^2}{n - K}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

# Generalized linear regression model: Heteroskedasticity

## consistent variance estimator: Intuition

Under **homoskedasticity**:

$$\text{Est. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\frac{\sum_{i=1}^n \hat{u}_i^2}{n - K}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Under **heteroskedasticity**:

$$\text{Est. Asy. Var} \left[ \hat{\beta}_1 \mid x \right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

Mind the deviation term for each  $i$ . Observations with regressors far from the mean contribute more to the variance estimate. This is **accounting** for heteroskedasticity.

## Generalized linear regression model: Heteroskedasticity consistent variance estimator: Intuition

$$\text{Asy. SEE} \left[ \hat{\beta}_1 \mid x \right] = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}}$$

The heteroskedasticity consistent SEE is often called White's estimator, or robust SEE, or simply robust SE.

# Generalized linear regression model: Heteroskedasticity consistent variance estimator: Example

```
. regress wage educ
```

Source	SS	df	MS	Number of obs	=	997
Model	7842.35455	1	7842.35455	F(1, 995)	=	251.46
Residual	31031.0745	995	31.1870095	Prob > F	=	0.0000
				R-squared	=	0.2017
				Adj R-squared	=	0.2009
Total	38873.429	996	39.0295472	Root MSE	=	5.5845

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	1.135645	.0716154	15.86	0.000	.9951106	1.27618
_cons	-4.860424	.9679821	-5.02	0.000	-6.759944	-2.960903



# Generalized linear regression model: Heteroskedasticity consistent variance estimator: Example

```
. regress wage educ, robust
```

Linear regression	Number of obs	=	997
	F(1, 995)	=	178.66
	Prob > F	=	0.0000
	R-squared	=	0.2017
	Root MSE	=	5.5845

wage	Robust					
	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	1.135645	.0849627	13.37	0.000	.9689186	1.302372
_cons	-4.860424	1.078429	-4.51	0.000	-6.976681	-2.744167

## Generalized linear regression model: Heteroskedasticity consistent variance estimator

We should always use the robust SEE because if there is no heteroskedasticity, it is equivalent to the standard SEE.

## GLS estimator or the HCV estimator?

If we detect heteroskedasticity, should we use the GLS estimator, which is a coefficient estimator, or the HCVE, which is a S.E. estimator? In the GLS approach, we alter the model and hence the coefficient estimator that is a function of the altered model. We also alter the standard errors since they depend on the altered coefficient estimator. The HCVE does not do anything to the coefficient estimator. It acknowledges heteroskedasticity and accounts for it in the standard error of the coefficient estimator. The advantage of the HCVE is that we do not need to figure out the covariance structure of the errors. It is difficult to obtain a reasonable estimate of the error covariance structure, and hence to use the GLS estimator.