

# Heteroskedasticity

## R Tutorials for Applied Statistics

**Note on required packages:** The following code requires the packages `lmtest`, `sandwich`, and `tidyverse`. The `sandwich` package contains procedures to estimate regression error variance that may change with the explanatory variables. The `lmtest` package contains procedures to conduct hypothesis tests when there is heteroskedasticity. If you have not done so, download and install the package `lmtest`, `sandwich`, and `tidyverse`. When the packages are installed, load these libraries.

```
# This only needs to be executed once for your machine
install.packages("lmtest")

# This only needs to be executed once for your machine
install.packages("sandwich")

# This only needs to be executed once for your machine
install.packages("tidyverse")

# This needs to be executed every time you load R
library("lmtest")

# This needs to be executed every time you load R
library("sandwich")

# This needs to be executed every time you load R
library("tidyverse")
```

---

## 1 Introduction

**Homoskedasticity** is the property that the variance of the error term of a regression (estimated by the variance of the residual in the sample) is the same across different values for the explanatory variables, or the same across time for time series models.

**Heteroskedasticity** is the property when the variance of the error term changes predictably with one or more of the explanatory variables.

Heteroskedasticity is common with financial variables related to income or spending. Suppose you are interested in predicting spending on housing and one of your explanatory variables is income. When income is low and the predicted value for spending on housing is low, the errors you make in the regression are also small. For people in low income groups, the amount each spends on housing varies little from other people in the low income group. For people in high income groups, the amount each spending on housing can vary greatly from other people in the high income group. The variance of your error term increases as the predicted value for  $\hat{y}$  increases. This is heteroskedasticity.

### Problems and non-problems with heteroskedasticity:

1. T-statistics and p-values from ordinary least squares (OLS) results assumes homoskedasticity.
2. Even with heteroskedasticity, the OLS estimates for your coefficients and predicted values for  $\hat{y}$  are still *unbiased*. OLS is still useful.
3. OLS estimates for the *variances* of your coefficients are biased with heteroskedasticity.
4. There is a straightforward correction to the OLS variances to allow for heteroskedasticity, holding on to your OLS estimates for the coefficients and predicted values.

## 2 Example: Factors affecting monthly earnings

Let us examine a data set that explores the relationship between total monthly earnings ( `MonthlyEarnings` ) and a number of factors that may influence monthly earnings including including each person's IQ ( `IQ` ), a measure of knowledge workplace environment ( `Knowledge` ), years of education ( `YearsEdu` ), years experience ( `YearsExperience` ), and years at current job ( `Tenure` ).

The code below downloads a CSV file that includes data on the above variables from 1980 for 663 individuals and assigns it to a data set called `df`.

```
load(url("http://murraylax.org/datasets/wage2.RData"))
```

The following call to `lm()` estimates a multiple regression predicting monthly earnings based on the five explanatory variables given above. The call to `summary()` displays some summary statistics from the regression.

```
lmwages <- lm(MonthlyEarnings
  ~ IQ + Knowledge + YearsEdu + YearsExperience + Tenure,
  data = df)
```

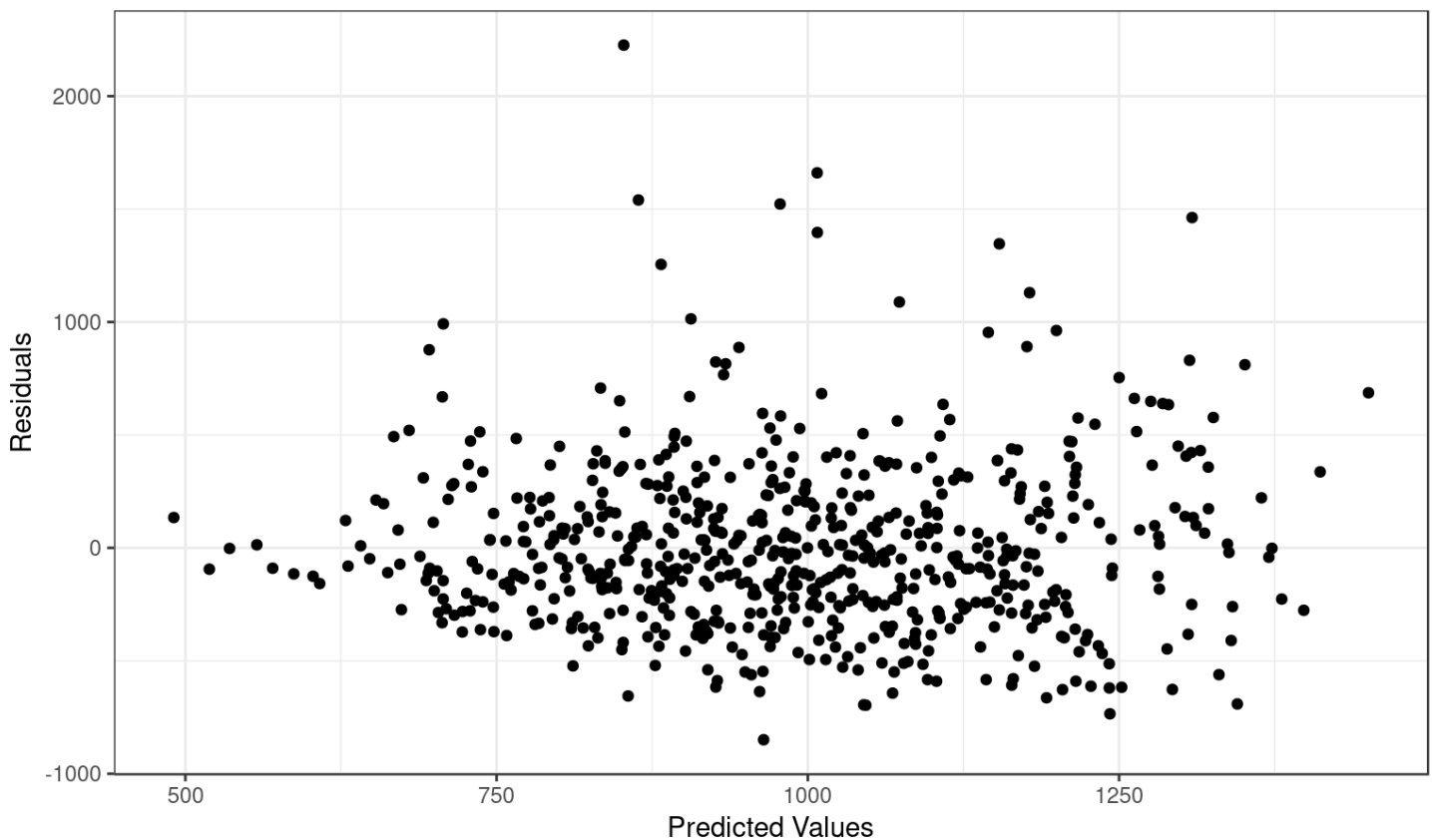
Let us save the residuals and predicted values from the regression into a new dataframe.

```
reg.df <- data.frame(resids=lmwages$residuals, predicts=lmwages$fitted.values)
```

Let's plot the residuals to see if they are more spread out for larger values of the predicted values:

```
ggplot(reg.df, aes(x=predicts, y=resids)) +
  geom_point() +
  theme_bw() +
  labs(title="Visual Inspection for Changing Variance",
       x="Predicted Values", y="Residuals")
```

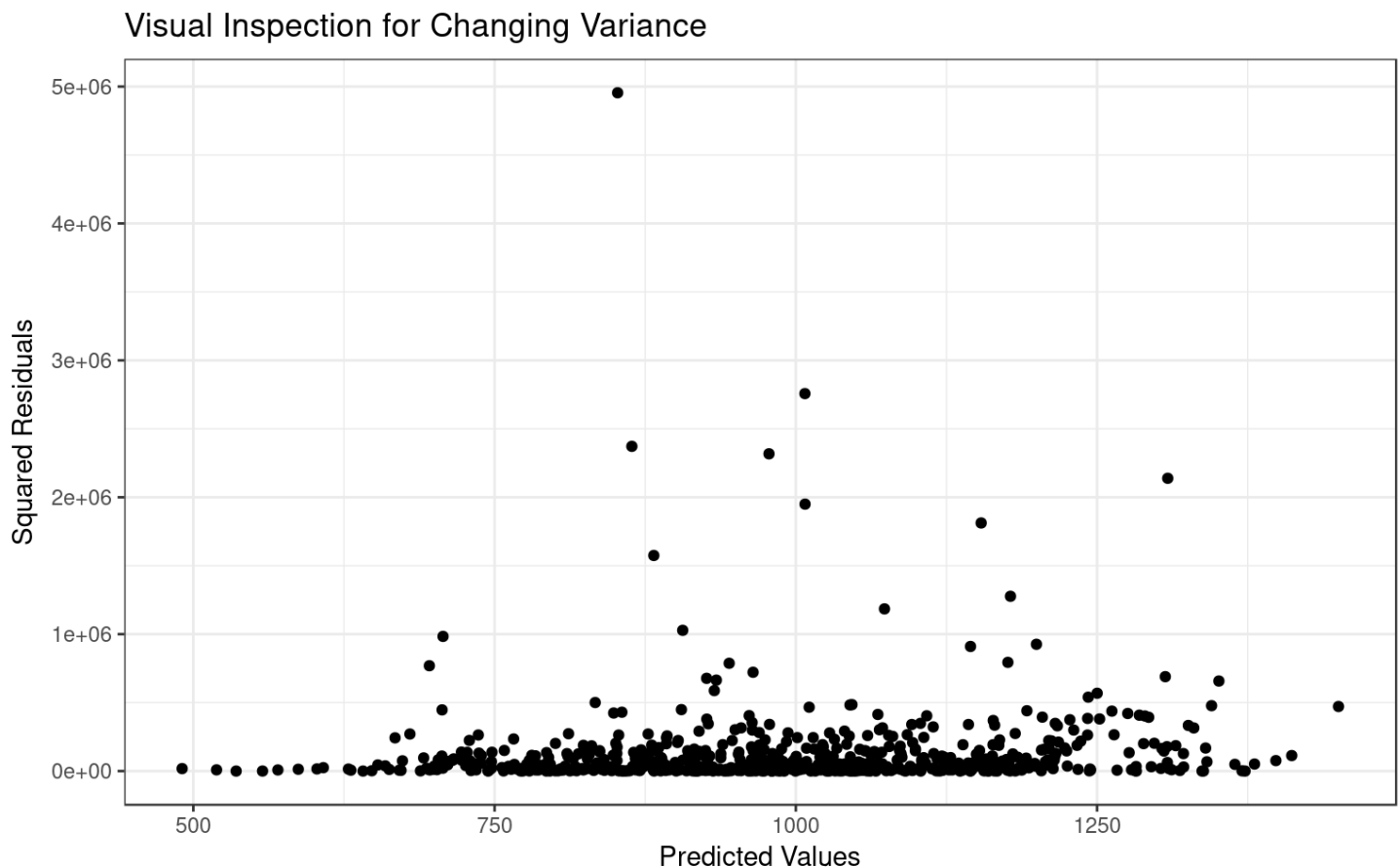
Visual Inspection for Changing Variance



It does appear that the distribution is less spread out at the left end. That is, at the lowest levels for predicted monthly earnings, the residuals are smaller.

Let us also plot the *squared* residual against the predicted values. The squared residual is an estimate of the variance of the error term:

```
ggplot(reg.df, aes(x=predicts, y=resids^2)) +  
  geom_point() +  
  theme_bw() +  
  labs(title="Visual Inspection for Changing Variance",  
        x="Predicted Values", y="Squared Residuals")
```



A visual inspection suggests there may be a pattern. As the predicted values increase, the variance of the error term increases. This suggests there may be heteroskedasticity.

Ordinary least squares assumes homoskedasticity. Said another way, ordinary least squares assumes the variance for the error term is the same for all predicted values. When we have heteroskedasticity, the variance for the residual on each observation is different, and given by each observation's squared residual shown in the graph above.

**In Section 3 below**, we learn how to correct the OLS estimates for the variances for the coefficients and predicted values making use of these squared residuals.

**In Section 4 below**, we learn how to more formally test for the presence for heteroskedasticity. It is based on the relationships we see in these graphs between the squared residuals and the predicted values.

### 3 Heteroskedasticity robust standard errors

**White's correction for heteroskedasticity** is a method for using OLS estimates for coefficients and predicted values (which are unbiased), but fixing the estimates for the variances of the coefficients (which are biased). It uses as an estimate for the *possibly changing* variance the squared residuals estimated from OLS which we computed and graphed above.

In R, we can compute the White heteroskedastic variance/covariance matrix for the coefficients with the call below to `vcovHC` from the `sandwich` package. VCOVHC stands for Variance / Covariance Heteroskedastic Consistent.

```
vv <- vcovHC(lmwages, type="HC1")
```

The first parameter in the call above is our original output from our call to `lm()` above. The second parameter `type="HC1"` tells the function to use the White estimate for the variance covariance matrix which uses as an estimate for the changing variance the squared residuals from the OLS call.

Then we can use this estimate for the variance / covariance to properly compute our standard errors, t-statistics, and p-values for the coefficients:

```
coeftest(lmwages, vcov = vv)
```

```
##
## t test of coefficients:
##
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -512.5397   135.5678  -3.7807 0.0001706 ***
## IQ              3.8887     1.1352   3.4256 0.0006517 ***
## Knowledge       8.0476     2.4932   3.2278 0.0013093 **
## YearsEdu       46.3084     8.6877   5.3303 1.349e-07 ***
## YearsExperience 13.4454     4.1497   3.2401 0.0012550 **
## Tenure          3.3915     3.0319   1.1186 0.2637138
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The first parameter to `coefrest` is the result of our original call to `lm()` and the second parameter is the updated variance / covariance matrix to use.

Let's compare the result to the OLS regression output:

```
summary(lmwages)
```

```
##
## Call:
## lm(formula = MonthlyEarnings ~ IQ + Knowledge + YearsEdu + YearsExperience +
##     Tenure, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -849.51 -244.91  -41.28  191.41 2225.88
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -512.540    139.180   -3.683 0.000250 ***
## IQ              3.889      1.204    3.230 0.001299 **
## Knowledge       8.048      2.246    3.582 0.000366 ***
## YearsEdu       46.308      8.833    5.243 2.13e-07 ***
## YearsExperience 13.445      4.064    3.309 0.000989 ***
## Tenure          3.392      3.016    1.124 0.261262
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 369.6 on 657 degrees of freedom
## Multiple R-squared:  0.1796, Adjusted R-squared:  0.1733
## F-statistic: 28.76 on 5 and 657 DF,  p-value: < 2.2e-16
```

The coefficients are exactly the same, but the estimates for the standard errors, the t-statistics, and p-values are slightly different.

One reason the standard errors and p-values are only slightly different is that there may actually be very little heteroskedasticity. The differences in variance of the error term may be small even for large differences in the predicted value for  $Xx$ .

The White heteroskedastic robust standard errors are valid for **either homoskedasticity or heteroskedasticity**, so it is always safe to use these estimates if you are not sure if the homoskedasticity assumption holds.

## 4 White test for heteroskedasticity

In section 2 we examined visual evidence that the magnitudes of the residuals or squared residuals were on average larger for larger predicted values for monthly earnings. In this section, we will test this relationship more formally.

The goal is to determine whether the squared residuals can be explained by the predicted values. Both, the squared residuals and predicted values from the OLS regression are unbiased, so we will use these values.

We run the following regression of the squared residuals on the predicted values and squared predicted values:

```
lmhet <- lm(I(resids^2) ~ predicts + I(predicts^2), data=reg.df)
summary(lmhet)
```

```
##
## Call:
## lm(formula = I(resids^2) ~ predicts + I(predicts^2), data = reg.df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -199840 -114791  -79308    8958 4840453
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.466e+04  3.244e+05   0.107   0.915
## predicts      4.887e+01  6.622e+02   0.074   0.941
## I(predicts^2) 5.205e-02  3.323e-01   0.157   0.876
##
## Residual standard error: 322800 on 660 degrees of freedom
## Multiple R-squared:  0.006593,    Adjusted R-squared:  0.003583
## F-statistic:  2.19 on 2 and 660 DF,  p-value: 0.1127
```

Do the predicted values from the regression explain the squared residuals? We answer this with the joint F-test.

$H_0 : \beta_{\text{predicts}} = 0; \beta_{\text{predicts}^2} = 0$  (i.e. homoskedasticity)

$H_0: \beta_{\text{predicts}} = 0; \beta_{\text{predicts}^2} = 0$  (i.e. homoskedasticity)

$H_A : \text{Either } \beta_{\text{predicts}} \neq 0 \text{ or } \beta_{\text{predicts}^2} \neq 0$  (i.e. heteroskedasticity)

$H_A: \text{Either } \beta_{\text{predicts}} \neq 0 \text{ or } \beta_{\text{predicts}^2} \neq 0$  (i.e. heteroskedasticity)

If the null hypothesis is true and both coefficients are equal to zero, then the predicted values of the regression do not explain the squared residuals. In this case we say we have homoskedasticity: the variance of the residuals do not change predictably with the predicted values.

If the alternative hypothesis is true then we have statistical evidence that either the predicted values or the squared predicted values helps explain the variance of the error term. In this case we have heteroskedasticity: As the predicted values change, so does the variance of the error term.

The p-value is 0.1127. We fail to find sufficient statistical evidence that the model exhibits heteroskedasticity. This may mean that we do not have heteroskedasticity, or it may mean we do not have enough power to find sufficient statistical evidence. The White heteroskedastic-robust test on the coefficients is safe with or without heteroskedasticity.

## 5 Functional Form



Heteroskedasticity on its own is a small problem. OLS estimates for coefficients are still unbiased, and it is easy to correct standard errors and p-values to allow for the possibility of heteroskedasticity.

However, sometimes heteroskedasticity is the result of a *misspecified model*. That is, you estimated a functional form for the model that is not appropriate, that is not true in the data.

It is common for financial variables and variables related to income, that larger numbers grow exponentially. When we take the natural log of these variables, the transformed values grow linearly, and linear regression becomes more applicable.

Let us revisit the model above, but instead use the *natural log of monthly earnings* as the outcome variable:

```
lmlogwages <- lm(log(MonthlyEarnings)
  ~ IQ + Knowledge + YearsEdu + YearsExperience + Tenure,
  data = df)
```

Let us compute the White test for heteroskedasticity:

```
logreg.df <- data.frame(resids=lmlogwages$residuals, predicts=lmlogwages$fitted.values)

lmhet <- lm(I(resids^2) ~ predicts + I(predicts^2), data=logreg.df)
summary(lmhet)
```

```
##
## Call:
## lm(formula = I(resids^2) ~ predicts + I(predicts^2), data = reg.df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -199840 -114791  -79308    8958 4840453
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.466e+04  3.244e+05   0.107   0.915
## predicts     4.887e+01  6.622e+02   0.074   0.941
## I(predicts^2) 5.205e-02  3.323e-01   0.157   0.876
##
## Residual standard error: 322800 on 660 degrees of freedom
## Multiple R-squared:  0.006593,    Adjusted R-squared:  0.003583
## F-statistic:  2.19 on 2 and 660 DF,  p-value: 0.1127
```

The p-value for the joint F-test for regression significance is 0.1127. Again, we fail to reject the null hypothesis. We fail to find statistical evidence that the variance of the residual changes with the predicted or squared predicted values. We failed to find statistical evidence for heteroskedasticity.