# MQE: Economic Inference from Data: Odds and Ends

Claire Duquennois

2020

# Odds and Ends

- Non-Standard Standard Errors
    - Robust standard errors
    - Clustered standard errors
    - Newey West Standard Errors
    - Conley Standard errors
- Confidence intervals for prediction

# Non-standard standard errors

A standard error estimates the uncertainty around an estimated parameter.

Formally we have

$$se = \sqrt{\widehat{Var(\hat{\beta})}}.$$

Just like calculating point estimates, it is incredibly important to get your standard errors right.

You have to know what you don't know!

- ▶ Robust standard errors
- ▶ Clustered standard errors

# Robust standard errors

Using the diamonds data set from ggplot2:

```
knitr::kable(head(diamonds))
```

| carat | cut | color | clarity | depth | table | price | x | y | z |
|---|---|---|---|---|---|---|---|---|---|
| 0.23 | Ideal | E | SI2 | 61.5 | 55 | 326 | 3.95 | 3.98 | 2.43 |
| 0.21 | Premium | E | SI1 | 59.8 | 61 | 326 | 3.89 | 3.84 | 2.31 |
| 0.23 | Good | E | VS1 | 56.9 | 65 | 327 | 4.05 | 4.07 | 2.31 |
| 0.29 | Premium | I | VS2 | 62.4 | 58 | 334 | 4.20 | 4.23 | 2.63 |
| 0.31 | Good | J | SI2 | 63.3 | 58 | 335 | 4.34 | 4.35 | 2.75 |
| 0.24 | Very Good | J | VVS2 | 62.8 | 57 | 336 | 3.94 | 3.96 | 2.48 |

# Robust standard errors

## Regress price on carats and depth.

```
reg1<-felm(price~carat+depth, diamonds)

summary(reg1)
```

```
##
## Call:
##    felm(formula = price ~ carat + depth, data = diamonds)
##
## Residuals:
##      Min      1Q  Median      3Q     Max
## -18238.9  -801.6   -19.6   546.3 12683.7
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 4045.333    286.205   14.13   <2e-16 ***
## carat       7765.141     14.009  554.28   <2e-16 ***
## depth       -102.165      4.635  -22.04   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1542 on 53937 degrees of freedom
## Multiple R-squared(full model): 0.8507    Adjusted R-squared: 0.8507
## Multiple R-squared(proj model): 0.8507    Adjusted R-squared: 0.8507
## F-statistic(full model):1.536e+05 on 2 and 53937 DF, p-value: < 2.2e-16
## F-statistic(proj model): 1.536e+05 on 2 and 53937 DF, p-value: < 2.2e-16
```
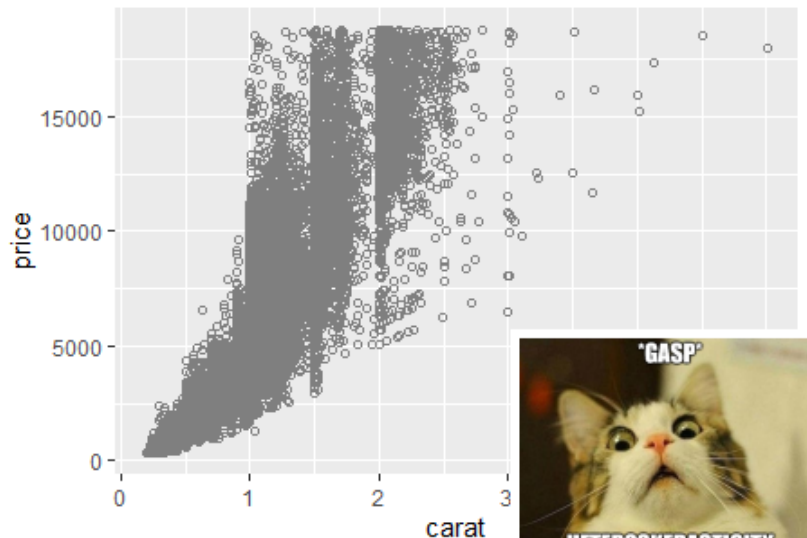
# Robust standard errors

Cool.

Plot the data to check OLS assumptions:

```
myPlot <- ggplot(data = diamonds, aes(y = price, x = carat)) +
geom_point(color = "gray50", shape = 21)
```

# Robust standard errors

# Robust standard errors

You should have the econometric heebie jeebies.

Homoskedastic assumption needed for OLS is not valid!

- ▶ The higher the carat, the greater the variance in price.
- ▶ ⇒ OLS standard errors are likely to be wrong.

Thankfully all is not lost!

# Robust standard errors

Lets relax the homoskedasticity assumption and allow for the variance to depend on the value of $x_i$.

We know that

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

With heteroskedasticity $\sigma^2$ is no longer constant and becomes a function of the particular value of $x_i$ an observation has, so

$$Var(u_i|x_i) = \sigma_i^2$$

Where are we going to find all these $\sigma_i^2$ for each individual observation?

# Eicker, Huber and White to the rescue!

Econometricians Eicker, Huber and White figured out a way to do this by basically using the square of the estimated residual of each observation, $\hat{u}_i^2$, as a stand-in for $\sigma_i^2$.

With this trick, a valid estimator for $Var(\hat{beta}_1)$, with heteroskedasticity of **any** form (including homoskedasticity), is

$$Var(\hat{\beta}_1) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{(\sum_{i=1}^n (x_i - \bar{x})^2)^2}$$

We commonly call the resulting standard errors "robust", or "heteroskedasticity-robust".

# Robust standard errors

## How can we find these in R?

```
reg1<-felm(price~carat+depth, diamonds)

summary(reg1, robust=TRUE)
```

```
##
## Call:
##    felm(formula = price ~ carat + depth, data = diamonds)
##
## Residuals:
##     Min      1Q   Median      3Q      Max
## -18238.9  -801.6   -19.6   546.3  12683.7
##
## Coefficients:
##             Estimate Robust s.e t value Pr(>|t|)
## (Intercept) 4045.333    369.176   10.96   <2e-16 ***
## carat       7765.141     25.105  309.31   <2e-16 ***
## depth       -102.165      5.946  -17.18   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1542 on 53937 degrees of freedom
## Multiple R-squared(full model): 0.8507   Adjusted R-squared: 0.8507
## Multiple R-squared(proj model): 0.8507   Adjusted R-squared: 0.8507
## F-statistic(full model, *iid*):1.536e+05 on 2 and 53937 DF, p-value: < 2.2e-16
## F-statistic(proj model): 4.878e+04 on 2 and 53937 DF, p-value: < 2.2e-16
```

## Robust standard errors

Or if you want to put them in a stargazer table:

```
stargazer(reg1, type = "latex" , se = list(reg1$rse), header=FALSE)
```

Table 2

|  | *Dependent variable:* |
| --- | --- |
|  | price |
| carat | 7,765.141$^{***}$ |
|  | (25.105) |
| depth | $-102.165^{***}$ |
|  | (5.946) |
| Constant | 4,045.333$^{***}$ |
|  | (369.176) |
| Observations | 53,940 |
| $R^2$ | 0.851 |
| Adjusted $R^2$ | 0.851 |
| Residual Std. Error | 1,541.649 (df = 53937) |
| *Note:* | $^{*}$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01 |

Note: robust standard errors are larger than regular standard errors, and thus more conservative (which is the right thing to be... you want to know what you don't know).

# Econometricians Haiku

*T-stats looks too good*

*Try cluster standard errors*

*significance gone.*

from Angrist and Pischke 2008

# Clustered standard errors

Suppose that every observation belongs to (only) one of G groups.

The assumption we make when we cluster:

- there is no correlation across groups
- we allow for arbitrary within-group correlation.

# Clustered standard errors

Example: consider individuals within a village.

It may be reasonable to think that individuals' error terms are:

- ▶ correlated within a village
- ▶ aren't correlated across villages

# Clustered standard errors

I will spare you the matrix math needed to dive deeper into this.

Suffice to say that "cluster-robust" estimates allow for a more complicated set of correlations to exist within observations within a cluster.

One thing to be aware of though is that you will need to have a fairly large number of clusters (40+) for the estimate to be credible.

# Clustered standard errors

Clustering in R:

I use the `NOxEmissions` dataset from the `robustbase` package.

- hourly $NO_x$ readings, including $NO_x$ concentration, auto emissions and windspeed.
- use the observation date as our cluster variable.

This allows for arbitrary dependence between observations in the same day, and zero correlation across days.

Is this reasonable? . . . Maybe. But we'll go with it for now:

# Clustered standard errors

```
nox <- as.data.frame(NOxEmissions) %>%mutate(ones = 1)
noClusters <- felm(data = nox, LNOx ~ sqrtWS )
Clusters <- felm(data = nox, LNOx ~ sqrtWS |0|0| julday)
stargazer(noClusters,Clusters, type = "latex" , header=FALSE, omit.stat = "all")
```

Table 3

|  | Dependent variable: | |
|---|---|---|
|  | LNOx | |
|  | (1) | (2) |
| sqrtWS | $-0.864^{***}$ | $-0.864^{***}$ |
|  | (0.020) | (0.048) |
| Constant | $5.559^{***}$ | $5.559^{***}$ |
|  | (0.029) | (0.065) |

Note: $^{*}p<0.1$; $^{**}p<0.05$; $^{***}p<0.01$

Here, the regular standard errors are smaller than the clustered standard errors.

This need not necessarily be the case and depends on the correlation between observations within a cluster.

# Newey West Standard Errors

For time series data.

# Conley Standard Errors

For spatial data.

# Confidence intervals for predictions

You know how to "predict" a value of the dependent variable, $y$, given certain values of the independent variables.

This prediction is just a guess, with uncertainty.

We can construct a confidence interval to give a range of possible values for this prediction.

There are two kinds of predictions we can make:

- ▶ A confidence interval for the **average** $y$ given $x_1, x_2 \ldots x_k$.
- ▶ A confidence interval for a **particular** $y$ given $x_1, x_2 \ldots x_k$.

# Confidence intervals for predictions

Using Wooldridge's birth weight data:

$$bweight = \beta_0 + \beta_1 lfaminc + \beta_2 meduc + \beta_3 parity + u$$

▶ bwght is birth weight in ounces,

▶ lfaminc is the log of family income in \$1000s,

▶ meduc is the education of the mother in years,

▶ parity is the birth order of the child.

# Confidence intervals for predictions

Estimating this equation in R, we get the following results:

```r
#using the bwght data from the wooldridge package
reg1<-lm(bwght~lfaminc+motheduc+parity, bwght)

summary(reg1)
```

```
##
## Call:
## lm(formula = bwght ~ lfaminc + motheduc + parity, data = bwght)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -94.533 -11.888   0.779  13.136 151.477
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 105.5652     3.3666  31.356  < 2e-16 ***
## lfaminc       2.1313     0.6506   3.276  0.00108 **
## motheduc      0.3172     0.2520   1.259  0.20829
## parity        1.5261     0.6119   2.494  0.01275 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 20.21 on 1383 degrees of freedom
##   (1 observation deleted due to missingness)
## Multiple R-squared:  0.01633,    Adjusted R-squared:  0.0142
## F-statistic: 7.654 on 3 and 1383 DF,  p-value: 4.482e-05
```

# Confidence intervals for predictions: for a specific average

Our model gives us the expected value:

$E[bweight|faminc, meduc, parity] = \beta_0 + \beta_1 log(faminc) + \beta_2 meduc + \beta_3 parity$

and our regression gives us an estimate of this:

$\hat{E}[bweight|faminc, meduc, parity] = \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 log(faminc) + \hat{\beta}_2 meduc + \hat{\beta}_3 parity$

$\hat{y}$ is the expected value of y given the particular values for the explanatory variables.

# Confidence intervals for predictions: for a specific average

Say we are interested in a confidence interval for the **average birthweight** for babies with:

- a family income of \$14,500 (ln(14.5)=2.674),
- mothers with 12 years of education,
- 2 older siblings (parity=3).

$$\hat{E}[bweight|faminc = 14.5, \, meduc = 12, \, parity = 3] = 105.66 + 2.13 ln(faminc) + 0.317 meduc + 1.53 parity$$

$$\hat{y}_{faminc=14.5, meduc=12, parity=3} = 105.66 + 2.13(2.674) + 0.317(12) + 1.53(3)$$

$$= 119.75 ounces$$

How do we find a standard error for $\hat{y}$ at these particular values of the explanatory variables?

This standard error is complicated because $\widehat{bweight}$ is a function of our $\hat{\beta}$'s which are all random variables.

To avoid this computation, we want to transform our data.

# Confidence intervals for predictions: for a specific average

Recall that we have the following regression in mind

$$bweight = \beta_0 + \beta_1 lfaminc + \beta_2 meduc + \beta_3 parity + u$$

Then

$$\hat{\beta}_0 = \hat{E}(bweight|lfaminc = 0, meduc = 0, parity = 0)$$

# Confidence intervals for predictions: for a specific average

If we modify the regression by subtracting our particular values from the independent variables, we get

$$bweight = \beta_0 + \beta_1(lfaminc - 2.674) + \beta_2(meduc - 12) + \beta_3(parity - 3) + u$$

Then

$$\hat{\beta}_0 = \hat{E}(bweight | lfaminc = 2.674, meduc = 12, parity = 3).$$

The new intercept is the predicted birthweight for babies with the particular values we are interested in.

If we run this regression in R, we can then grab the standard errors for the intercept.

# Confidence intervals for predictions: for a specific average

So step by step we need to:

1) Generate new variables: $\tilde{x}_j = x_j - \alpha_j$

2) Run the regression: $y = \tilde{\beta}_0 + \tilde{\beta}_1 \tilde{x}_1 + ... + \tilde{\beta}_k \tilde{x}_k + \tilde{u}$

3) Then $\hat{E}[y|x_1 = \alpha_1, ..., x_k = \alpha_k] = \tilde{\beta}_0$

4) Plug these values into the formula for confidence intervals and interpret.

## Confidence intervals for predictions: for a specific average

```
#Step 1: generate new variables
bwght$lfaminc_0<-bwght$lfaminc-2.674
bwght$motheduc_0<-bwght$motheduc-12
bwght$parity_0<-bwght$parity-3

#step 2: run the new regression
reg2<-lm(bwght~lfaminc_0+motheduc_0+parity_0,bwght)

summary(reg2)
```

```
##
## Call:
## lm(formula = bwght ~ lfaminc_0 + motheduc_0 + parity_0, data = bwght)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -94.533 -11.888   0.779  13.136 151.477
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 119.6491     1.0066 118.864  < 2e-16 ***
## lfaminc_0     2.1313     0.6506   3.276  0.00108 **
## motheduc_0    0.3172     0.2520   1.259  0.20829
## parity_0      1.5261     0.6119   2.494  0.01275 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 20.21 on 1383 degrees of freedom
##   (1 observation deleted due to missingness)
## Multiple R-squared:  0.01633,	Adjusted R-squared:  0.0142
## F-statistic: 7.654 on 3 and 1383 DF,  p-value: 4.482e-05
```

# Confidence intervals for predictions: for a specific average

The 95% confidence interval for the average birthweight for babies given a family income of $14,500, a mother with 12 years of education and with 2 older siblings is:

$$[119.64 - 1.96(1.007), 119.64 + 1.96(1.007)] = [117.6653, 121.6158]$$

A confidence interval for a particular average is not the same as a confidence interval for a particular individual.

For individual observations, we must account for the variance in the unobserved error, $u_i$, which measures our ignorance of the unobserved factors that affect $y_i$.

# Confidence Interval for prediction: a specific unit

We want a confidence interval for $bweight_{i=1}$, the birthweight of baby $i = 1$, with

$$bweight_{i=1} = \beta_0 + \beta_1 lfaminc_{i=1} + \beta_2 meduc_{i=1} + \beta_3 parity_{i=1} + u_{i=1}$$

Our best prediction of $bweight_{i=1}$ is $\widehat{bwieght}_{i=1}$ where

$$\widehat{bweight}_{i=1} = \hat{\beta}_0 + \hat{\beta}_1 lfaminc_{i=1} + \hat{\beta}_2 meduc_{i=1} + \hat{\beta}_3 parity_{i=1}$$

# Confidence Interval for prediction: a specific unit

There is some error, $\hat{u}_{i=1}$, associated with using $\widehat{bweight}_{i=1}$ to predict $bweight_{i=1}$ where

$$
\begin{aligned}
\hat{u}_{i=1} &= bweight_{i=1} - \widehat{bweight}_{i=1} \\
&= (\beta_0 + \beta_1 lfaminc_{i=1} + \beta_2 meduc_{i=1} + \beta_3 parity_{i=1} + u_{i=1}) \\
&\quad - (\hat{\beta}_0 + \hat{\beta}_1 lfaminc_{i=1} + \hat{\beta}_2 meduc_{i=1} + \hat{\beta}_3 parity_{i=1})
\end{aligned}
$$

Finding the expected value, we get:

$$
\begin{aligned}
E[\hat{u}_{i=1}] &= E[bweight_{i=1} - \widehat{bweight}_{i=1}] \\
&= (\beta_0 + \beta_1 lfaminc_{i=1} + \beta_2 meduc_{i=1} + \beta_3 parity_{i=1} + E[u_{i=1}]) \\
&\quad - (E[\hat{\beta}_0] + E[\hat{\beta}_1] lfaminc_{i=1} + E[\hat{\beta}_2] meduc_{i=1} + E[\hat{\beta}_3] parity_{i=1}) \\
&= 0
\end{aligned}
$$

# Confidence Interval for prediction: a specific unit

Finding the variance we get

$$Var(\hat{u}_{i=1}) = Var(bweight_{i=1} - \widehat{bweight}_{i=1})$$

$$= Var(\beta_0 + \beta_1 lfaminc_{i=1} + \beta_2 meduc_{i=1} + \beta_3 parity_{i=1} + u_{i=1} - \widehat{bweight}_{i=1})$$

$$= Var(\widehat{bweight}_{i=1}) + Var(u_{i=1})$$

$$= Var(\widehat{bweight}_{i=1}) + \sigma^2$$

$$\widehat{Var}(\hat{u}_{i=1}) = Var(\widehat{bweight}_{i=1}) + \hat{\sigma}^2$$

# Confidence Interval for prediction: a specific unit

There are two sources of variation in $\hat{u}_{i=1}$.

- the sampling error in $\widehat{bweight}_{i=1}$ which arises because we have estimated the population parameters $\beta$.

- the variance of the error in the population ($u_{i=1}$).

We can compute:

- $Var(\widehat{bweight}_{i=1})$ exactly the way we did before.

- $\hat{\sigma}^2$ from our regression output.

The 95% confidence interval for $bweight_{i=1}$ is then

$$\hat{y} \pm 1.96 * se(\hat{u}_{i=1})$$

# Confidence Interval for prediction: a specific unit

Steps in computing a confidence interval for a particular $y$ when $x_j = \alpha_j$:

1) Generate new variables: $\tilde{x}_j = x_j - \alpha_j$

2) Run the regression: $y = \tilde{\beta}_0 + \tilde{\beta}_1 \tilde{x}_1 + ... + \tilde{\beta}_k \tilde{x}_k + \tilde{u}$

3) Then $\hat{E}[y | x_1 = \alpha_1, ..., x_k = \alpha_k] = \tilde{\beta}_0$ and the standard error of the estimate is $se(\tilde{\beta}_0)$

4) Get an estimate for the variance of $\hat{u} = \hat{\sigma}^2$ from the R output

5) compute the standard error: $\sqrt{se(\tilde{\beta}_0)^2 + \hat{\sigma}^2}$

6) Plug these values into the formula for confidence intervals and interpret.

# Confidence Interval for prediction: a specific unit

```
#Step 1: generate new variables
bwght$lfaminc_0<-bwght$lfaminc-2.674
bwght$motheduc_0<-bwght$motheduc-12
bwght$parity_0<-bwght$parity-3

#step 2: run the new regression
reg2<-lm(bwght~lfaminc_0+motheduc_0+parity_0,bwght)
summary(reg2)
```

```
##
## Call:
## lm(formula = bwght ~ lfaminc_0 + motheduc_0 + parity_0, data = bwght)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -94.533 -11.888   0.779  13.136 151.477
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 119.6491     1.0066 118.864  < 2e-16 ***
## lfaminc_0     2.1313     0.6506   3.276  0.00108 **
## motheduc_0    0.3172     0.2520   1.259  0.20829
## parity_0      1.5261     0.6119   2.494  0.01275 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 20.21 on 1383 degrees of freedom
##   (1 observation deleted due to missingness)
## Multiple R-squared:  0.01633,    Adjusted R-squared:  0.0142
## F-statistic: 7.654 on 3 and 1383 DF,  p-value: 4.482e-05
```

# Confidence Interval for prediction: a specific unit

```
#step 4: get the estimate of the variance
summary(lm(bwght~lfaminc_0+motheduc_0+parity_0,bwght))$sigma^2
```

## [1] 408.5987

The 95% confidence interval for a particular baby's birthweight with:

- ▶ family income of $14,500 (ln(14.5=2.674)),

- ▶ a mother with 12 years of education

- ▶ 2 older siblings is:

$$SE = \sqrt{se(\tilde{\beta}_0)^2 + \hat{\sigma}^2} = \sqrt{(1.007^2) + 408.59} = 20.239$$
$$CI = [119.64 - 1.96 * (20.239); 119.64 + 1.96 * (20.239)]$$
$$= [79.972; 159.308]$$