

# Econometrics II TA Session #8

Hiroki Kato

## 1 Empirical Application of Panel Data Model: Earnings Equation

### 1.1 Background

A researcher wants to estimate the effect of full-time work experience on wages. He uses a *balanced* panel of 595 individuals from 1976 to 1982, taken from the Panel Study of Income Dynamics (PSID). The *balanced* panel data means that we can observe all individuals every year.

```
dt <- read.csv("./data/wages.csv")
head(dt, 14)
```

##	exp	wks	bluecol	ind	south	smsa	married	sex	union	ed	black	lwage	id	time
## 1	3	32	no	0	yes	no	yes	male	no	9	no	5.56068	1	1
## 2	4	43	no	0	yes	no	yes	male	no	9	no	5.72031	1	2
## 3	5	40	no	0	yes	no	yes	male	no	9	no	5.99645	1	3
## 4	6	39	no	0	yes	no	yes	male	no	9	no	5.99645	1	4
## 5	7	42	no	1	yes	no	yes	male	no	9	no	6.06146	1	5
## 6	8	35	no	1	yes	no	yes	male	no	9	no	6.17379	1	6
## 7	9	32	no	1	yes	no	yes	male	no	9	no	6.24417	1	7
## 8	30	34	yes	0	no	no	yes	male	no	11	no	6.16331	2	1
## 9	31	27	yes	0	no	no	yes	male	no	11	no	6.21461	2	2
## 10	32	33	yes	1	no	no	yes	male	yes	11	no	6.26340	2	3
## 11	33	30	yes	1	no	no	yes	male	no	11	no	6.54391	2	4
## 12	34	30	yes	1	no	no	yes	male	no	11	no	6.69703	2	5
## 13	35	37	yes	1	no	no	yes	male	no	11	no	6.79122	2	6
## 14	36	30	yes	1	no	no	yes	male	no	11	no	6.81564	2	7

The variable `id` and `time` indicate individual and time indexes. We use these two variables to apply panel data models. Additionally, we use the following variables:

- `exp`: years of full-time work experience
- `sqexp`: squared value of `exp`
- `sex`: a dummy variable taking 1 if an individual is female
- `ed`: years of education

- `lwage`: logarithm of wage

```
dt <- dt[,c("id", "time", "exp", "lwage")]
dt$sqexp <- dt$exp^2
summary(dt)
```

```
##          id          time          exp          lwage          sqexp
##  Min.    : 1    Min.    :1    Min.    : 1.00    Min.    :4.605    Min.    : 1.0
##  1st Qu.:149    1st Qu.:2    1st Qu.:11.00    1st Qu.:6.395    1st Qu.: 121.0
##  Median :298    Median :4    Median :18.00    Median :6.685    Median : 324.0
##  Mean   :298    Mean   :4    Mean   :19.85    Mean   :6.676    Mean   : 514.4
##  3rd Qu.:447    3rd Qu.:6    3rd Qu.:29.00    3rd Qu.:6.953    3rd Qu.: 841.0
##  Max.   :595    Max.   :7    Max.   :51.00    Max.   :8.537    Max.   :2601.0
```

## 1.2 Pooled OLS

Using the OLS method, we want to estimate the following linear panel data model:

$$\text{lwage}_{it} = \alpha + \beta_1 \cdot \text{exp}_{it} + \beta_2 \cdot \text{sqexp}_{it} + \beta_3 \cdot \text{sex}_{it} + \beta_4 \cdot \text{ed}_{it} + u_{it}.$$

We will discuss assumptions for applying the OLS method. Let  $\mathbf{X}_{it}$  be a  $1 \times K$  (stochastic) explanatory vector. This vector contains `exp`, `sqexp`, `sex` and `ed`. Let  $Y_{it}$  be a random variable of outcome, that is `lwage`. The balanced panel data is given by

	$i = 1$	$i = 2$	...	$i = n$
$t = 1$	$(Y_{11}, \mathbf{X}_{11})$	$(Y_{21}, \mathbf{X}_{21})$	...	$(Y_{n1}, \mathbf{X}_{n1})$
$t = 2$	$(Y_{12}, \mathbf{X}_{12})$	$(Y_{22}, \mathbf{X}_{22})$	...	$(Y_{n2}, \mathbf{X}_{n2})$
$\vdots$	$\vdots$	$\vdots$	...	$\vdots$
$t = T$	$(Y_{1T}, \mathbf{X}_{1T})$	$(Y_{2T}, \mathbf{X}_{2T})$	...	$(Y_{nT}, \mathbf{X}_{nT})$

Then, the linear panel data model can be rewritten as follows:

$$Y_{it} = \mathbf{X}_{it}\beta + u_{it}, \quad t = 1, \dots, T, \quad i = 1, \dots, n.$$

Using notations  $\underline{\mathbf{X}}_i = (\mathbf{X}'_{i1}, \dots, \mathbf{X}'_{iT})'$  and  $\underline{Y}_i = (Y_{i1}, \dots, Y_{iT})'$ , and  $\underline{u}_i = (u_{i1}, \dots, u_{iT})'$ , we can reformulate this model as follows:

$$\underline{Y}_i = \underline{\mathbf{X}}_i\beta + \underline{u}_i, \quad \forall i.$$

Now, we assume

1.  $E[\mathbf{X}'_{it}u_{it}] = 0, \forall i, t$ . This assumption, called (*contempraneous*) *exogeneity assumption*, implies that  $u_{it}$  and  $\mathbf{X}_{it}$  are orthogonal in the conditional mean sence,  $E[u_{it}|\mathbf{X}_{it}] = 0$ .

However, this assumption does not imply  $u_{it}$  is uncorrelated with the explanatory variables in all time periods (strictly exogeneity), that is,  $E[u_{it}|\mathbf{X}_{i1}, \dots, \mathbf{X}_{iT}] = 0$ . This assumption places no restriction on the relationship between  $\mathbf{X}_{is}$  and  $u_{it}$  for  $s \neq t$ .

2.  $E[\mathbf{X}_i' \mathbf{X}_i] \succ 0$ .

Under these two assumptions, the true parameter can be identified by

$$\beta = E[\mathbf{X}_i' \mathbf{X}_i]^{-1} E[\mathbf{X}_i' Y_i].$$

Hence, the OLSE (pooled OLSE) is given by

$$\hat{\beta} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i' \mathbf{X}_i \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i' Y_i \right) = \left( \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^T \mathbf{X}_{it}' \mathbf{X}_{it} \right) \left( \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^T \mathbf{X}_{it}' Y_{it} \right).$$

The pooled OLS estimator is consistent and asymptotically normally distributed.

$$\sqrt{n}(\hat{\beta} - \beta) \sim N(0, A^{-1} B A^{-1}),$$

where  $A = E[\mathbf{X}_i' \mathbf{X}_i]$  and  $B = E[\mathbf{X}_i' u_i u_i' \mathbf{X}_i]$ . The consistent estimator of the asymptotic variance covariance matrix is given by

$$\hat{A}^{-1} \hat{B} \hat{A}^{-1} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i' \mathbf{X}_i \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i' u_i u_i' \mathbf{X}_i \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i' \mathbf{X}_i \right)^{-1}$$

The standard errors calculated by this matrix is called *robust standard errors clustered by individuals*.

In R, the pooled OLSE can be obtained by `lm` function. However, the `lm` function does not return the cluster-robust standard errors. Thus, you need to calculate them by yourself. Here is a sample code.

```
# OLSE
pool <- lm(lwage ~ -1 + exp + sqexp, data = dt)

# Clustered SE
X <- model.matrix(pool); uhat <- pool$residuals
uhatset <- matrix(0, nrow = nrow(X), ncol = nrow(X))

i_from <- 1; j_from <- 1
for (i in 1:max(dt$id)) {
  x <- as.numeric(rownames(dt))[dt$id == i]
  usq <- uhat[x] %*% t(uhat[x])
  i_to <- i_from + nrow(usq) - 1
  j_to <- j_from + ncol(usq) - 1
  uhatset[i_from:i_to, j_from:j_to] <- usq
  i_from <- i_to + 1; j_from <- j_to + 1
}
```

```

}

Ahat <- t(X) %*% X
Bhat <- t(X) %*% uhatset %*% X
clust_vcov <- solve(Ahat) %*% Bhat %*% solve(Ahat)
clust_se <- sqrt(diag(clust_vcov))

print("Pooled OLSE"); coef(pool)

```

```
## [1] "Pooled OLSE"
```

```
##          exp          sqexp
## 0.64570881 -0.01279755
```

```
print("SE of pooled OLSE"); clust_se
```

```
## [1] "SE of pooled OLSE"
```

```
##          exp          sqexp
## 0.0107859273 0.0003765058
```

Alternatively, using the `plm` function (the package `plm`) and the `coeftest` function (the package `lmtest`), you can obtain the asymptotic variance covariance matrix of pooled OLS easily. The `plm` function provides the panel data model. When you want to estimate pooled OLS, you need to specify `model = "pooling"`. Moreover, you should specify individual and time index using `index` argument. This argument passes `index = c("individual index", "time index")`. After estimating the pooled OLS by the `plm` function, you must use the `coeftest` function to obtain the cluster-robust standard errors. To calculate the clustered standard errors, you should use the `vcovHC` function in the `vcov` argument.

```

library(plm)
library(lmtest)
library(sandwich)
test <- plm(lwage ~ -1 + exp + sqexp, data = dt, model = "pooling", index = c("id", "time"))
coeftest(test, vcov = vcovHC(test, type = "HC0", cluster = "group"))

```

```
##
```

```
## t test of coefficients:
```

```
##
```

```
##          Estimate Std. Error t value Pr(>|t|)
## exp      0.64570881 0.01078593  59.866 < 2.2e-16 ***
## sqexp    -0.01279755 0.00037651 -33.990 < 2.2e-16 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# OLS
```

```
pool <- lm(lwage ~ -1 + exp + sqexp, data = dt)
uhat <- pool$residuals
```

```

omega_sum <- matrix(0, ncol = max(dt$time), nrow = max(dt$time))
for (i in 1:max(dt$id)) {
  x <- as.numeric(rownames(dt))[dt$id == i]
  omega_sum <- uhat[x] %*% t(uhat[x]) + omega_sum
}
omega <- omega_sum/max(dt$id)

# FGLS
X <- model.matrix(pool)
Y <- dt$lwage
Iomega <- diag(max(dt$id)) %x% solve(omega)
bfgls <- solve(t(X) %*% Iomega %*% X) %*% (t(X) %*% Iomega %*% Y)

# vcov of FGLS
ufgls <- Y - X %*% bfgls
uhatset <- matrix(0, nrow = nrow(X), ncol = nrow(X))
i_from <- 1; j_from <- 1
for (i in 1:max(dt$id)) {
  x <- as.numeric(rownames(dt))[dt$id == i]
  usq <- uhat[x] %*% t(uhat[x])
  i_to <- i_from + nrow(usq) - 1
  j_to <- j_from + ncol(usq) - 1
  uhatset[i_from:i_to, j_from:j_to] <- usq
  i_from <- i_to + 1; j_from <- j_to + 1
}

Ahat <- t(X) %*% Iomega %*% X
Bhat <- t(X) %*% Iomega %*% uhatset %*% Iomega %*% X
vcov_fgls <- solve(Ahat) %*% Bhat %*% solve(Ahat)
se_fgls <- sqrt(diag(vcov_fgls))

# estimate
i <- rep(1, max(dt$time))
Qt <- diag(max(dt$time)) - i %*% solve(t(i) %*% i) %*% t(i)
Ybar <- diag(max(dt$id)) %x% Qt %*% Y
Xbar <- diag(max(dt$id)) %x% Qt %*% X
bfe <- solve(t(Xbar) %*% Xbar) %*% t(Xbar) %*% Ybar

# inference
uhat <- Ybar - Xbar %*% bfe
sigmahat <- sum(uhat^2)/(max(dt$id)*(max(dt$time)-1)-2)
vcovfe <- sigmahat * solve(t(Xbar) %*% Xbar)
sefe <- sqrt(diag(vcovfe))

```

```

library(plm)
summary(plm(lwage ~ 1 + exp + sqexp, data = dt, index = c("id", "time"), model = "within",

## Oneway (individual) effect Within Model
##
## Call:
## plm(formula = lwage ~ 1 + exp + sqexp, data = dt, model = "within",
##      index = c("id", "time"))
##
## Balanced Panel: n = 595, T = 7, N = 4165
##
## Residuals:
##      Min.      1st Qu.      Median      3rd Qu.      Max.
## -1.8119015 -0.0506647  0.0041017  0.0607943  1.9430281
##
## Coefficients:
##      Estimate Std. Error t-value Pr(>|t|)
## exp      0.11398290  0.00246524  46.236 < 2.2e-16 ***
## sqexp    -0.00042940  0.00005452  -7.876 4.452e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Total Sum of Squares:      240.65
## Residual Sum of Squares: 82.677
## R-Squared:      0.65644
## Adj. R-Squared: 0.59906
## F-statistic: 3408.73 on 2 and 3568 DF, p-value: < 2.22e-16

# pooled OLS and estimator of sigma
n <- max(dt$id); t <- max(dt$time)
vhat <- lm(lwage ~ -1 + exp + sqexp, data = dt)$residuals
sigmav <- sum(vhat^2)/(n*t - 2)
vdt <- data.frame(vhat = vhat, id = dt$id, time = dt$time)
vdt$time1 <- ifelse(vdt$time > 1, 1, 0)
vdt$time2 <- ifelse(vdt$time > 2, 1, 0)
vdt$time3 <- ifelse(vdt$time > 3, 1, 0)
vdt$time4 <- ifelse(vdt$time > 4, 1, 0)
vdt$time5 <- ifelse(vdt$time > 5, 1, 0)
vdt$time6 <- ifelse(vdt$time > 6, 1, 0)

library(tidyverse)
for (i in 1:n) {
  vdt <- vdt %>%
    mutate(
      dm_u = case_when(

```

```

        id == i & time == 1 ~ vhat * time1,
        id == i & time == 2 ~ vhat * time2,
        id == i & time == 3 ~ vhat * time3,
        id == i & time == 4 ~ vhat * time4,
        id == i & time == 5 ~ vhat * time5,
        id == i & time == 6 ~ vhat * time6
    )
)
}

library(plm)
summary(plm(lwage ~ -1 + exp + sqexp, data = dt, index = c("id", "time"), model = "random",

## Oneway (individual) effect Random Effect Model
##   (Swamy-Arora's transformation)
##
## Call:
## plm(formula = lwage ~ -1 + exp + sqexp, data = dt, model = "random",
##     index = c("id", "time"))
##
## Balanced Panel: n = 595, T = 7, N = 4165
##
## Effects:
##               var std.dev share
## idiosyncratic 0.02317 0.15222 0.009
## individual    2.62039 1.61876 0.991
## theta: 0.9645
##
## Residuals:
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -1.7247  0.0639  0.1400  0.1422  0.2258  2.1452
##
## Coefficients:
##               Estimate Std. Error z-value Pr(>|z|)
## exp      1.6583e-01  3.2167e-03  51.552 < 2.2e-16 ***
## sqexp -1.2025e-03  7.3838e-05 -16.285 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Total Sum of Squares:    241.47
## Residual Sum of Squares: 186.73
## R-Squared:    0.62742
## Adj. R-Squared: 0.62733
## Chisq: 6442.14 on 2 DF, p-value: < 2.22e-16

```