

**Problem Set 1 Solutions**  
**EE426 Semester 2/2014**

**1. (Difference-in-differences)**

1.1 If building the incinerator reduces the value of homes closer to the site, the sign of  $\delta_1$  should be positive, that is the longer distance of homes from the incinerator site (that was built in 1981), the higher value of the home. >> **the difference in elasticity of price of house with respect to the distance to the incinerator site between 1981 and 1978**

$\beta_1$  measures the location effect that is not due to the presence of the incinerator. ( $\delta_1$  explains the location effect that is due to the presence of the incinerator since the incinerator was built in 1981.) If  $\beta_1 > 0$ , on average with or without the incinerator, the prices of the house that is located far away from the supposed incinerator site increase with the distance from the supposed incinerator site.

1.2 The results are reported in column (1). We see that, with the presence of the incinerator in 1981, there is a positively relationship between the distance and the housing prices. If the house is located one percent far away from the incinerator, the housing price increases 0.048 percent in year 1981 higher than the increase in year 1978. However, the coefficient is not statistically significant even at a 10% significance level.

	(1) b/se	(2) b/se
y81	-0.011 (0.8051)	-0.225 (0.4947)
ldist	0.317*** (0.0515)	0.001 (0.0446)
y81ldist	0.048 (0.0818)	0.062 (0.0503)
age		-0.008*** (0.0014)
agesq		0.000*** (0.0000)
rooms		0.046*** (0.0173)
baths		0.101*** (0.0278)
lintst		-0.060* (0.0317)
lland		0.095*** (0.0247)
larea		0.351*** (0.0519)
_cons	8.058*** (0.5084)	7.674*** (0.5016)
r2	0.396	0.787
N	321	321

```

use KIELMC.dta

reg lprice y81 ldist y81ldist
est sto reg1

reg lprice y81 ldist y81ldist age agesq rooms baths lintst lland larea
est sto reg2

estout reg1 reg2, cells(b(star fmt(3)) se(par fmt(4))) stats(r2 N, fmt(3
0)) starlevels(* 0.10 ** 0.05 *** 0.01)

```

1.3 Column (2) reports the results when adding other control variables. When we control for other factors variables, the coefficient of  $y81 \cdot \log(\text{dist})$  becomes 0.062. With other things being equal, one percent of the house location away from the incinerator increases the housing price in 1981 for 0.062 percent higher than that in 1978. Still, the coefficient is not statistically significant. However, the main explanations on the housing price are from other control variables as they are all statistically significant. Hence, model (1) seems to omit some important variables.

## 2. Pooled OLS vs. First Differences

### 2.1 Pooled OLS:

$$\log(\text{wage})_{it} = \beta_0 + \beta_1 \text{educ}_{it} + \beta_2 \text{black}_{it} + \beta_3 \text{hisp}_{it} + \beta_4 \text{exper}_{it} + \beta_5 \text{married}_{it} + \beta_6 \text{union}_{it} + \gamma_1 d81 + \gamma_2 d82 + \gamma_3 d83 + \gamma_4 d84 + \gamma_5 d85 + \gamma_6 d86 + \gamma_7 d87 + u_{it}$$

STATA command:

```

reg lwage educ black hisp exper married union d81- d87
est sto pool1

```

If they are married, wage will increase by 11.1 percent on average, compared to those who are not married with other controls being equal. Additionally, if they have a membership in a union, their wage will be 18.6 percent higher than those who are not in the union, ceteris paribus. These coefficients are all statistically significant at 1% significance level.

### 2.2-2.3 Do the first difference:

$$\text{Equation: } \Delta \log(\text{wage})_{it} = \beta_0 + \beta_5 \Delta \text{married}_{it} + \beta_6 \Delta \text{union}_{it} + \delta_2 d82 + \delta_3 d83 + \delta_4 d84 + \delta_5 d85 + \delta_6 d86 + \delta_7 d87 + \Delta u_{it}$$

```

sort nr year
by nr: gen l_lwage = lwage[_n-1]
by nr: gen l_educ = educ[_n-1]
by nr: gen l_black = black[_n-1]
by nr: gen l_hisp = hisp[_n-1]
by nr: gen l_exper = exper[_n-1]
by nr: gen l_married = married[_n-1]
by nr: gen l_union = union[_n-1]

```

```

gen d_lwage = lwage - l_lwage
gen d_educ = educ - l_educ
gen d_black = black - l_black
gen d_hisp = hisp - l_hisp
gen d_exper = exper - l_exper
gen d_married = married - l_married
gen d_union = union - l_union

```

```

reg d_lwage d_educ d_black d_hisp d_exper d_married d_union d82- d87
est sto fdiff1
note: d_educ omitted because of collinearity
note: d_black omitted because of collinearity
note: d_hisp omitted because of collinearity
note: d_exper omitted because of collinearity

```

You can see that these time-invariant variables are all dropped out of the estimation because of collinearity after doing the first difference. You can browse to see how these variables look like:

```

browse year nr d_educ d_black d_hisp d_exper
>> They are all constant across time.

```

	(1) b/se	(2) b/se
educ	0.093*** (0.0052)	
black	-0.137*** (0.0236)	
hisp	0.014 (0.0208)	
exper	0.030*** (0.0055)	
married	0.111*** (0.0157)	
union	0.186*** (0.0171)	
d_married		0.040* (0.0229)
d_union		0.042** (0.0197)
d81	0.078*** (0.0296)	
d82	0.097*** (0.0311)	-0.060** (0.0269)
d83	0.107*** (0.0334)	-0.071*** (0.0269)
d84	0.140*** (0.0363)	-0.047* (0.0269)
d85	0.160*** (0.0398)	-0.067** (0.0269)
d86	0.188*** (0.0437)	-0.056** (0.0269)
d87	0.211*** (0.0478)	-0.052* (0.0269)
_cons	0.156** (0.0753)	0.115*** (0.0191)
r2	0.188	0.005
N	4360	3815

For your reference, what if we do the changes in year dummy variables.

```
*using loop
forvalues y = 1(1)7 {
  by nr: gen l_d8`y' = d8`y'[_n-1]
  gen d_d8`y' = d8`y' - l_d8`y'
}
```

```
browse d_d8*
```

\*Take a look on how it looks like when we do changes in year dummy:  $\Delta d_{81}, \Delta d_{82}, \dots, \Delta d_{87}$

When using the changes in year dummy, don't forget that we don't have a constant. (Compare to our lecture note in class on "Differencing with more than 2 periods")

```
reg d_lwage d_married d_union d_d81 d_d82 d_d83 d_d84 d_d85 d_d86 d_d87,
nocons
est sto fdiff2
```

```
estout fdiff1 fdiff2, cells(b(star fmt(3)) se(par fmt(4))) stats(r2 N,
fmt(3 0)) starlevels(* 0.10 ** 0.05 *** 0.01)
```

	fdiff1 b/se	fdiff2 b/se
d_married	0.040* (0.0229)	0.040* (0.0229)
d_union	0.042** (0.0197)	0.042** (0.0197)
d82	-0.060** (0.0269)	
d83	-0.071*** (0.0269)	
d84	-0.047* (0.0269)	
d85	-0.067** (0.0269)	
d86	-0.056** (0.0269)	
d87	-0.052* (0.0269)	
d_d81		0.115*** (0.0191)
d_d82		0.171*** (0.0271)
d_d83		0.215*** (0.0334)
d_d84		0.284*** (0.0387)
d_d85		0.333*** (0.0432)
d_d86		0.392*** (0.0474)
d_d87		0.455*** (0.0512)
_cons	0.115*** (0.0191)	

r2	0.005	0.027
N	3815	3815

If you predict the fitted values from these two models (fdiff1 and fdiff2), you will see that the fitted values are exactly the same.

```
reg d_lwage d_married d_union d82- d87
predict y1
```

```
reg d_lwage d_married d_union d_d81 d_d82 d_d83 d_d84 d_d85 d_d86 d_d87,
nocons
predict y2
```

```
sum y1 y2
```

Variable	Obs	Mean	Std. Dev.	Min	Max
y1	3815	.0675718	.0300554	-.0335739	.1975348
y2	3815	.0675718	.0300554	-.0335739	.1975348

2.4 For the first difference, if they are married, wage will increase by 4.03 percent, compared to those who are not married with other controls being equal. Additionally, if they have a membership in a union, their wage will be 4.2 percent higher than those who are not in the union, ceteris paribus. (The interpretation on coefficients for FD is still the same as we use level. See Example 13.5 and Section 13.4)

The first difference estimation provides smaller estimates compared to the levels estimates. We may say that without taking into account individual unobserved effects (that are constant across time), Pooled OLS overestimates the premium on marriage and the effect of union membership on log(wage).

### 3. FE, RE, and pooled OLS

#### 3.1-3.2 Pooled OLS vs. RE

##### \*3.1 Pooled OLS

```
reg lwage educ black hisp
est sto ols1
```

##### \*3.2 RE

```
xtset nr year
```

```
xtreg lwage educ black hisp, re
est sto re1
```

```
estout ols1 re1, cells(b(star fmt(3)) se(par fmt(4))) stats(r2 N, fmt(3
0)) starlevels(* 0.10 ** 0.05 *** 0.01)
```

	ols1 b/se	re1 b/se
educ	0.077*** (0.0046)	0.077*** (0.0092)
black	-0.123*** (0.0247)	-0.123** (0.0497)
hisp	0.025 (0.0222)	0.025 (0.0447)
_cons	0.752*** (0.0554)	0.752*** (0.1114)
r2	0.070	
N	4360	4360

The RE and pooled OLS estimates of  $\beta_j$  are all the same.

3.3 The RE standard errors are larger than the pooled OLS standard errors. If you do BP test (xttest0), it rejects the null hypothesis that  $\text{var}(a_i) = 0$ . So, it means that we have significant heterogeneous (unobserved) effects  $a_i$ , and that transformation on the data (GLS for RE) is preferred.

xttest0

Breusch and Pagan Lagrangian multiplier test for random effects

$$\text{lwage}[nr,t] = Xb + u[nr] + e[nr,t]$$

Estimated results:

	Var	sd = sqrt(Var)
lwage	.2836728	.5326094
e	.1499484	.3872317
u	.1148867	.3389494

Test:  $\text{Var}(u) = 0$

chibar2(01) = 2843.93  
Prob > chibar2 = 0.0000

3.4-3.5 Add a full set of year dummies, do pooled OLS, RE, and FE.

\*3.4

```
reg lwage educ black hisp d81- d87
est sto ols2
```

```
xtreg lwage educ black hisp d81- d87, re
est sto re2
```

\*3.5

```
xtreg lwage educ black hisp d81- d87, fe
est sto fe
```

```
estout ols2 re2 fe, cells(b(star fmt(3)) se(par fmt(4))) stats(r2 N, fmt(3
0)) starlevels(* 0.10 ** 0.05 *** 0.01)
```

	ols2 b/se	re2 b/se	fe b/se
educ	0.077*** (0.0044)	0.077*** (0.0092)	
black	-0.123*** (0.0237)	-0.123** (0.0497)	
hisp	0.025 (0.0213)	0.025 (0.0447)	
d81	0.119*** (0.0299)	0.119*** (0.0215)	0.119*** (0.0215)
d82	0.178*** (0.0299)	0.178*** (0.0215)	0.178*** (0.0215)
d83	0.226*** (0.0299)	0.226*** (0.0215)	0.226*** (0.0215)
d84	0.297*** (0.0299)	0.297*** (0.0215)	0.297*** (0.0215)
d85	0.346*** (0.0299)	0.346*** (0.0215)	0.346*** (0.0215)
d86	0.406*** (0.0299)	0.406*** (0.0215)	0.406*** (0.0215)
d87	0.473*** (0.0299)	0.473*** (0.0215)	0.473*** (0.0215)
_cons	0.497*** (0.0567)	0.497*** (0.1123)	1.393*** (0.0152)
r2	0.145		0.163
N	4360	4360	4360

- The coefficients of pooled OLS and RE on educ, black, and hisp, after adding year dummy, are still the same, but the constant is smaller as all year dummies are statistically significant.

- With FE, it drops all explanatory variables, except year dummies. The FE coefficients on the year dummies are the same as the RE coefficients on the year dummies, but the constant coefficients are largely different.

### 3.6

\* year dummies as intercepts and interactions with educ

```
xtreg lwage d81- d87 c.educ#d81 c.educ#d82 c.educ#d83 c.educ#d84 c.educ#d85
c.educ#d86 c.educ#d87 , fe
```

```
Fixed-effects (within) regression                Number of obs   =   4360
Group variable: nr                             Number of groups =   545

R-sq:  within = 0.1647                          Obs per group:  min =    8
          between = 0.1183                        avg   =   8.0
          overall = 0.1056                        max   =    8

corr(u_i, Xb) = 0.0588                          F(14,3801)     =   53.54
                                                Prob > F       =   0.0000
```

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
d81	-.0283826	.1463373	-0.19	0.846	-.3152897 .2585246
d82	-.0129925	.1463373	-0.09	0.929	-.2998996 .2739146



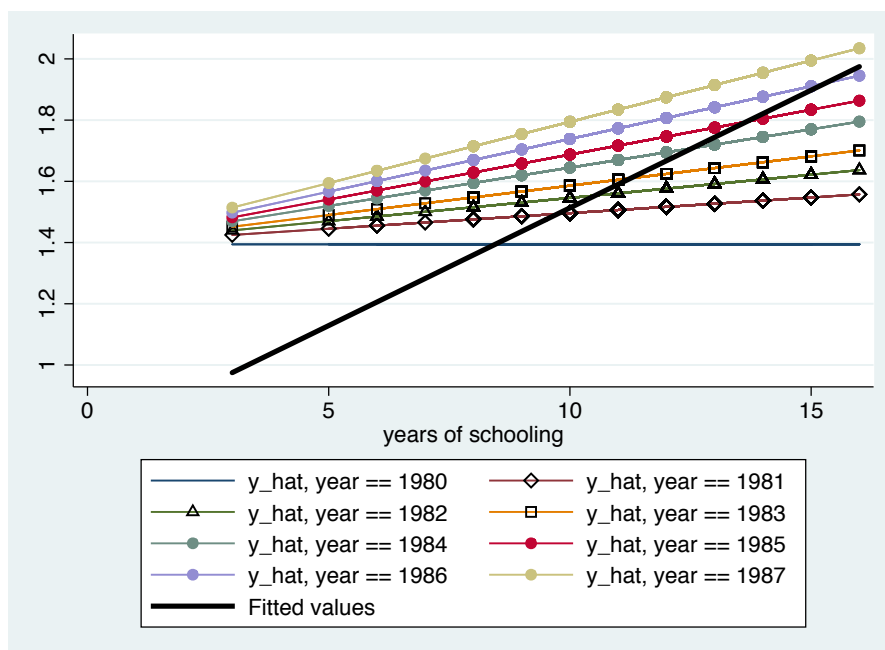
d85#c.educ	1	.0293259	.0018042	16.25	0.000	.0257886	.0328633
d86#c.educ	1	.0344277	.0018042	19.08	0.000	.0308904	.037965
d87#c.educ	1	.040027	.0018042	22.19	0.000	.0364897	.0435644
_cons		1.394139	.0150326	92.74	0.000	1.364666	1.423611
-----							
sigma_u		.37939091					
sigma_e		.35429438					
rho		.53416615	(fraction of variance due to u_i)				
-----							
F test that all u_i=0:		F(544, 3808) =	9.10	Prob > F =		0.0000	

Graph from the latter model (without year dummies, only interactions)

```

predict y_hat
separate y_hat, by(year)
twoway connected y_hat1980- y_hat1987 educ, msymbol(none diamond_hollow
triangle_hollow square_hollow) msize(medium) mcolor(black black black black) ||
lfit lwage educ, clwidth(thick) clcolor(black)

```



With both year dummy variables and their interactions with education, there is no effect of year as intercepts on  $\log(\text{wage})$ . Only from year 1985 onwards that education shows positive relationship with  $\log(\text{wage})$ . With only interactions model, we see increasing association between each year education and  $\log(\text{wage})$ , reflecting the effects of schooling on  $\log(\text{wage})$  has changed over time.

### 3.7 General conclusions?

From BP test, we prefer RE model. Since we have all time-invariant variables, using FE will drop all these variables we want to estimate. Hence, in this case, RE will be our preferred model.

4. Show that  $\text{corr}(\Delta u_{it}, \Delta u_{i,t+1}) = -0.5$  when  $u_{it} \sim \text{iid}N(0, \sigma_u^2)$ , hence we have  $E(u_{it}, u_{is}) = 0$  for any  $t \neq s$

$$\text{corr}(\Delta u_{it}, \Delta u_{i,t+1}) = \frac{E(u_{it} - u_{i,t-1}, u_{i,t+1} - u_{it})}{[E(u_{it} - u_{i,t-1})^2 E(u_{i,t+1} - u_{it})^2]^{1/2}}$$

$$E(u_{it} - u_{i,t-1}, u_{i,t+1} - u_{it}) = E[u_{it}u_{i,t+1} - u_{it}^2 - u_{i,t-1}u_{i,t+1} + u_{i,t+1}u_{it}] = -E[u_{it}^2] = -\sigma_u^2$$

$$E(u_{it} - u_{i,t-1})^2 = E[u_{it}^2 - 2u_{it}u_{i,t-1} + u_{i,t-1}^2] = 2\sigma_u^2$$

$$E(u_{i,t+1} - u_{it})^2 = E[u_{i,t+1}^2 - 2u_{i,t+1}u_{it} + u_{it}^2] = 2\sigma_u^2$$

$$\therefore \text{corr}(\Delta u_{it}, \Delta u_{i,t+1}) = \frac{-\sigma_u^2}{(2\sigma_u^2 \cdot 2\sigma_u^2)^{1/2}} = -\frac{1}{2}$$

5. The composite error  $v_{it} = a_i + u_{it}$ ,  $E[a_i, u_{it}] = 0$ ,  $E[u_{it}^2] = \sigma_u^2$ , and  $E(u_{it}, u_{is}) = 0$  for any  $t \neq s$ . Define  $e_{it} = v_{it} - \lambda \bar{v}_i$ , where  $\lambda = 1 - [\sigma_u^2 / (\sigma_u^2 + T\sigma_a^2)]^{1/2}$ .

- (5.1) Show that  $E(e_{it}) = 0$ .

$$\begin{aligned} E[v_{it} - \lambda \bar{v}_i] &= E[v_{it}] - \lambda E[\bar{v}_i] = E[a_i + u_{it}] - \lambda E[\bar{a}_i + \bar{u}_i] \\ &= E[a_i] + E[u_{it}] - \lambda E[\bar{a}] - \lambda E[\bar{u}] = 0 \end{aligned}$$

- (5.2) Show that  $\text{Var}(e_{it}) = \sigma_u^2$ ,  $t = 1, \dots, T$ .

$$\text{Var}(e_{it}) = E[(v_{it} - \lambda \bar{v}_i)^2] = E(v_{it}^2) - 2\lambda E(v_{it}, \bar{v}_i) + \lambda^2 E(\bar{v}_i^2)$$

$$E(v_{it}^2) = E[a_i^2 + 2a_i u_{it} + u_{it}^2] = \sigma_a^2 + \sigma_u^2$$

$$E(\bar{v}_i^2) = E[\bar{a}_i^2 + 2\bar{a}_i \bar{u}_i + \bar{u}_i^2] = \sigma_a^2 + \frac{1}{T^2} \sum E(u_{it}^2) = \sigma_a^2 + \frac{\sigma_u^2}{T}$$

$$E(v_{it}, \bar{v}_i) = E[a_i + u_{it}, \bar{a}_i + \bar{u}_i] = E[a_i \bar{a}] + E[a_i \bar{u}] + E[u_{it} \bar{a}] + E[u_{it} \bar{u}]$$

$$= \sigma_a^2 + \frac{\sigma_u^2}{T}$$

**Note:** (i)  $a_i$  is the same for every  $t$  within each cross-sectional unit  $i$ , so  $a_i = \bar{a}$ .

$$E[a_i \bar{a}] = \text{Var}(a_i) = \sigma_a^2$$

$$(ii) E[\bar{u}_i^2] = E\left(\frac{\sum u_{it}}{T}\right)^2 = \frac{1}{T^2} \sum (E u_{it}^2) = \frac{T}{T^2} \sigma_u^2 = \frac{\sigma_u^2}{T}$$

$$\begin{aligned} \therefore \text{Var}(e_{it}) &= \sigma_a^2 + \sigma_u^2 - 2\lambda \left[ \sigma_a^2 + \frac{\sigma_u^2}{T} \right] + \lambda^2 \left[ \sigma_a^2 + \frac{\sigma_u^2}{T} \right] \\ &= (1 - 2\lambda + \lambda^2) \sigma_a^2 + (1 - 2\lambda + \lambda^2) \frac{\sigma_u^2}{T} + (T - 1) \frac{\sigma_u^2}{T} \end{aligned}$$

$$\begin{aligned}
&= (1 - \lambda)^2 \left( \sigma_a^2 + \frac{\sigma_u^2}{T} \right) + (T - 1) \frac{\sigma_u^2}{T} \\
&= \left( 1 - 1 + \left( \frac{\sigma_u^2}{\sigma_u^2 + T\sigma_a^2} \right)^{1/2} \right)^2 \left( \frac{T\sigma_a^2 + \sigma_u^2}{T} \right) + \sigma_u^2 - \frac{\sigma_u^2}{T} \\
&= \frac{\sigma_u^2}{T} + \sigma_u^2 - \frac{\sigma_u^2}{T} = \sigma_u^2
\end{aligned}$$

**(5.3)** Show that for  $t \neq s$ ,  $Cov(e_{it}, e_{is}) = 0$ .

$$Cov(e_{it}, e_{is}) = E[(v_{it} - \lambda \bar{v}_i), (v_{is} - \lambda \bar{v}_i)] = E[v_{it}v_{is} - \lambda v_{it}\bar{v}_i - \lambda v_{is}\bar{v}_i + \lambda^2 \bar{v}_i^2]$$

$$E[v_{it}v_{is}] = E[a_i^2 + a_i u_{is} + a_i u_{it} + u_{it}u_{is}] = \sigma_a^2$$

$$\begin{aligned}
\therefore Cov(e_{it}, e_{is}) &= \sigma_a^2 - 2\lambda \left( \sigma_a^2 + \frac{\sigma_u^2}{T} \right) + \lambda^2 \left( \sigma_a^2 + \frac{\sigma_u^2}{T} \right) = (1 - 2\lambda + \lambda^2) \left( \sigma_a^2 + \frac{\sigma_u^2}{T} \right) - \frac{\sigma_u^2}{T} \\
&= \frac{\sigma_u^2}{\sigma_u^2 + T\sigma_a^2} \cdot \frac{T\sigma_a^2 + \sigma_u^2}{T} - \frac{\sigma_u^2}{T} = 0
\end{aligned}$$