

# Four Special Topics

Interaction Terms

Standardized Regression

Decomposing Regression Effects

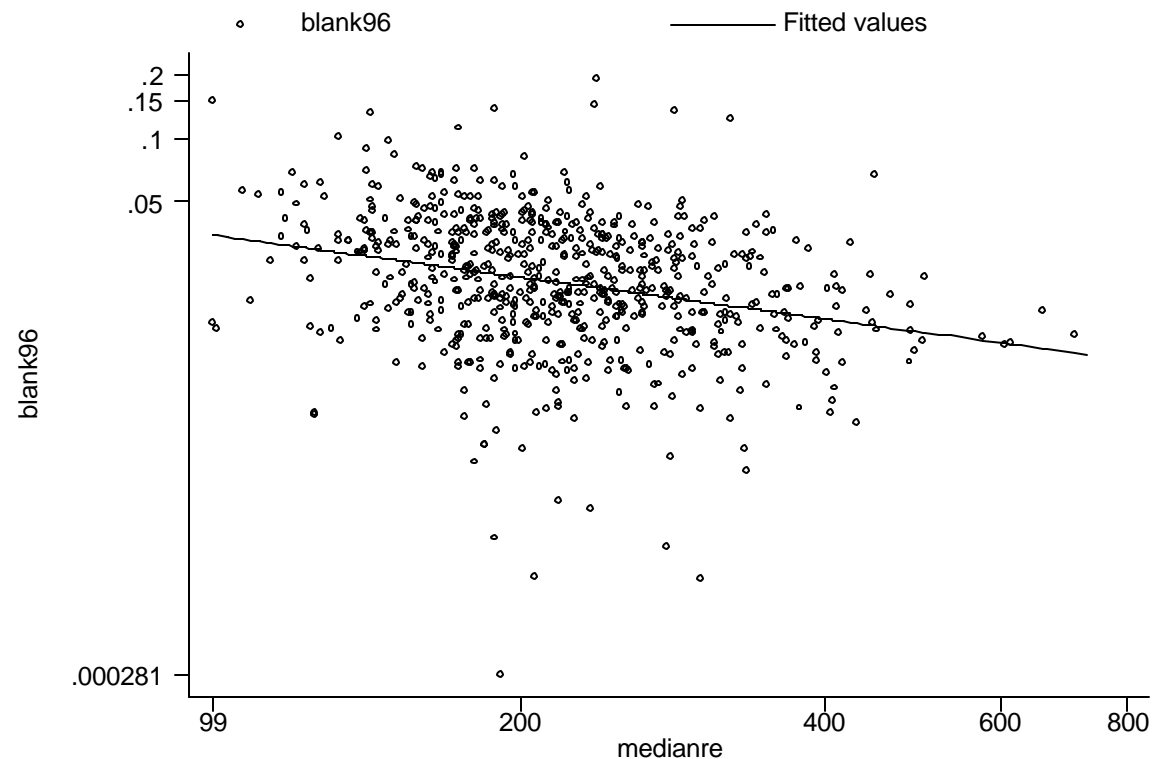
Measurement Error

# Interaction Terms

What Happens When Different  
Models Apply in Different  
Situations?

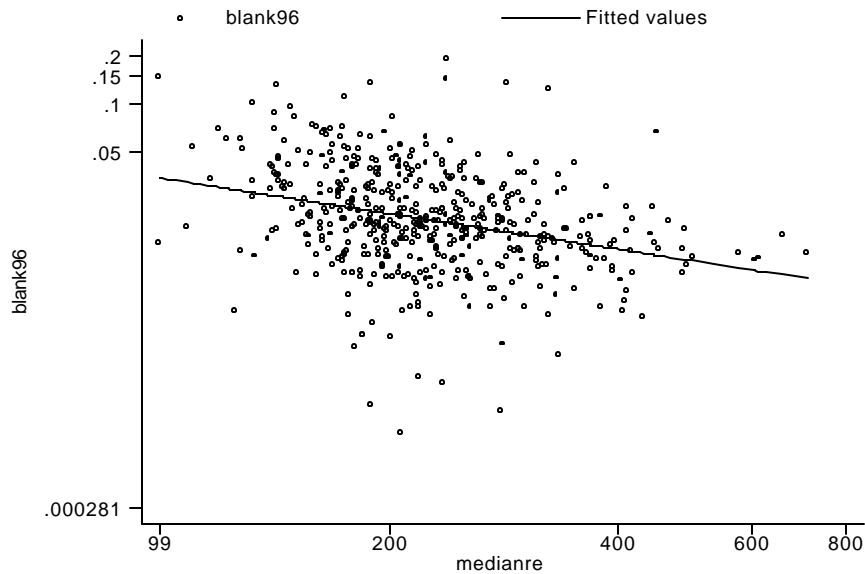
# Regression of Blank Ballots (1996) on Median Rent (1990)

	Coeff.	s.e.
Intercept	-0.36	0.48
Slope	-0.65	0.088
N	663	
R <sup>2</sup>	.077	



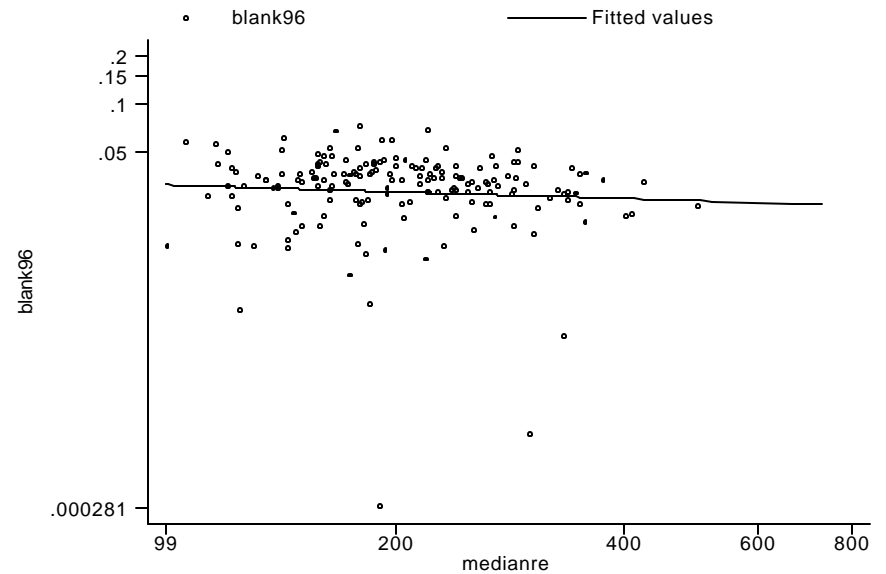
# Regression of Blank Ballots (1996) on Median Rent (1990), By Ballot Type

## Scanning



	Coeff.	s.e.
Intercept	0.040	0.48
Slope	-0.74	0.088
N	491	
R <sup>2</sup>	.10	

## Electronic



	Coeff.	s.e.
Intercept	-2.83	0.82
Slope	-0.14	0.15
N	172	
R <sup>2</sup>	.005	

## What to do?

- Run two separate regressions
  - Advantage: conceptually simple
  - Disadvantage: hypothesis testing cumbersome
- Interaction terms
  - Advantage: hypothesis testing facilitated
  - Disadvantage: conceptually complex

Interaction terms generally

$$y = \mathbf{b}_0 + \mathbf{b}_1 X_1 + \mathbf{b}_2 X_2 + \mathbf{b}_3 X_1 X_2 + \mathbf{e}$$

Rewriting,

$$y = \mathbf{b}_0 + (\mathbf{b}_1 + \mathbf{b}_3 X_2) X_1 + \mathbf{b}_2 X_2 + \mathbf{e}$$

# Interaction terms in the voting machine example

- Define  $S_c = 1$  if the county uses optical scanning, 0 otherwise
- Run this regression:

$$\text{blankpct}_c = \mathbf{b}_0 + \mathbf{b}_1 \times \text{rent}_c + \mathbf{b}_2 \times S_c + \mathbf{b}_3 \times S_c \times \text{rent}_c + \mathbf{e}_c$$

---

Note that if  $S_c = 0$  (i.e., electronic county), we have

$$\text{blankpct}_c = \mathbf{b}_0 + \mathbf{b}_1 \times \text{rent}_c + \mathbf{e}_c$$

If  $S_c = 1$  (i.e., scanned county), we have

$$\text{blankpct}_c = \mathbf{b}_0 + \mathbf{b}_1 \times \text{rent}_c + \mathbf{b}_2 + \mathbf{b}_3 \times \text{rent}_c + \mathbf{e}_c \text{ or}$$

$$\text{blankpct}_c = (\mathbf{b}_0 + \mathbf{b}_2) + (\mathbf{b}_1 + \mathbf{b}_3) \times \text{rent}_c + \mathbf{e}_c$$

# Doing this in *STATA*

```
. gen scan=ve96_cod=="5"
. gen s=scan
. gen scanrent=scan*rent
. reg blank rent scan scanrent
```

Source	SS	df	MS	Number of obs =	663
Model	48.597571	3	16.1991903	F( 3, 659) =	32.71
Residual	326.36878	659	.495248527	Prob > F =	0.0000
Total	374.966351	662	.566414427	R-squared =	0.1296
				Adj R-squared =	0.1256
				Root MSE =	.70374

blank	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
rent	-.14256	.1686598	-0.85	0.398	-.4737353 .1886153
scan	2.866141	1.051092	2.73	0.007	.8022468 4.930035
scanrent	-.6017711	.196494	-3.06	0.002	-.987601 -.2159413
_cons	-2.826596	.8975356	-3.15	0.002	-4.58897 -1.064222



# Standardized Regression

Comparing (Standardized) Apples  
with (Standardized) Oranges

# Which “matters” more in determining vote outcomes, popularity or the economy?

```
. reg vote drdi gallup
```

Source	SS	df	MS	Number of obs = 13		
Model	.038942217	2	.019471109	F( 2, 10)	=	20.01
Residual	.009732889	10	.000973289	Prob > F	=	0.0003
Total	.048675106	12	.004056259	R-squared	=	0.8000
				Adj R-squared	=	0.7601
				Root MSE	=	.0312

vote	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
drdi	1.908849	.545243	3.50	0.006	.6939719	3.123726
gallup	.2554055	.07231	3.53	0.005	.0942888	.4165223
_cons	.3422054	.0350065	9.78	0.000	.2642061	.4202047

# Solutions I

- Normalize into percentages
  - Take logs of everything
  - Advantage: elegant
  - Disadvantages:
    - Not always appropriate transform
    - Zero, negative numbers

$$\ln(y) = \mathbf{b}_0 + \mathbf{b}_1 \ln(x_1) + \mathbf{b}_2 \ln(x_2) + \mathbf{e}$$

Calculate  $\partial y / \partial x_1$  and  $\partial y / \partial x_2$  and rearrange terms :

$$\mathbf{b}_1 = \frac{\partial y / y}{\partial x_1 / x_1}, \mathbf{b}_2 = \frac{\partial y / y}{\partial x_2 / x_2}$$

# Solutions II

- Transform the variables into unit deviates (I.e., mean 0, s.d. 1)
  - Subtract each variable from its mean and divide by its standard deviation, I.e.:

$$z_{i,j} = \frac{(Z_{i,j} - \bar{Z}_i)}{s_{Z_i}}$$

# Doing this in *STATA*

```
. reg vote drdi gallup,beta
```

Source	SS	df	MS	Number of obs = 13	
Model	.038942217	2	.019471109	F( 2, 10) =	20.01
Residual	.009732889	10	.000973289	Prob > F =	0.0003
Total	.048675106	12	.004056259	R-squared =	0.8000
				Adj R-squared =	0.7601
				Root MSE =	.0312

vote	Coef.	Std. Err.	t	P> t	Beta
drdi	1.908849	.545243	3.50	0.006	.535644
gallup	.2554055	.07231	3.53	0.005	.5404141
_cons	.3422054	.0350065	9.78	0.000	.

# Decomposing Regression Effects

Direct and Indirect Effects

Recall the OLS solution

If

$$Y_i = \mathbf{b}_0 + \mathbf{b}_1 X_{1,i} + \mathbf{b}_2 X_{2,i} + \mathbf{e}_i$$

then

$$\hat{\mathbf{b}}_1 = \frac{\text{cov}(X_1, Y)}{\text{var}(X_1)} - \hat{\mathbf{b}}_2 \frac{\text{cov}(X_1, X_2)}{\text{var}(X_1)} \text{ and}$$

$$\hat{\mathbf{b}}_2 = \frac{\text{cov}(X_2, Y)}{\text{var}(X_2)} - \hat{\mathbf{b}}_1 \frac{\text{cov}(X_1, X_2)}{\text{var}(X_2)}$$

# Rearrange the first line

$$\hat{b}_1 = \frac{\text{cov}(X_1, Y)}{\text{var}(X_1)} - \hat{b}_2 \frac{\text{cov}(X_1, X_2)}{\text{var}(X_1)} \text{ or}$$

$$\frac{\text{cov}(X_1, Y)}{\text{var}(X_1)} = \hat{b}_1 + \hat{b}_2 \frac{\text{cov}(X_1, X_2)}{\text{var}(X_1)} \text{ or}$$

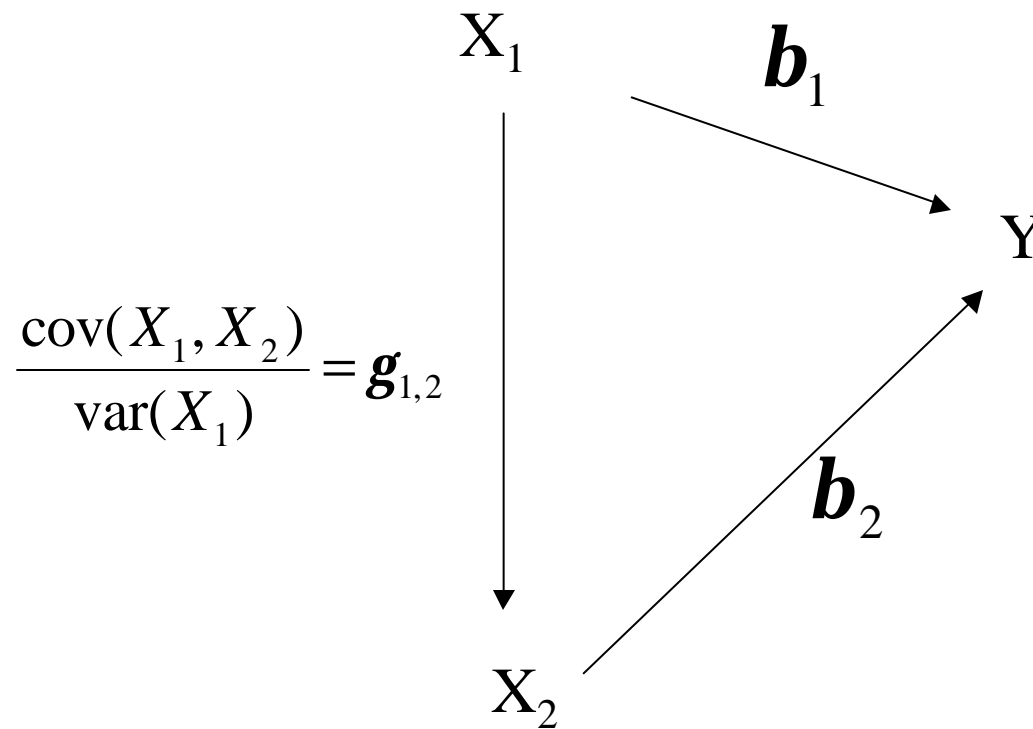
(Overall association b/t  $X_1$  and  $Y$ ) =

(Direct effect of  $X_1$  on  $Y$ ) +

(Direct effect of  $X_2$  on  $Y$ )  $\times$  (Bivariate effect of  $X_1$  on  $X_2$ )



# Graphically



(Overall association b/t  $X_1$  and  $Y$ ) =  
 (Direct effect of  $X_1$  on  $Y$ ) +  
 (Direct effect of  $X_2$  on  $Y$ )  $\times$  (Bivariate effect of  $X_1$  on  $X_2$ )

# Decomposing the effects of popularity and the economy on the vote

Effect	Bivariate	Direct	Indirect
Gallup	0.35 (100%)	0.26 (74%)	0.097 (26%)
Economy	2.64 (100%)	1.91 (72%)	0.74 (28%)

# Measurement Error

What Happens When You Can't  
Measure Things Perfectly?

# Suppose we measure $x$ with error?

Instead of observing  $x$ , we observe  $x' = x + e$   
( $e$  is random with mean  $\bar{e}$  and variance  $v_e$ )

$\therefore$  instead of doing the regression

$$y = \mathbf{a} + \mathbf{b}x + \mathbf{e},$$

we do the regression

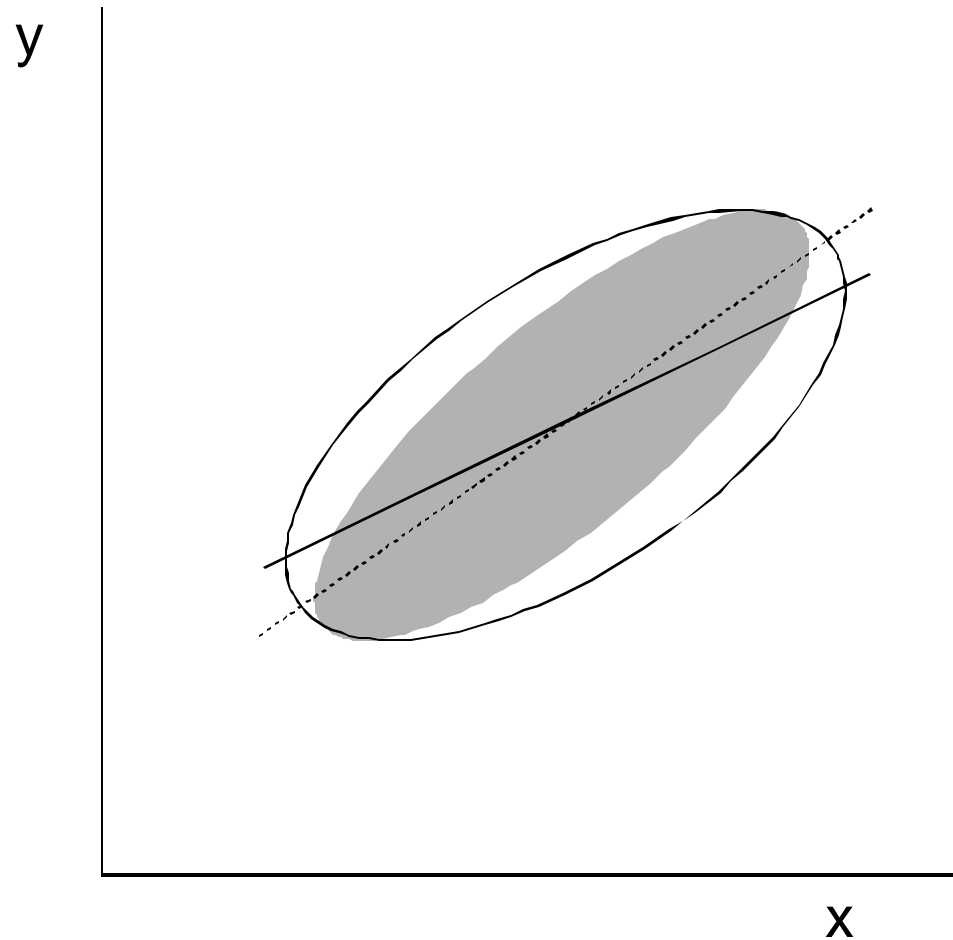
$$y = \mathbf{a} + \mathbf{b}'x' + \mathbf{e}.$$

What is the relationship between  $\mathbf{b}$  and  $\mathbf{b}'$ ?

Answer

$$\mathbf{b}' = \frac{\text{cov}(x, y)}{\text{var}(x) + \text{var}(e)}$$

# Errors in Independent Variables: The Picture



# Suppose we measure $y$ with error

Instead of observing  $y$ , we observe  $y' = y + e$   
( $e$  is random with mean  $\bar{e}$  and variance  $v_e$ )

$\therefore$  instead of doing the regression

$$y = \mathbf{a} + \mathbf{b}x + \mathbf{e},$$

we do the regression

$$y' = \mathbf{a} + \mathbf{b}'x + \mathbf{e}.$$

What is the relationship between  $\mathbf{b}$  and  $\mathbf{b}'$ ?

The answer

$$\mathbf{b}' = \frac{\text{cov}(x, y)}{\text{var}(x)} = \mathbf{b}$$

But...

- Standard errors and s.e.r. inflated
- $R^2$  deflated



# Errors in Dependent Variables: The Picture

