## Moving from correlation to causation

ECON 30331
Bill Evans
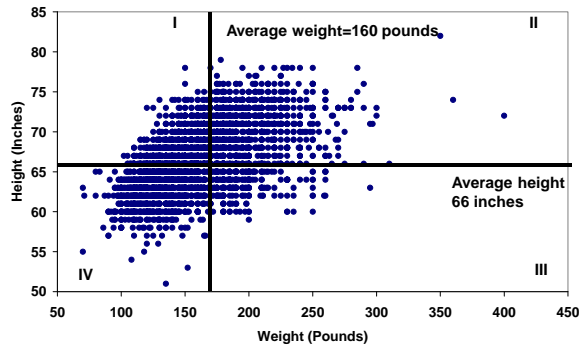
1

## Scatter plot

- Sample of N observations
  - Students, doctors, state, countries etc.

- For each observation, 2 pieces of data (X,Y)

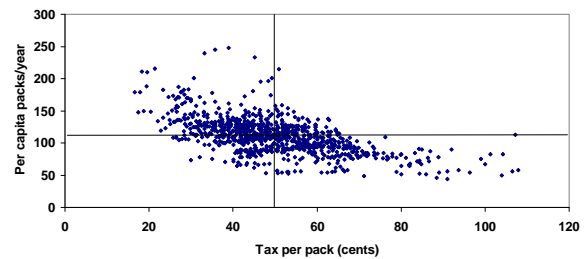- Plot each point for all observations in sample

2



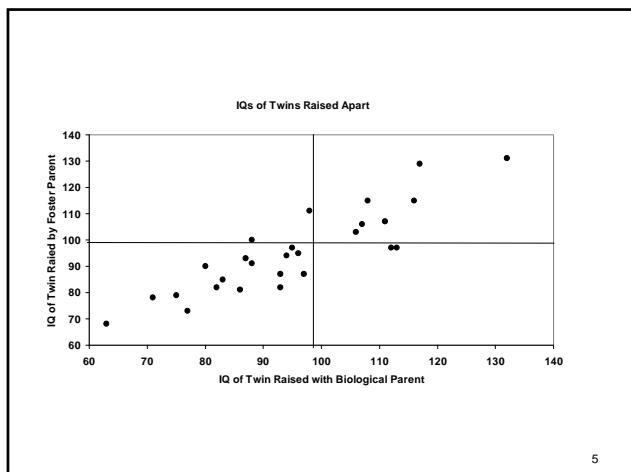Scatter Plot: Height and Weight of Adult Females

3



Cigarette Consumption and Taxes

4

**IQs of Twins Raised Apart**



5

## Covariance

- Measure of co-movement between variables

- Does the realization that X is above average convey any information about the likely value of Y?

- Identifies whether variables are 'statistically' related

6

## Covariance

- x and y are random variables

- $E[x] = \mu_x$      $Var(x) = \sigma^2_x$
- $E[y] = \mu_y$      $Var(y) = \sigma^2_y$

- $Cov(x,y) = E[(x - \mu_x)(y - \mu_y)] = \sigma_{xy}$
  $$= E[xy] - \mu_x\mu_y = \sigma_{xy}$$

7

$$If\ cov(x, y) > 0\ and\ y > \overline{y},$$
$$then,\ on\ average,\ x > \overline{x}$$

$$If\ cov(x, y) < 0\ and\ y > \overline{y},$$
$$then,\ on\ average,\ x < \overline{x}$$

8

## Problem

- Covariance is scale dependent

  - Covariance between height and weight will differ if measured in centimeters & kilograms or inches & pounds

- Not an attractive property for a measure of co-movement

9

## Demonstrate: Can show yourself

$$\text{cov}(x, y) = E[(x - \mu_x)(y - \mu_y)] = \sigma_{xy}$$
$$define: z = a + bx$$
$$\text{cov}(z, y) = E[(z - \mu_z)(y - \mu_y)]$$
$$z = a + bx$$
$$\mu_z = a + b\mu_x$$
$$(z - \mu_z) = b(x - \mu_x)$$
$$\text{cov}(z, y) = E[b(x - \mu_x)(y - \mu_y)]$$
$$= bE[(x - \mu_x)(y - \mu_y)] = b\sigma_{xy}$$

10

## Correlation coefficient

$$\rho(x, y) = \sigma_{xy} / (\sigma_x \sigma_y)$$

$$-1 \leq \rho(x, y) \leq 1$$

- Unlike the covariance, the correlation coefficient is NOT scale dependent

- The value is the same regardless of how x and y are measured
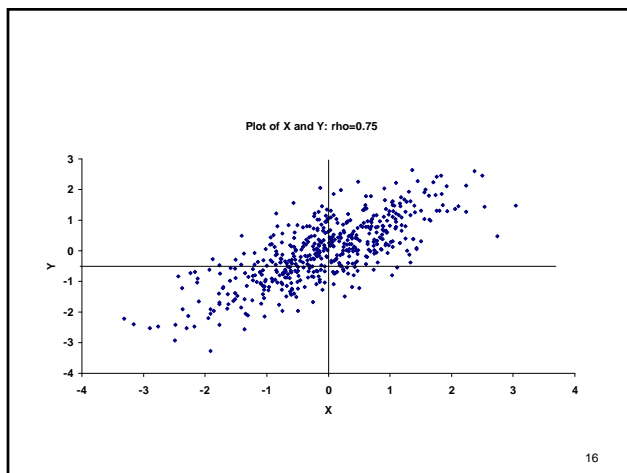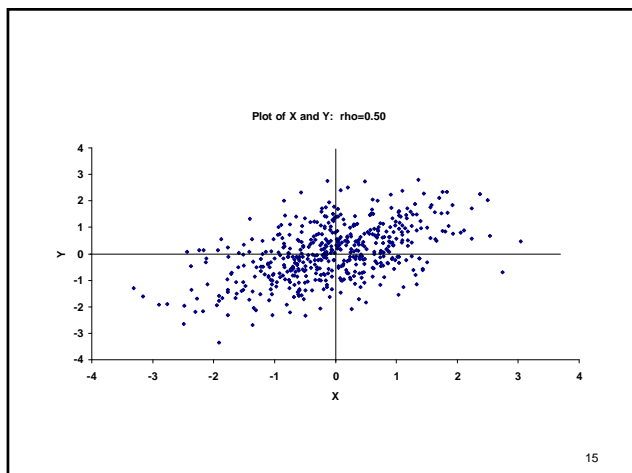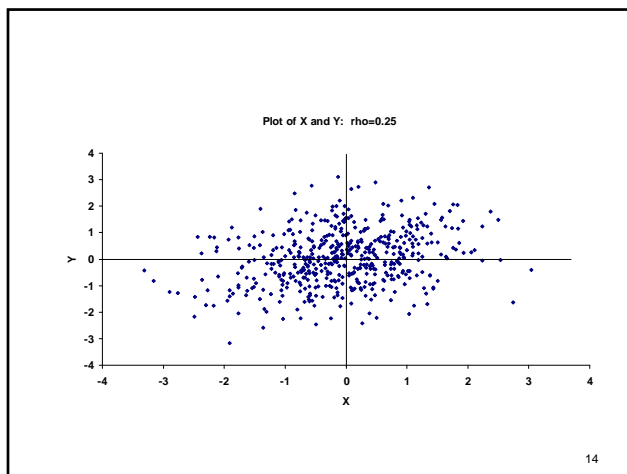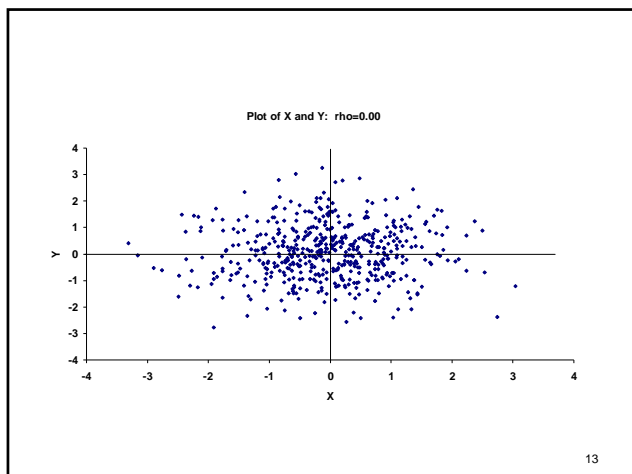
11

## Sample estimates

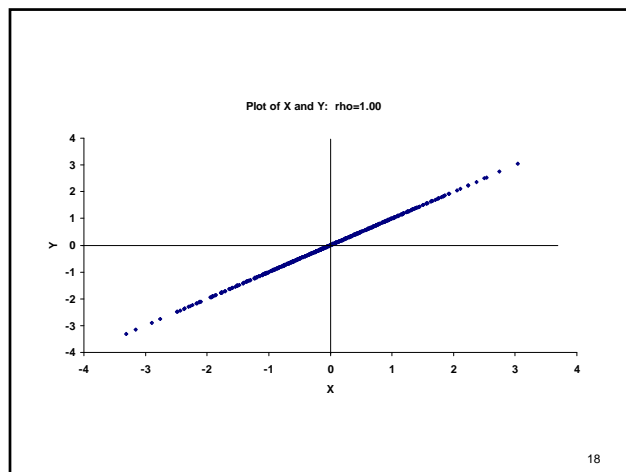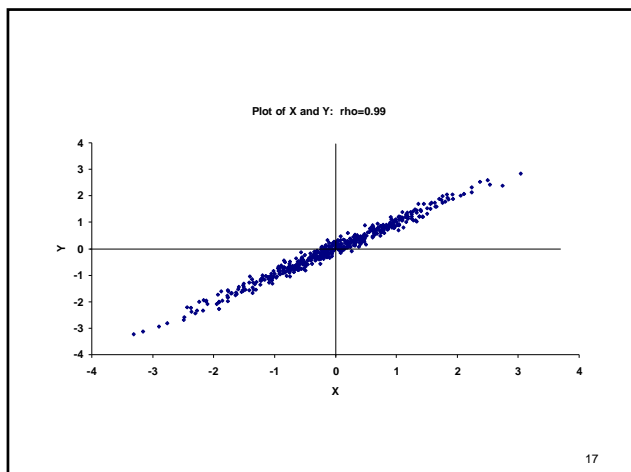$$\hat{\sigma}_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\hat{\sigma}_y^2 = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i - \bar{Y})^2$$

$$\hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

$$\hat{\rho} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x \hat{\sigma}_y}$$

12

Plot of X and Y: rho=0.00

13



Plot of X and Y: rho=0.25

14



Plot of X and Y: rho=0.50

15



Plot of X and Y: rho=0.75

16

Plot of X and Y: rho=0.99

17



Plot of X and Y: rho=1.00

18

## Cross-Sectional data

- Height and weight, men

- Height/weight, women

- Log(wages)/educ (m)

- Log(wage)/age (m)

19

## Cross-Sectional Data

- Husband/wife age

- Husband/wife educ

- Husband/wife height

- Father/son income

- Father/son educ.

20

## Cross-Sectional Data

- IQ's of Identical twins

- IQ's of fraternal twins

- IQ's of identical twins raised apart

- IQ's of siblings

- IQ's of unrelated children reared together

21

## Among undergrads in Intro Micro

- Math SAT/verbal SAT

- HS rank/total SAT

- GPA in micro/SAT

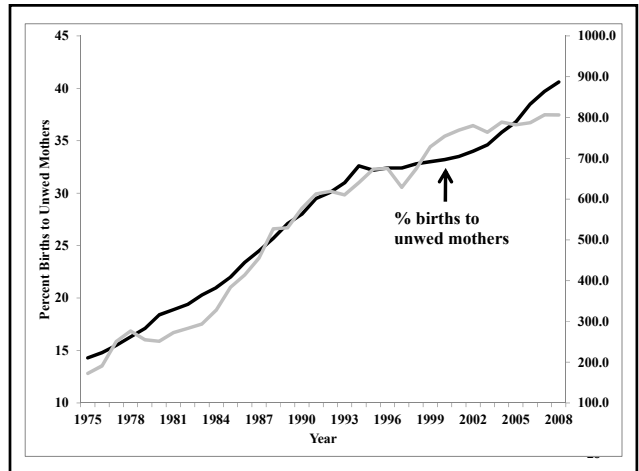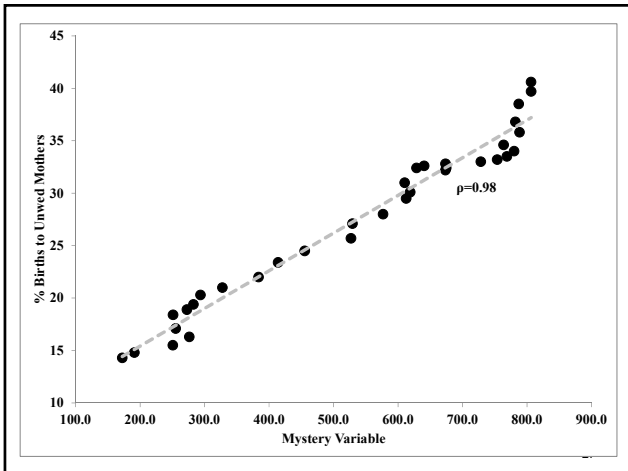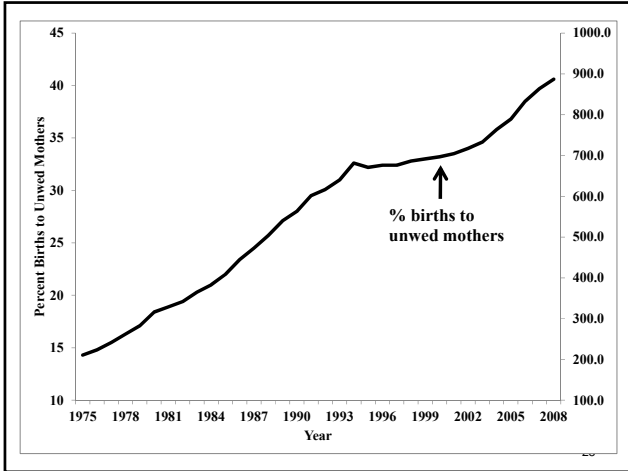- GPA in micro/HS percentile

22

## Limitation

- Correlation coefficient is a convenient way to measure a statistical relationship between two variables
- It does not however signify anything more than statistical observation
- It also does no get us any closer to saying whether something is causally related
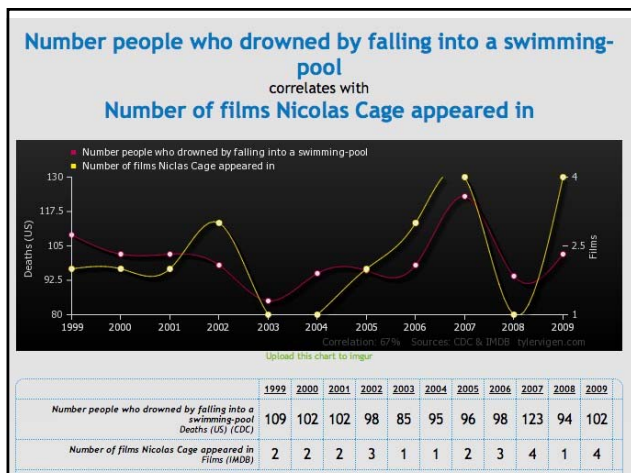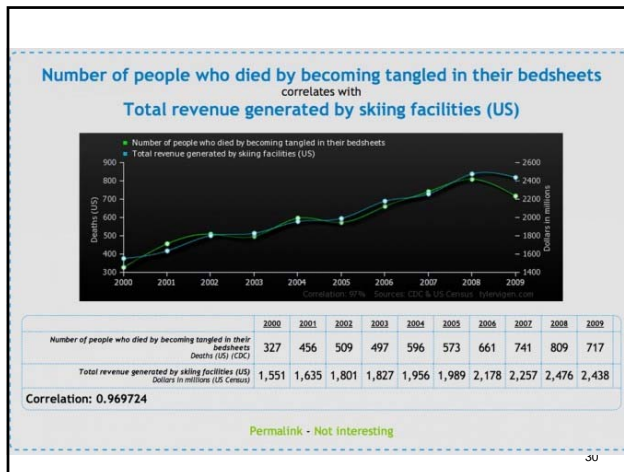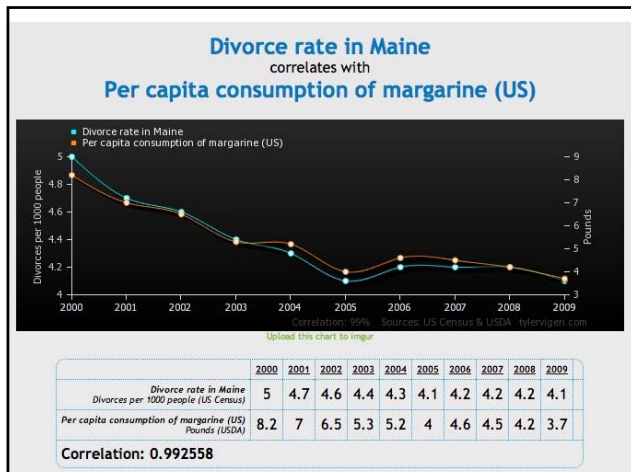- Correlation does not equal causation

23

## Births to unwed mothers

- Risen from 5% in 1960 to 37% in 2006
- Predictive of many child outcomes
  - Low birth weight, increased mortality, poor performance in schools, etc.
- Many potential explanations
  - Poor performance of male wages, rising divorce, availability of abortion
- Is there a magic bullet explanation?

24

## Basic economic model

- People/firms/organizations are purposeful

- Examples
  - Firms maximize profits
  - People maximize happiness/utility

- There are however limits or constraints on behavior
  - Consumers must pay prevailing prices
  - Firms have competitors

33

## Break variables into 2 groups

- Exogenous (external conditions)
  - Constraints on behavior
  - "Treatments" Factors that can be altered
  - "Independent" variables
- Endogenous outcomes
  - Choice variables
  - Outcomes of systems
  - "Dependent" variables

34

## Link between models/data

- Basic economic model has a prediction:
  - How quickly will demand fall when prices rise
  - What happens to outcomes (endogenous) when an external condition is changed (exogenous)

- Statistical goal: estimate the slope of the demand curve $\partial X / \partial P_x$

35

## Theory of Demand

- Core model of intermediate micro
- Model set up
  - Consumers derive utility from consumption of 2 goods (x,y)
    - $U = U(x,y)$
    - Utility function has specific properties
  - Pick utility maximizing bundle of (x,y) subject to constraints
    - Fixed prices for goods: $P_x$ and $P_y$
    - Fixed income, I

36

- Two implicit functions:

  $X = f(P_x, P_y, I)$

  $Y = g(P_x, P_y, I)$

- 3 "exogenous" variables: $P_x$, $P_y$ and I
- 2 "endogenous" variables: x and y

- Comparative statics: $\partial X/\partial P_x$ or $\partial X/\partial I$

37

---

- To build a statistical model that will allow us to predict the changes in outcomes, we need to assume a direction of causation
  - Prices alter how much you will purchase
  - Hours of study impact grades
  - Years of education alter earnings ability
- Our model will only accurately measure the impact of "x on y" if this assumption is correct

38

---

# Basic model: OLS

- Ordinary least squares regression

- Maybe 95% of statistics in social sciences

- Highly stylized models with tremendous capacity
  - Capacity comes from assumptions
  - If assumptions are correct – huge rewards
  - If assumptions are wrong, model is piece of junk

39

---

# Example

- State running a budget deficit
- Can raise taxes on cigarettes to cover shortfall
- Problem: when tax rate (t) increase, demand falls (Q) and will impact revenues
- Rev = tQ
- $\partial Rev/\partial t = t[\partial Q/\partial t] + Q$
- Key question: what is $\partial Q/\partial t$

40

## Slide 41

**Cigarette Consumption and Taxes**
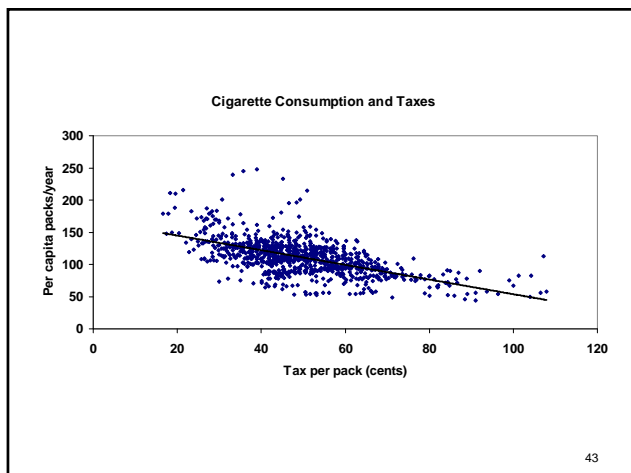


41

## Model

- $Y_i = \beta_0 + X_i \beta_1 + \varepsilon_i$
  - Linear
  - One input/one output
  - Y=quantity of cigarettes
  - X=taxes on cigarettes

- Parameter of interest
  - $\partial Y / \partial X = \beta_1$

42

## Slide 43

**Cigarette Consumption and Taxes**



43

## Problem

- Can always estimate basic model
  $$Y_i = \beta_0 + X_i \beta_1 + \varepsilon_i$$
- This does not mean the estimate for $\beta_1$ is any good
- Two typical problems that invalidate the estimate of $\beta_1$
  - Reverse causation (x may cause y but y may also cause x)
  - Omitted variables bias (some third factor may explain both y and x and hence, explain at least part of the reason why they are statistically related).

44

## Reverse Causation:
## An Economic Example

- Public finance economists are interested in the productivity of government spending

- Two largest components of local spending are schools and public safety
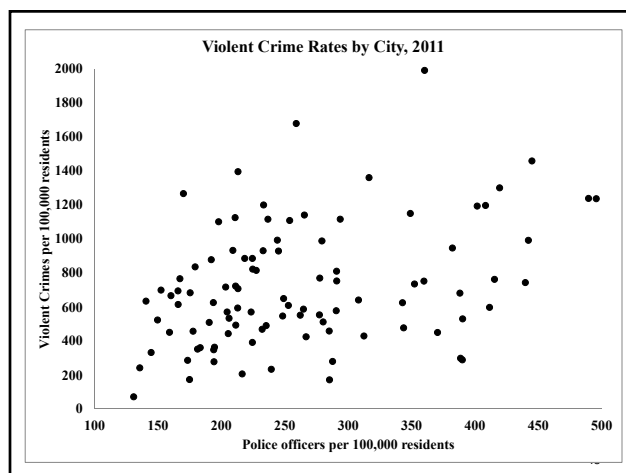
- Will hiring more police reduce crime?

45

- Let y=crime rate (crime per person)
- Let x=police employed per person

- Interested in estimating the gradient
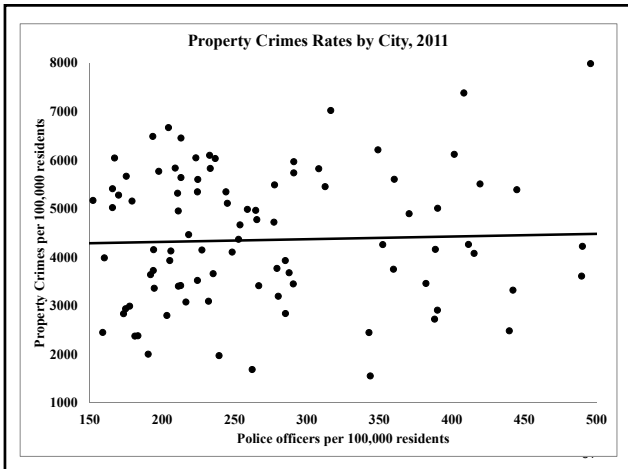- $\partial y / \partial x$ how will crime change when a city hires more police

46

- Collect data on a cross section of cities
  - 61 cities with populations in excess of 250K

- Estimate basic model
  $Y_i = \beta_0 + X_i\beta_1 + \varepsilon_i$

- What do you think is the most frequent sign (+ or -) on police?

47



Violent Crime Rates by City, 2011

Violent Crime Rates by City, 2011



Property Crimes Rates by City, 2011



Property Crimes Rates by City, 2011

## Highest violent crime rates, largest 100 cities

- Crime Rank

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

52

## Omitted variables bias

- Teen childbearing is associated with a number of poor economics outcomes later in life
  - Lower education
  - Lower earnings
  - Higher rates of welfare participation

## Outcomes of women aged 30-34 by Teen motherhood status

| Outcome | Teen mother | Not a teen mother |
|---|---|---|
| < a HS degree | 19.8% | 6.6% |
| ≥ college degree | 9.0% | 43.0% |
| In poverty | 30.9% | 13.0% |
| On welfare | 6.9% | 2.6% |
| Income from work | $23,884 | $36,206 |

54

## Omitted variables bias

- Teen childbearing is associated with a number of poor economics outcomes later in life
  - Lower education
  - Lower earnings
  - Higher rates of welfare participation
- Teen moms are not an random sample of the population – more likely from
  - Poor schools
  - Families with lower-educated moms
  - Families with teen mothers themselves

## Washington Post, August 15, 1997, page A3

*Lasting Effects Found From Spanking Children Antisocial Behavior Is Increased, Study Says*

Spanking children is apt to cause more long-term behavioral problems than most parents who use that approach to discipline may realize, a new study reports.

56

Children who get spanked regularly are more likely over time to cheat or lie, to be disobedient at school and to bully others, and have less remorse for what they do wrong, according to the study by researchers at the University of New Hampshire. It is being published this month in the medical journal Archives of Pediatrics and Adolescent Medicine. "When parents use corporal punishment to reduce antisocial behavior, the long-term effect tends to be the opposite," the study concludes.

57

## 4 tasks

- Outline basic statistical models
  - How do we get the estimates?
- Demonstrate properties – we want to know
  - When do we get "good" estimates?
  - When do we not??
- Illustrate how they are used in research
  - Do the estimates provide good internal and external validity
- Demonstrate how to obtain results using STATA

58

## Take away skills

- Some will use these techniques in the future – make your professor proud

- Some will not – your job is then to be a critical reader of the newspaper

59