# BIAS FROM SPECIFICATION OR MEASUREMENT

# ANSCOMBE'S QUARTET
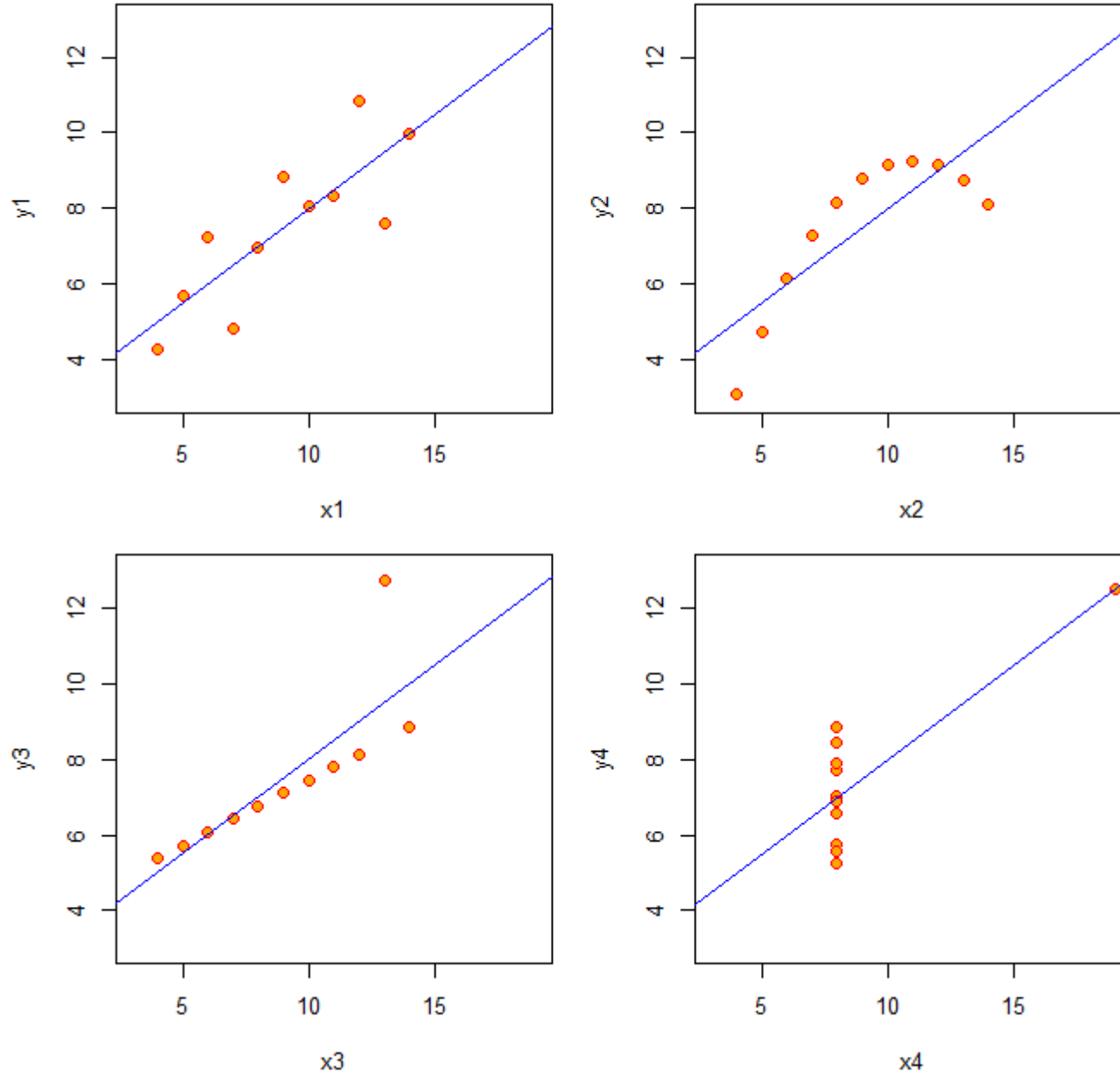## A LESSON IN MODEL FIT

# Anscombe's Quartet

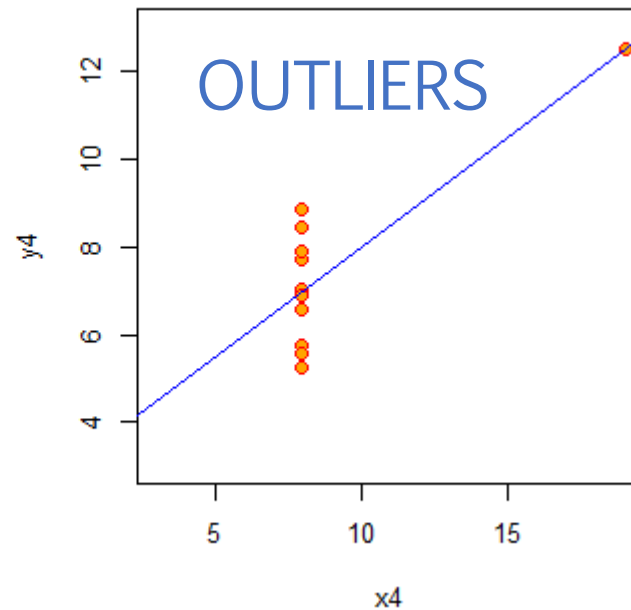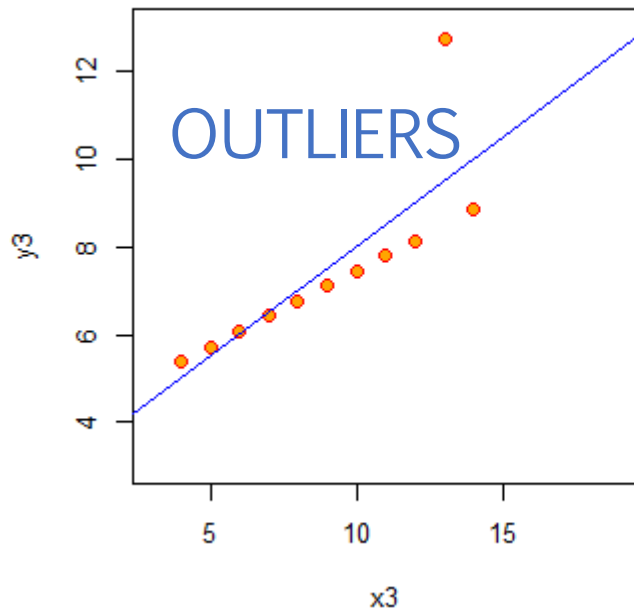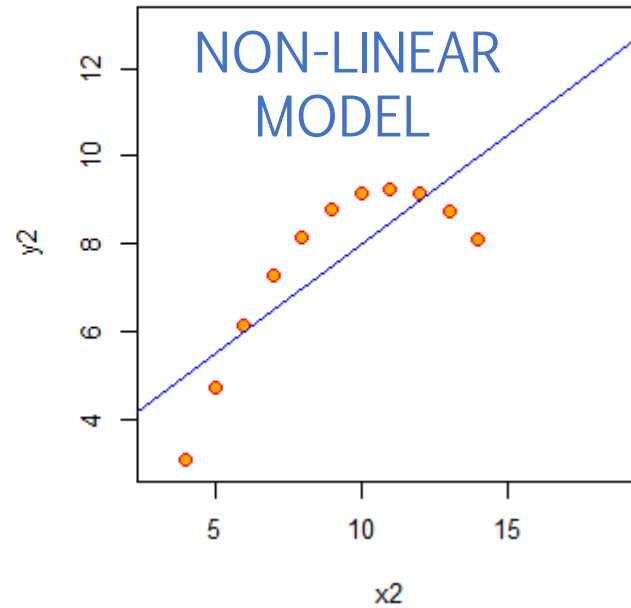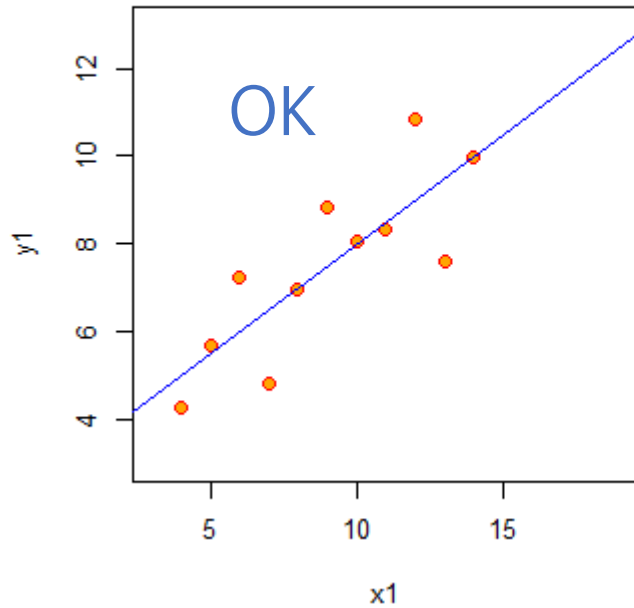| | x1 | y1 | x2 | y2 | x3 | y3 | x4 | y4 |
|---|---|---|---|---|---|---|---|---|
| 1 | 10 | 8.04 | 10 | 9.14 | 10 | 7.46 | 8 | 6.58 |
| 2 | 8 | 6.95 | 8 | 8.14 | 8 | 6.77 | 8 | 5.76 |
| 3 | 13 | 7.58 | 13 | 8.74 | 13 | 12.74 | 8 | 7.71 |
| 4 | 9 | 8.81 | 9 | 8.77 | 9 | 7.11 | 8 | 8.84 |
| 5 | 11 | 8.33 | 11 | 9.26 | 11 | 7.81 | 8 | 8.47 |
| 6 | 14 | 9.96 | 14 | 8.1 | 14 | 8.84 | 8 | 7.04 |
| 7 | 6 | 7.24 | 6 | 6.13 | 6 | 6.08 | 8 | 5.25 |
| 8 | 4 | 4.26 | 4 | 3.1 | 4 | 5.39 | 19 | 12.5 |
| 9 | 12 | 10.84 | 12 | 9.13 | 12 | 8.15 | 8 | 5.56 |
| 10 | 7 | 4.82 | 7 | 7.26 | 7 | 6.42 | 8 | 7.91 |
| 11 | 5 | 5.68 | 5 | 4.74 | 5 | 5.73 | 8 | 6.89 |
| | | | | | | | | |
| Mean | 9 | 7.5 | 9 | 7.5 | 9 | 7.5 | 9 | 7.5 |
| Variance | 11 | 4.12 | 11 | 4.12 | 11 | 4.12 | 11 | 4.12 |
| | | | | | | | | |
| Correlation | 0.816 | | 0.816 | | 0.816 | | 0.816 | |
| | | | | | | | | |
| Regression | $y = 3 + 0.5x$ | | $y = 3 + 0.5x$ | | $y = 3 + 0.5x$ | | $y = 3 + 0.5x$ | |

Four datasets that produce IDENTICAL descriptive stats, correlations, and regression models

Anscombe's 4 Regression data sets

BUT THEY ARE VERY DIFFERENT RELATIONSHIPS!

Anscombe's 4 Regression data sets

OK

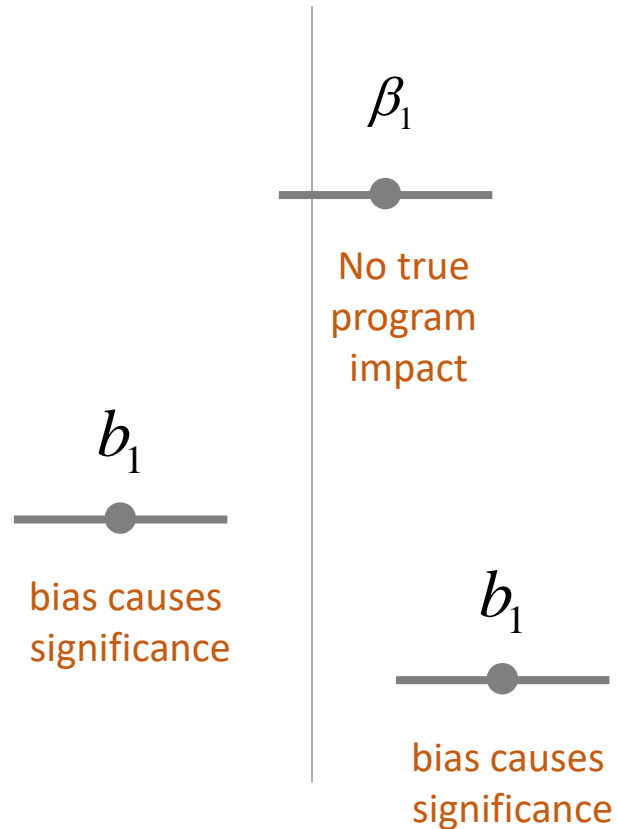NON-LINEAR MODEL

OUTLIERS

OUTLIERS

Anscombe's Quartet is often cited because (a) whoever created this example is a genius, and (b) it is a vivid demonstration of causes and consequences of SPECIFICATION BIAS.

We will consider what happens to slopes when outliers are present, or we use a linear specification when the relationship is non-linear.
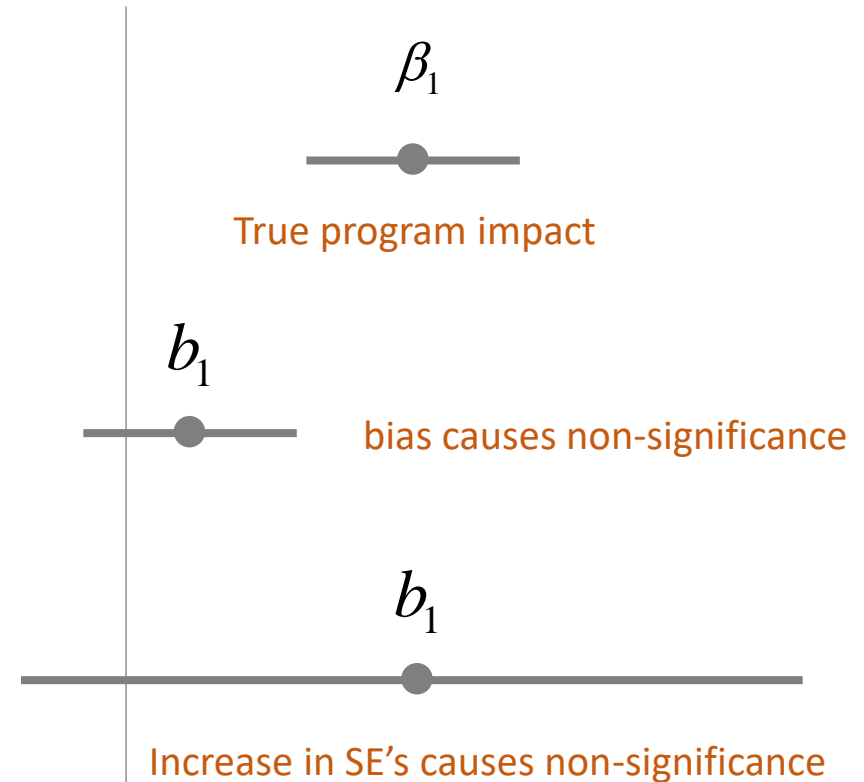
# CLASSES OF INFERENTIAL FAILURE
## TYPE I AND TYPE II ERRORS

# TYPE I ERROR
## FALSE POSITIVE
### CLAIMING PROGRAM HAS IMPACT
### WHEN IT DOESN'T

$\beta_1$

No true program impact

$b_1$

bias causes significance

$b_1$

bias causes significance

Type I errors are typically caused by OVB

# TYPE II ERROR
## FALSE NEGATIVE
### FAILING TO IDENTIFY TRUE
### PROGRAM IMPACT

$\beta_1$

True program impact

$b_1$

bias causes non-significance

$b_1$

Increase in SE's causes non-significance

Type II Errors can be caused by bias or inflated standard errors
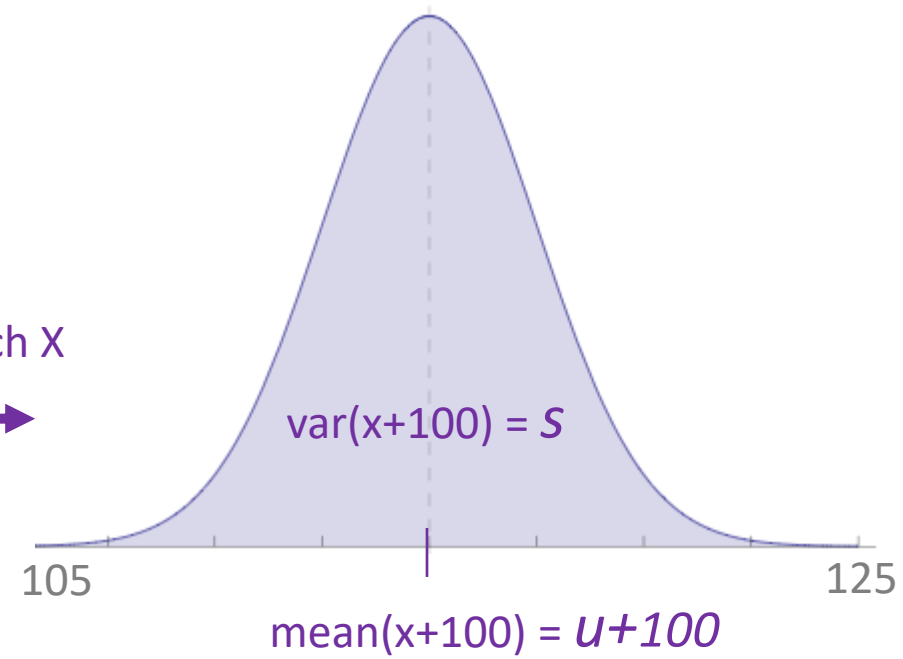
# IMPLICATIONS OF
## MEASUREMENT ERROR

var(x) = *S*

Add 100 to every X

→

var(x+100) = *S*

mean(x) = *u*

mean(x+100) = *u+100*

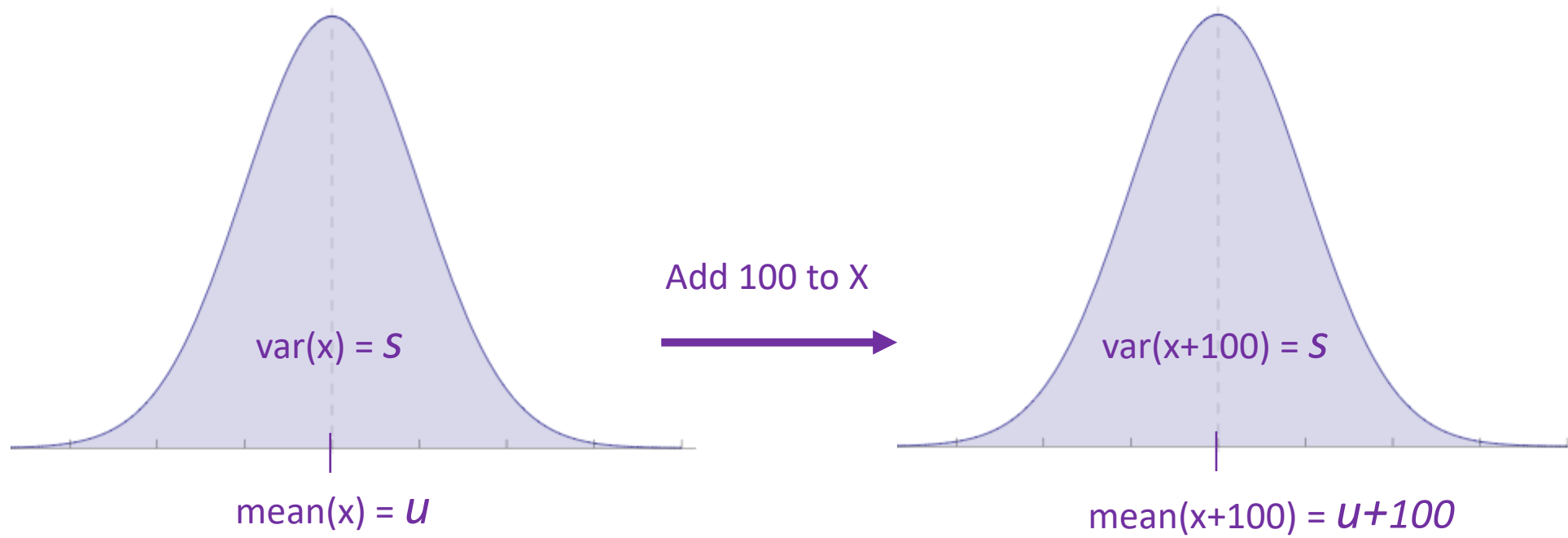"Linear Transformations"

X2 = X1 + 100

Variance of X is unchanged

var(x) = *S*

Add 100 to each X

var(x+100) = *S*

5                    25

105                    125

mean(x) = *u*

mean(x+100) = *u+100*

Y=b0+**b1**(X)

Y=b0+**b1**(X+100)

After **linear transformations** slopes b1 are identical

$var(x2) = var(x1)$ → slopes and standard errors same
$mean(x2) = mean(x1) + 100$ → x-axis moves right
*Intercept b0 (y when x=0) will be different*

var(x) = $S$

Add 100 to X

var(x+100) = $S$

mean(x) = $u$

mean(x+100) = $u+100$

"Linear Transformations"

X2 = X1 + **100**

Must add the **same constant** to every value of X

Just moves the distribution to right or left

X1

var(x1) = $S$

mean(x1) = $u$

X2

var(x2) = $S + \varepsilon$

mean(x2) = $u$

Measurement Error

X2 = X1 + $\varepsilon$

Add random error to every X.

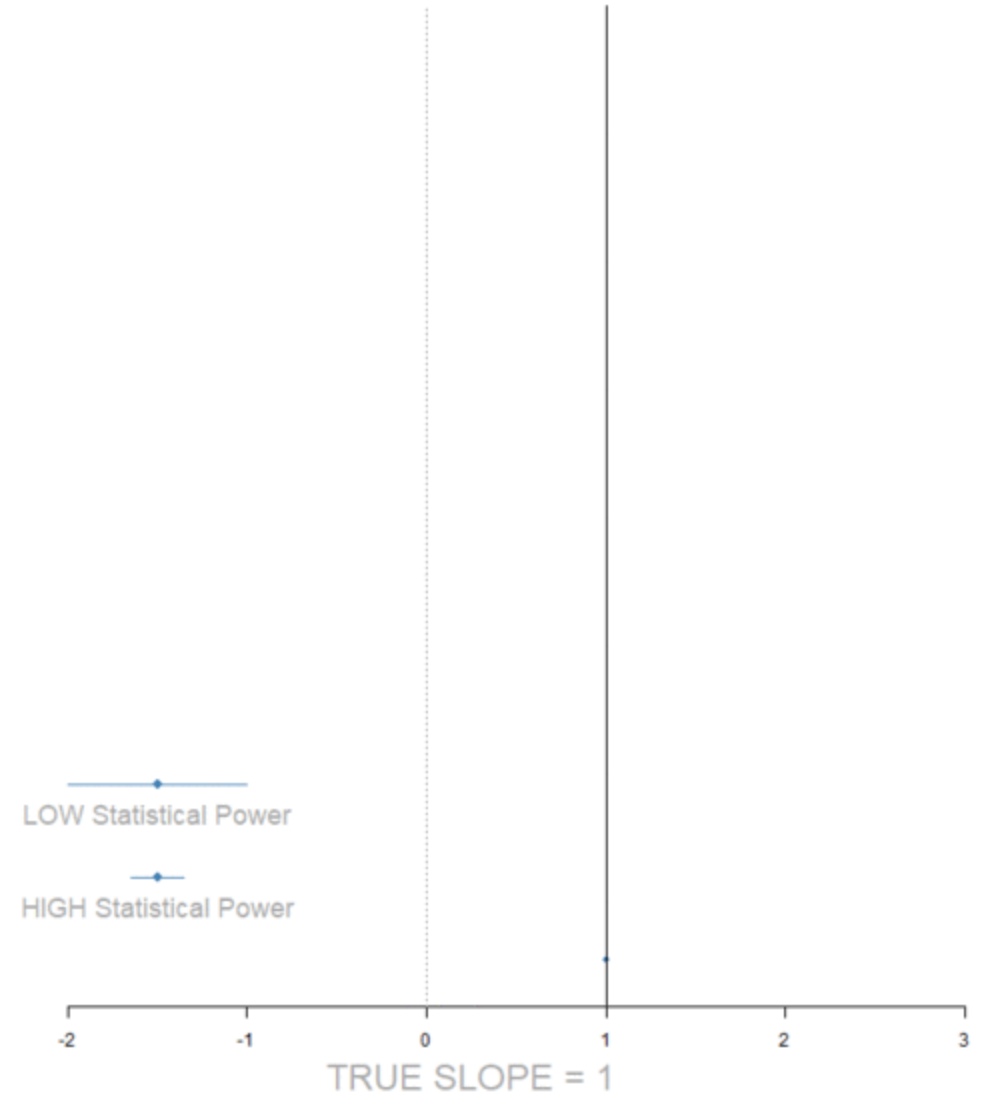*Random means each X is equally likely to be over-measured as under-measured.*

X2 has the same mean as X1, but more variance

# ADDING MEASUREMENT ERROR TO THE DV



(no bias)

Increase in standard errors

Type II Error risk

**Confidence Intervals for Slope Estimates**

LOW Statistical Power

HIGH Statistical Power

TRUE SLOPE = 1

# ADDING MEASUREMENT ERROR TO THE INDEPENDENT VARIABLE: "ATTENUATION BIAS"

$\beta_1$

True program impact

$b_1$

Smaller slope
Smaller SE

slope with measurement error

$$b_1 \downarrow = \frac{\text{cov}(x_1, y)}{\text{var}(x_1) \uparrow}$$

$$SE_{b1} \downarrow = \frac{residual}{samplesize \cdot \text{var}(x_1) \uparrow}$$

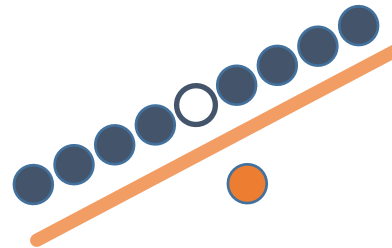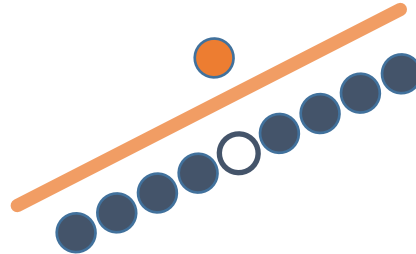**Measurement Error Added to the Independent Variable: Attenuation Bias**

Toward b1=zero

**Measurement Error Added to the Independent Variable: Attenuation Bias**

Toward b1=zero

14

# OUTLIERS

# SLOPES TOO LARGE
## SE LARGER

# SLOPES OK
## SE LARGER

# SLOPES TOO SMALL
## SE LARGER

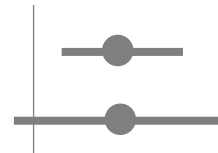Extreme of X:
Risk of bias in slope ↑
Risk of false positive
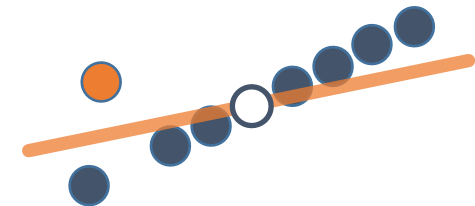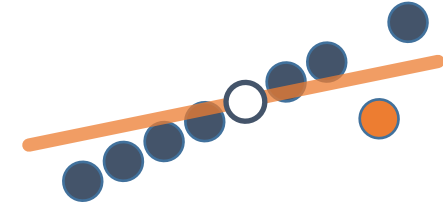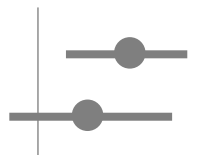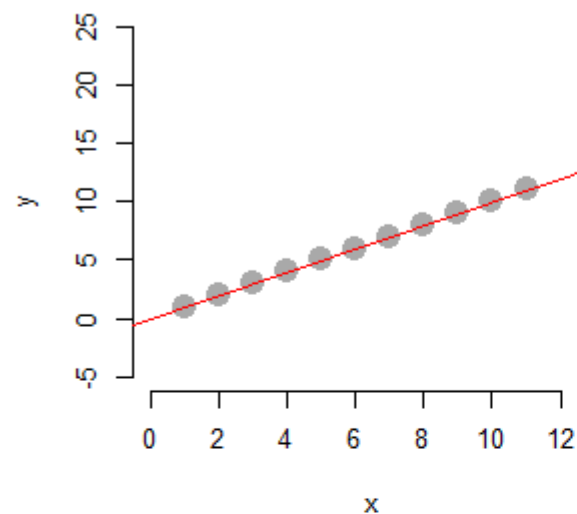
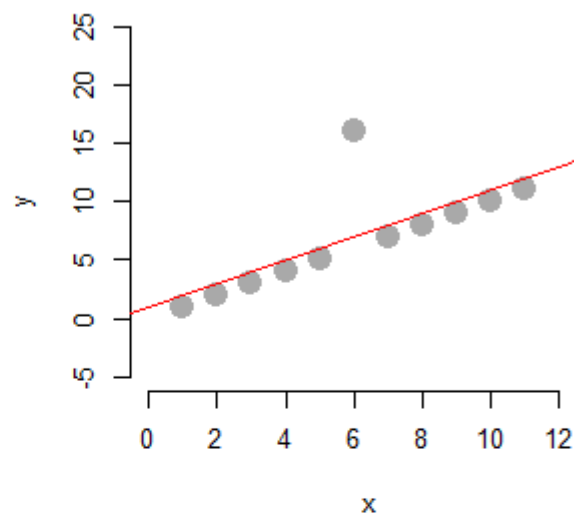Middle of X:
Don't bias slope
Increased risk of false negative

Extreme of X:
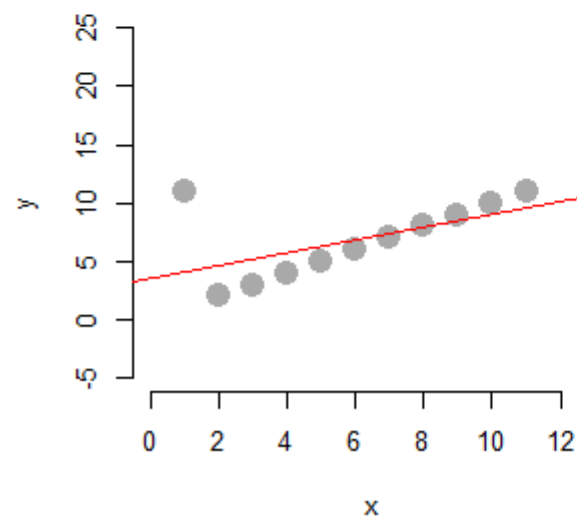Risk of bias in slope ↓
Increased risk of false negative

**Case 1**

**Case 2**

**Case 3**

**Case 4**

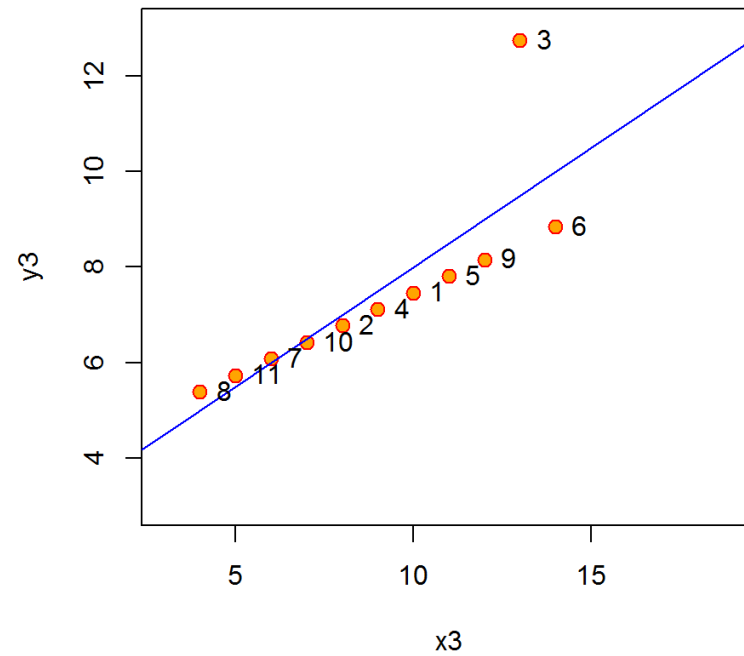| | Dependent variable: | | | |
|---|---|---|---|---|
| | y | | | |
| | (1) | (2) | (3) | (4) |
| x | 1.00*** | 1.00*** | 0.55* | 1.45*** |
| | (0.00) | (0.30) | (0.26) | (0.26) |
| Constant | 0.00*** | 0.91 | 3.64* | -1.82 |
| | (0.00) | (2.06) | (1.78) | (1.78) |

# IDENTIFYING OUTLIERS USING RESIDUALS AND COOK'S DISTANCE

| | Dependent variable: | |
|---|---|---|
| | y3 | |
| | (1) | (2) |
| x3 | 0.50*** | 0.35*** |
| | (0.12) | (0.0003) |
| Constant | 3.00** | 4.01*** |
| | (1.12) | (0.003) |
| Observations | 11 | 10 |
| $R^2$ | 0.67 | 1.00 |
| Adjusted $R^2$ | 0.63 | 1.00 |
| Note: | $p<0.1$; **$p<0.05$**; $p<0.01$ | |



Original Data



After Outlier Removal

# LOGGED
## REGRESSION MODELS

**Level-Level Model**

Number of Nonprofits

**Log-Linear Model**

**Log-Log Model**

Number of Nonprofits

Metro Population (logged)

Metro Population (logged)

The other common type of outlier comes from highly-skewed data. These problems can be fixed using log transformation of variables to convert data into a linear format for better model fit.

**Level-Level Model**

Highly-Skewed Data

**Linear-Log Model**

**Level-Level Model**

**Log-Log Model**

Linear Relationship

# NON-LINEAR RELATIONSHIPS
## QUADRATIC MODELS

Linear: $Y = b0 + b1(X_1) + e$

Quadratic: $Y = b0 + b1(X_1) + b2(X_1)^2 + e$



**Linear Fit**

**Quadratic Fit**