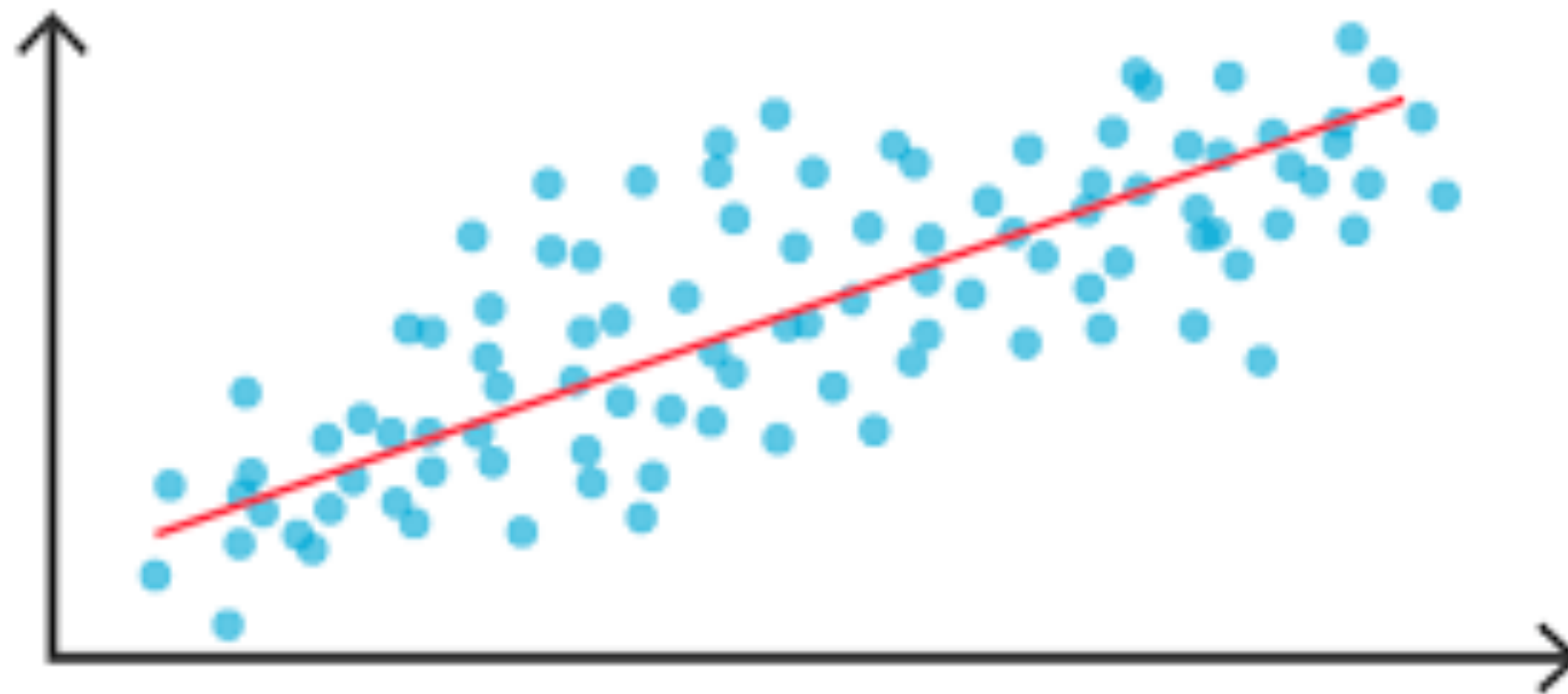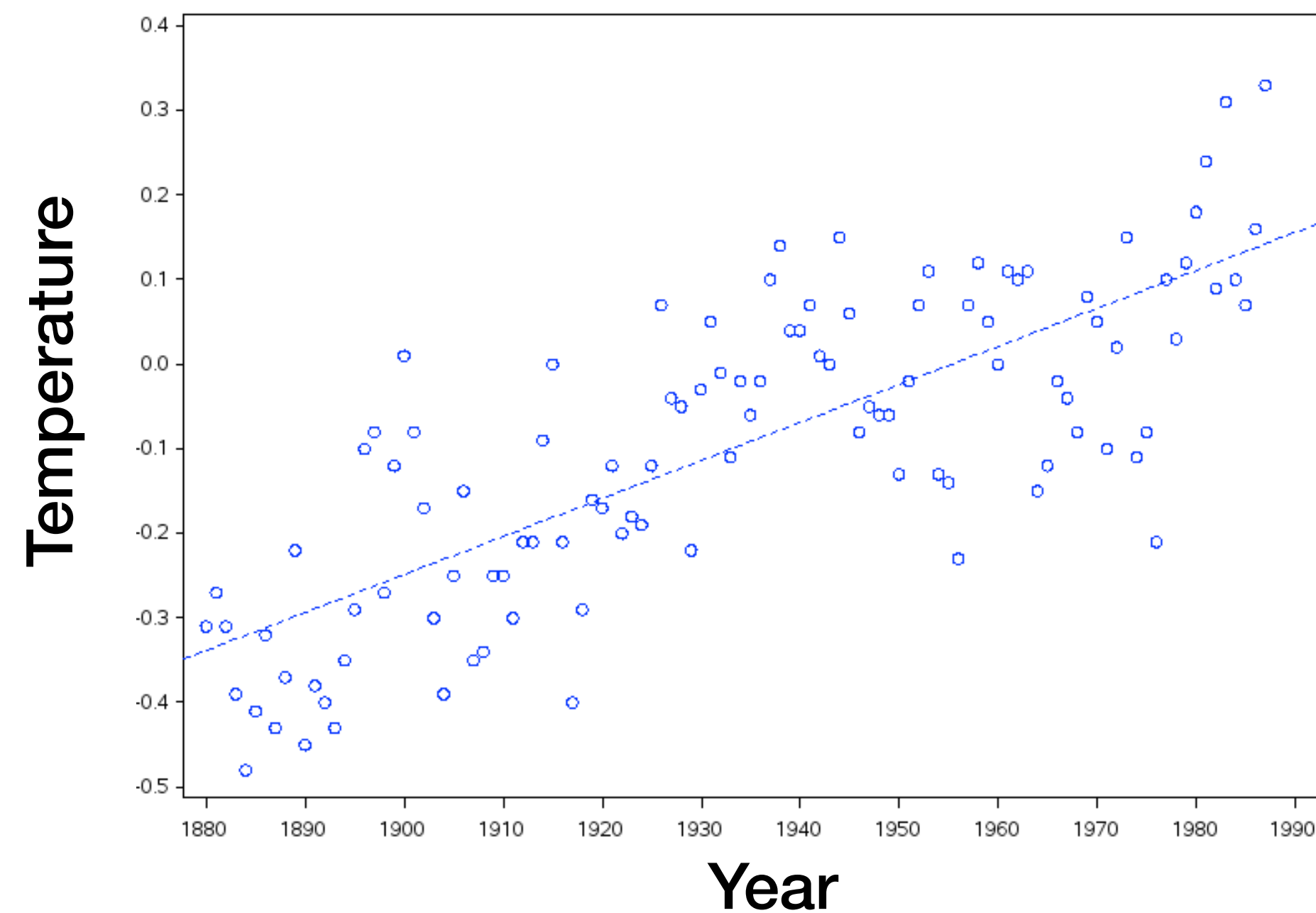# Foundational Statistics
## Simple Linear Regression

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$
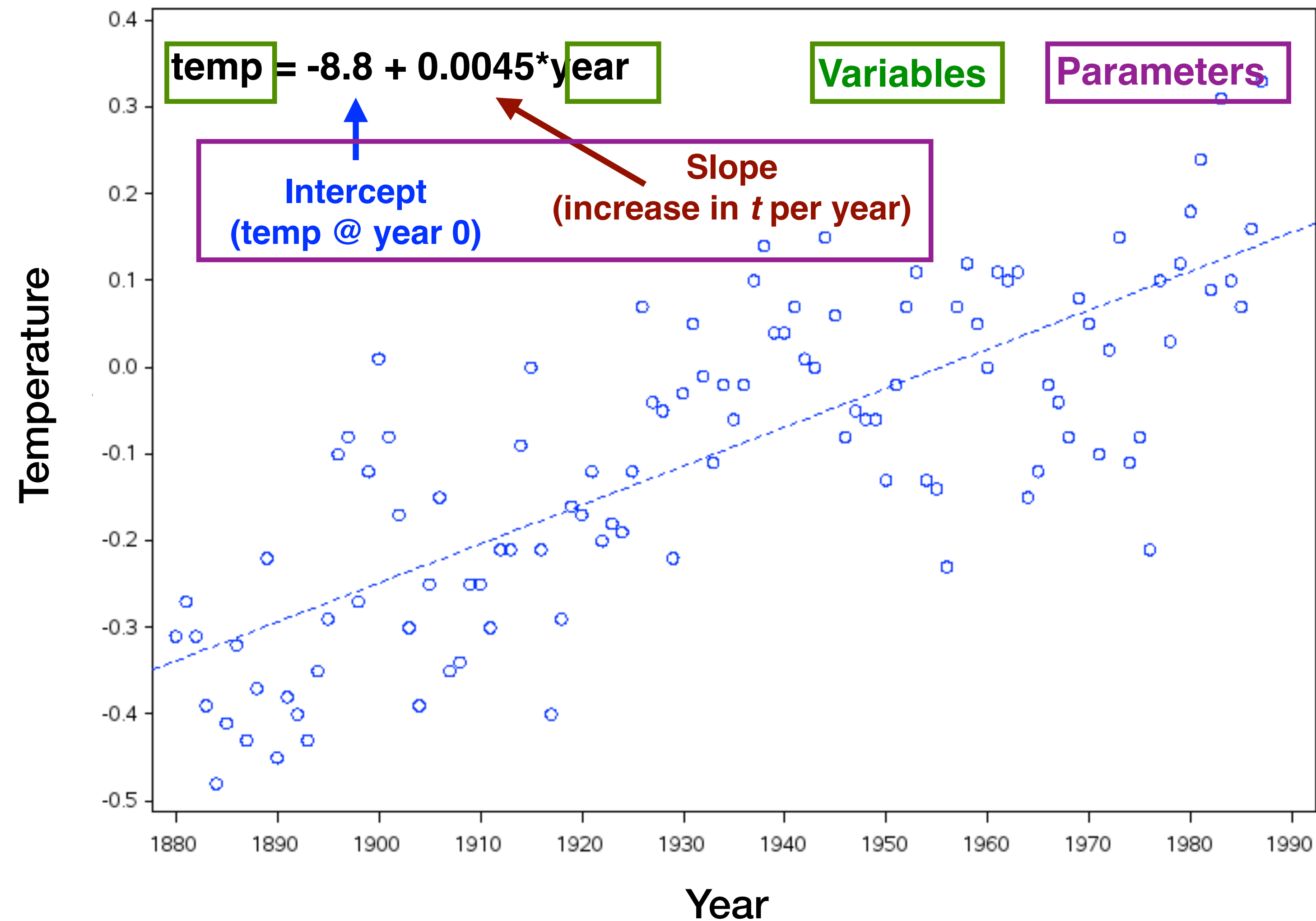
# A linear model for two numeric variables

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

| $y_i$ | $=$ | $\beta_0$ | $+$ | $\beta_1$ | $\times$ | $x_i$ | $+$ | $\varepsilon_i$ |
|---|---|---|---|---|---|---|---|---|
| response variable | $=$ | population intercept | $+$ | population slope | $\times$ | predictor variable | $+$ | error |

intercept term

slope term

model

# A generic linear model for two numeric variables

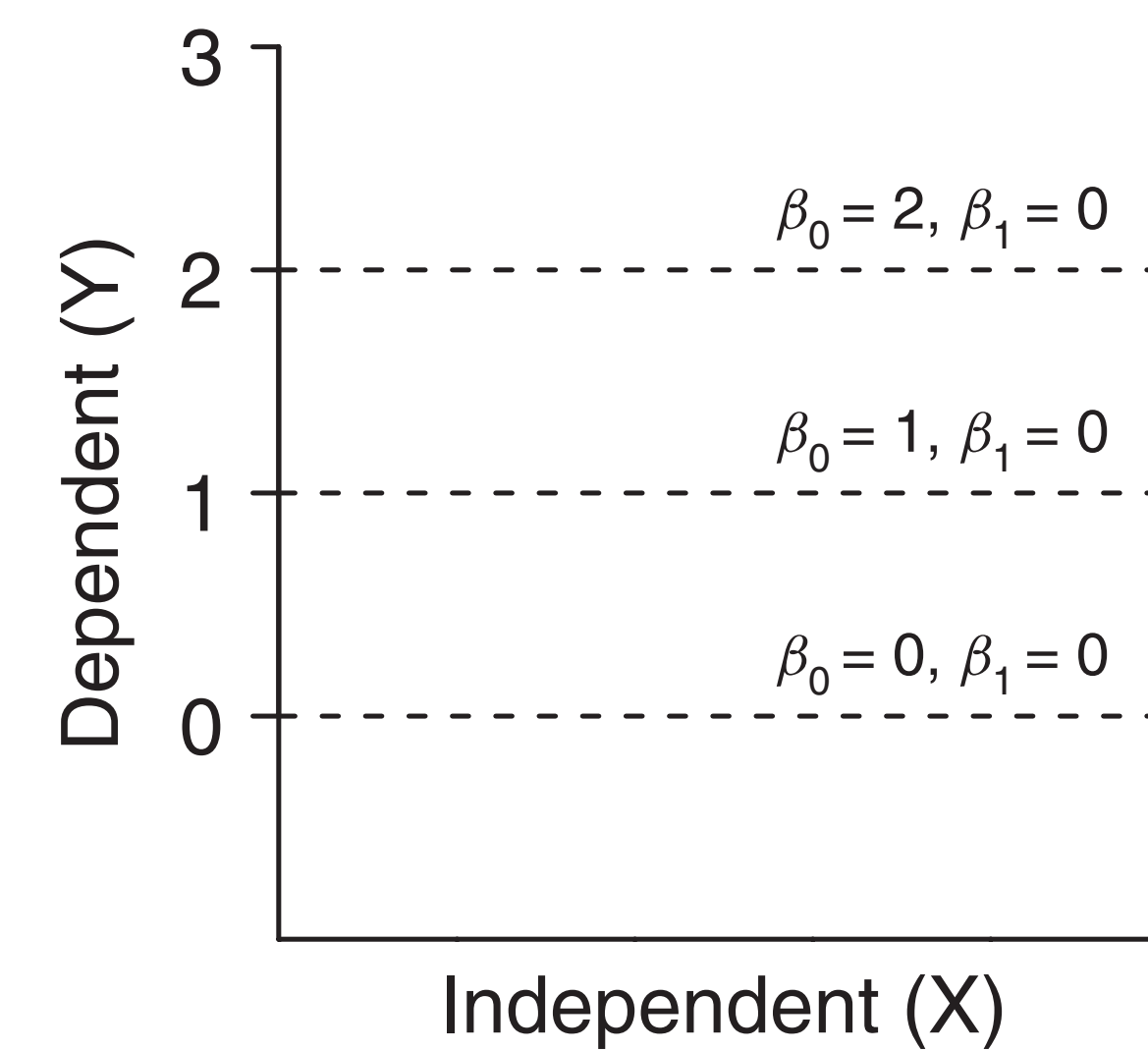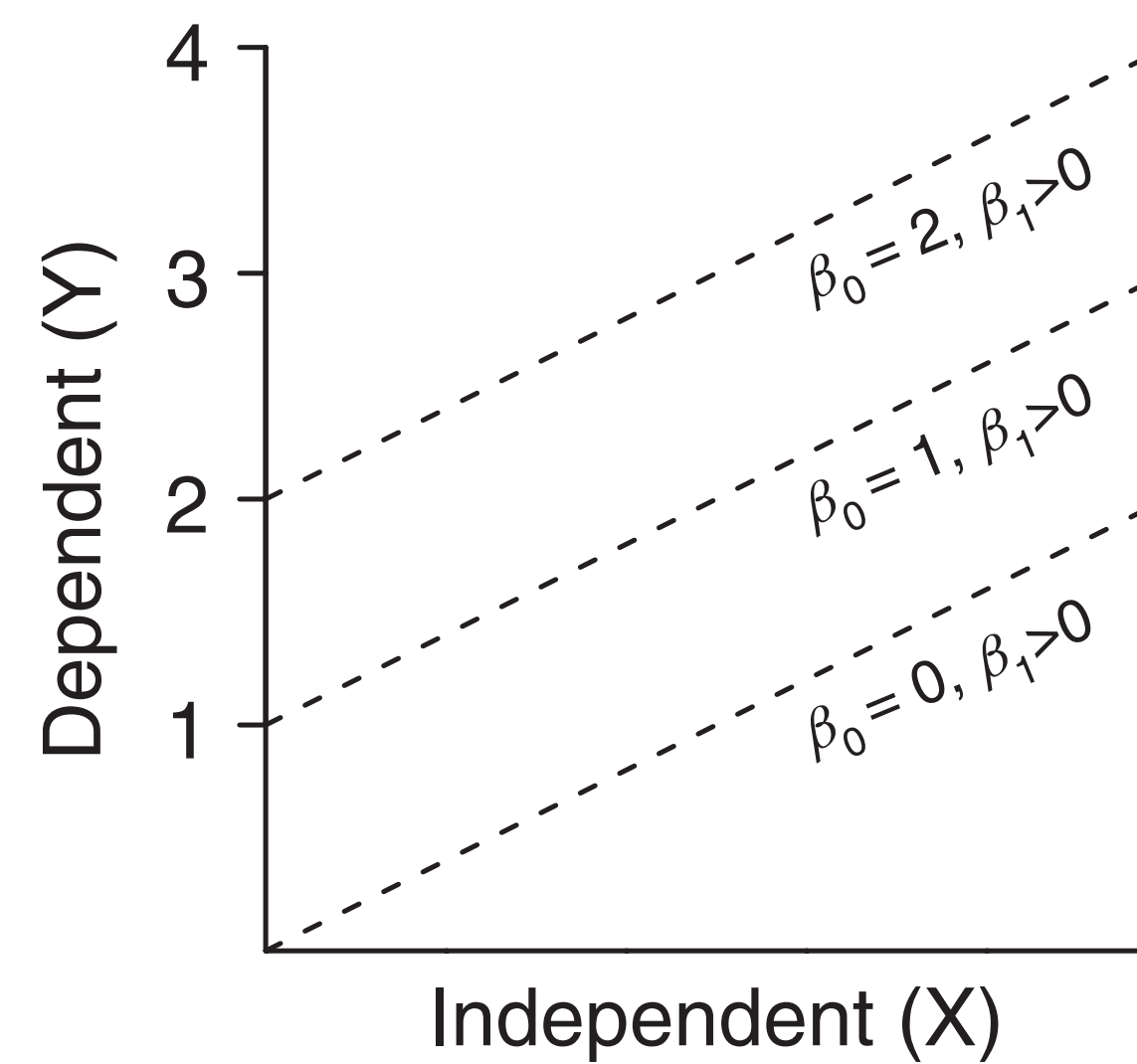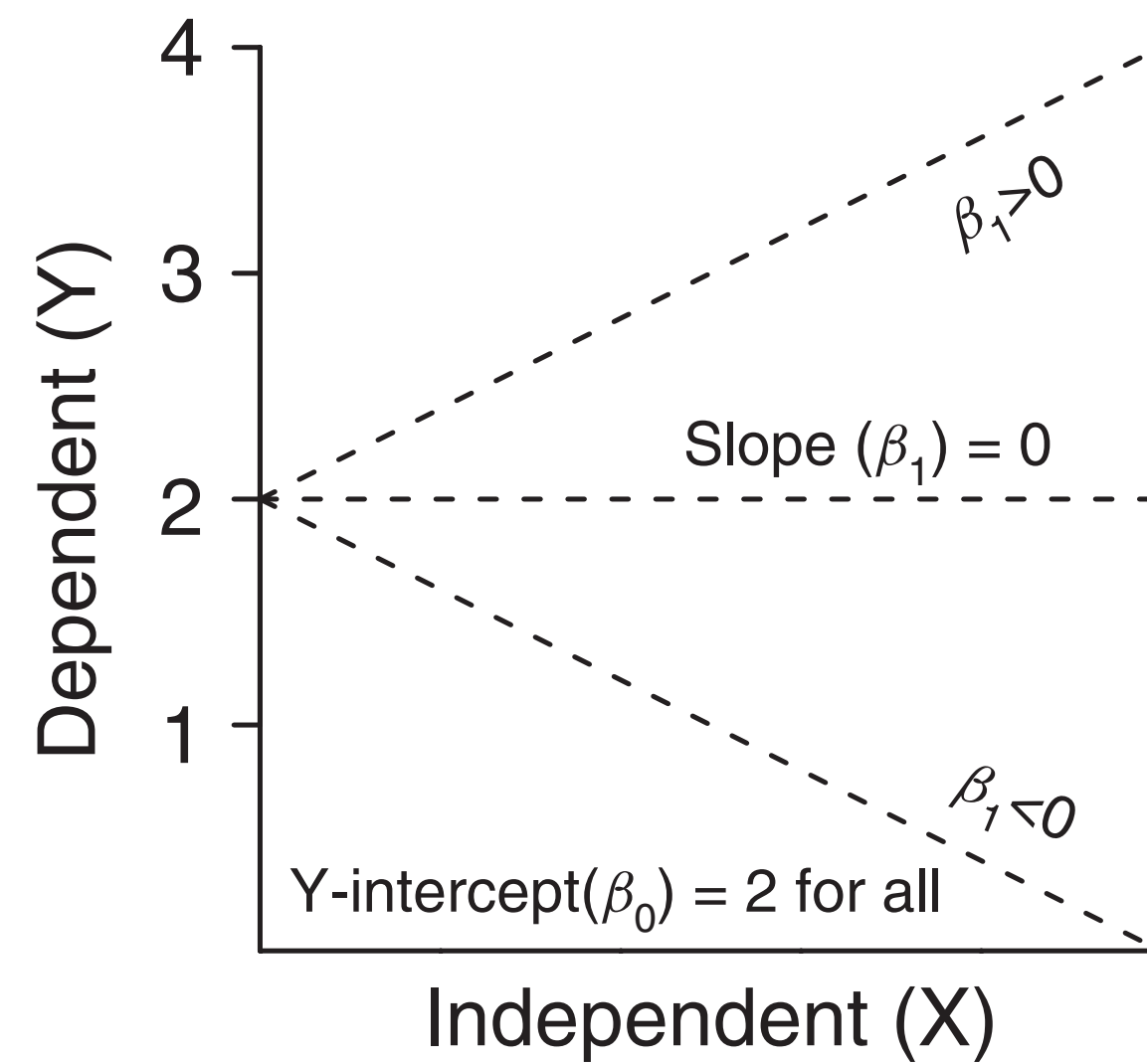**Ordinary Least Squares (OLS):** one method for estimating $\beta_0$ and $\beta_1$ from a random sample.

- Construct a "best fit" linear function to model variation in *y*

- Function is derived such that the total vertical distance between observed *y*-values and the line are minimized

- The *y*-intercept and slope of the line are our sample-based estimates for the population $\beta_0$ and $\beta_1$

- This approach (also called "Model I regression") assumes that the *x*-variable is measured without error

# **Ordinary Least Squares (OLS):** How it works

# Hypothesis tests in linear regression

$$H_0 : \beta_1 = 0$$  (the population slope equals zero)

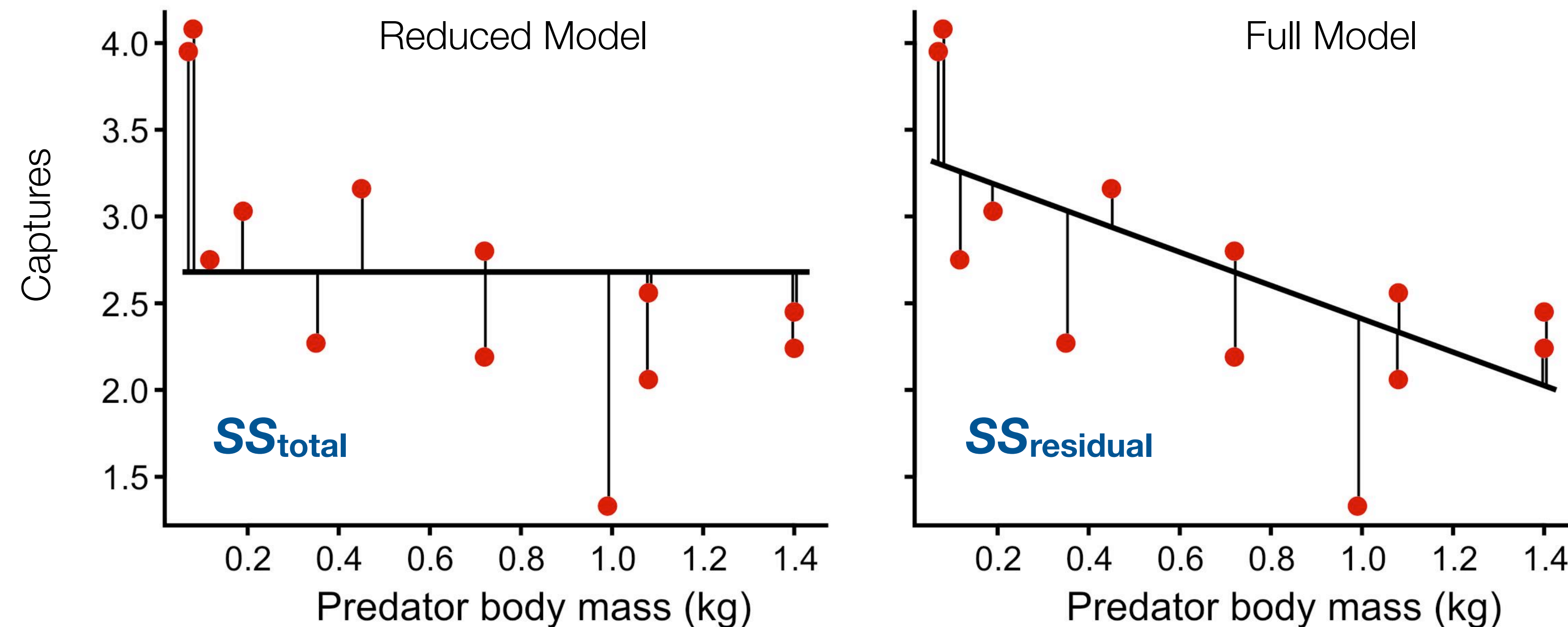$$H_0 : \beta_0 = 0$$  (the population y-intercept equals zero)

Comparing <u>full and reduced models</u> to test the
Null Hypothesis of $\beta_1 = 0$

$full\ model\ (H_A) -$     $y_i = \beta_0 + \beta_1 x_i + error_i$

$reduced\ model\ (H_0) -$     $y_i = \beta_0 + 0x_i + error_i$

$= \boxed{\beta_0 + error_i}$

**Model with no slope:**
**No linear effect of *x* on *y***

# Comparing <u>full and reduced models</u> to test the Null Hypothesis of $\beta_1 = 0$



<u>To test null hypothesis:</u>

1) fit a "reduced" model without slope term (fit under $H_0$)

2) fit the "full" model with slope term added back (fit under $H_A$)

**3) use the full and reduced models to calculate a <u>test statistic that reflects the ratio of explained to unexplained variation</u> by the full model**

# Comparing <u>full and reduced models</u> to test the Null Hypothesis of $\beta_1 = 0$

(i) the variation that is explained by the model ($SS_{Model}$)

$$SS_{Model} = \quad SS_{Total} \ (reduced \ model) - SS_{Residual} \ (full \ model)$$

<span style="color:#8B1A4A">The bigger this difference, the better our full model is, relative to the null model</span>

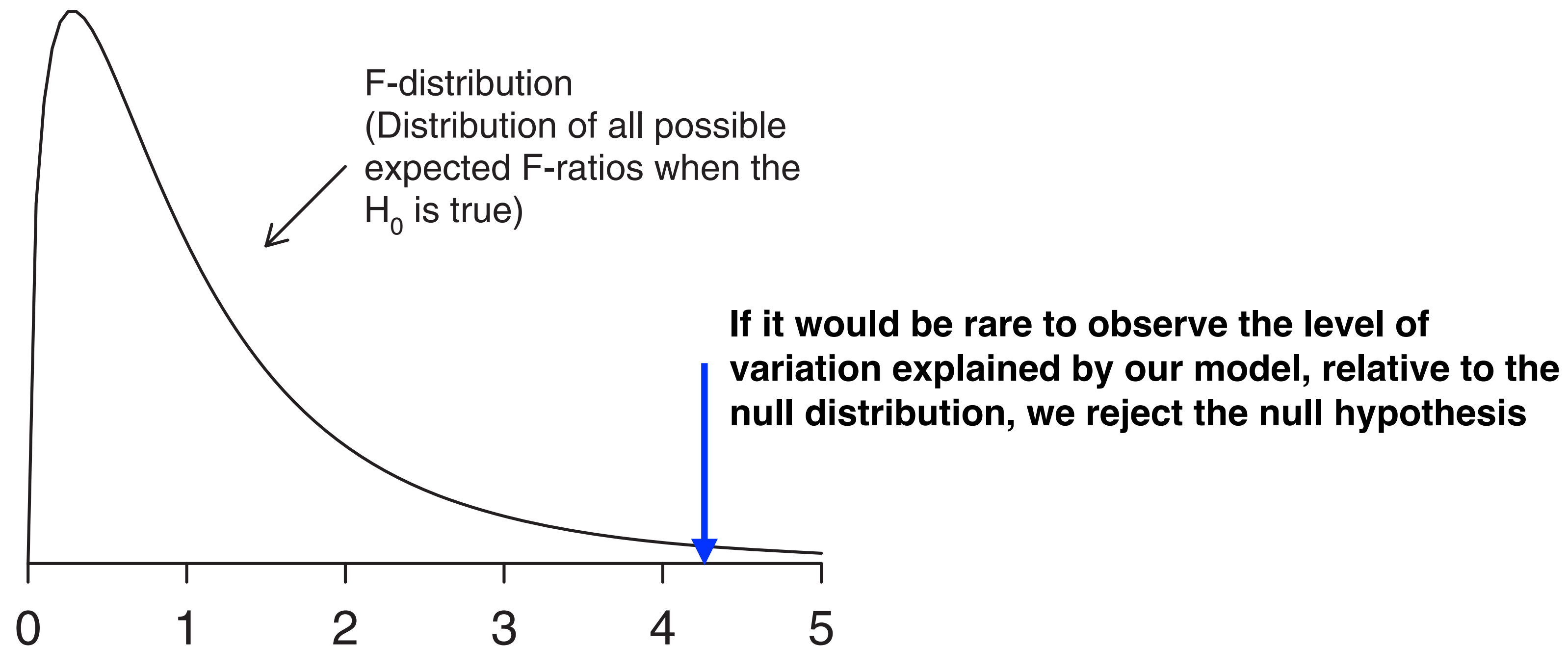(ii) the variation that is unexplained by the model ($SS_{Residual}$)

$$SS_{Residual} \ (full \ model)$$

<span style="color:#8B1A4A">If this value is small, our model explains much of the variation in $y$</span>

**So, If our model's ratio of explained to unexplained variation is high, we reject the null hypothesis of our predictor variables having zero explanatory power**

# Comparing <u>full and reduced models</u> to test the Null Hypothesis of $\beta_1 = 0$

**$F$-ratio (our test statistic)** $= \dfrac{SS_{reduced} - SS_{full}}{SS_{full}} = \dfrac{Var_{explained}}{Var_{unexplained}}$

F-distribution
(Distribution of all possible
expected F-ratios when the
$H_0$ is true)

**If it would be rare to observe the level of variation explained by our model, relative to the null distribution, we reject the null hypothesis**

0    1    2    3    4    5

# Assumptions of the *F*-ratio test for $\beta_1 = 0$

**1. Linear relationship between *y* and *x*, under *H*$_A$**

(Check using scatter plot)

**2. Bivariate normally distributed**

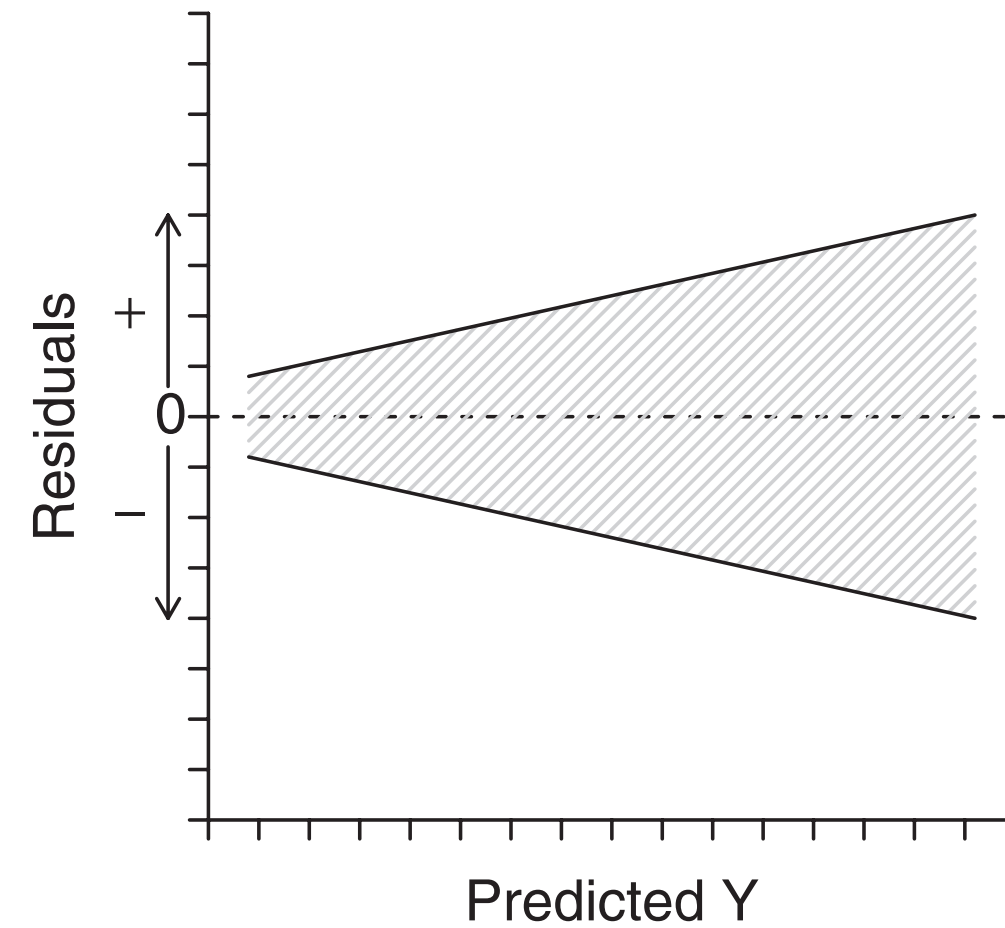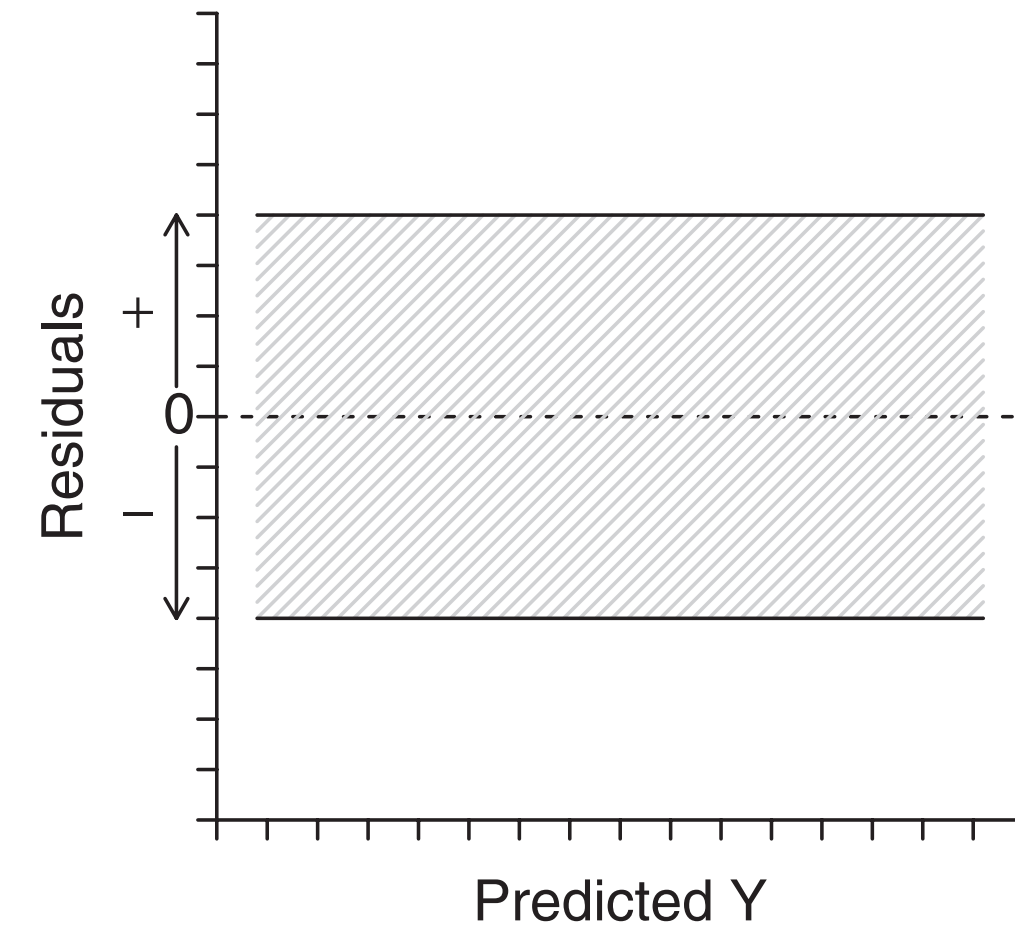(Check using histograms, boxplots, etc.)

**3. Variance of residuals is homogeneous across all values of *x***

(Check using residuals vs. predicted-*y* plot)

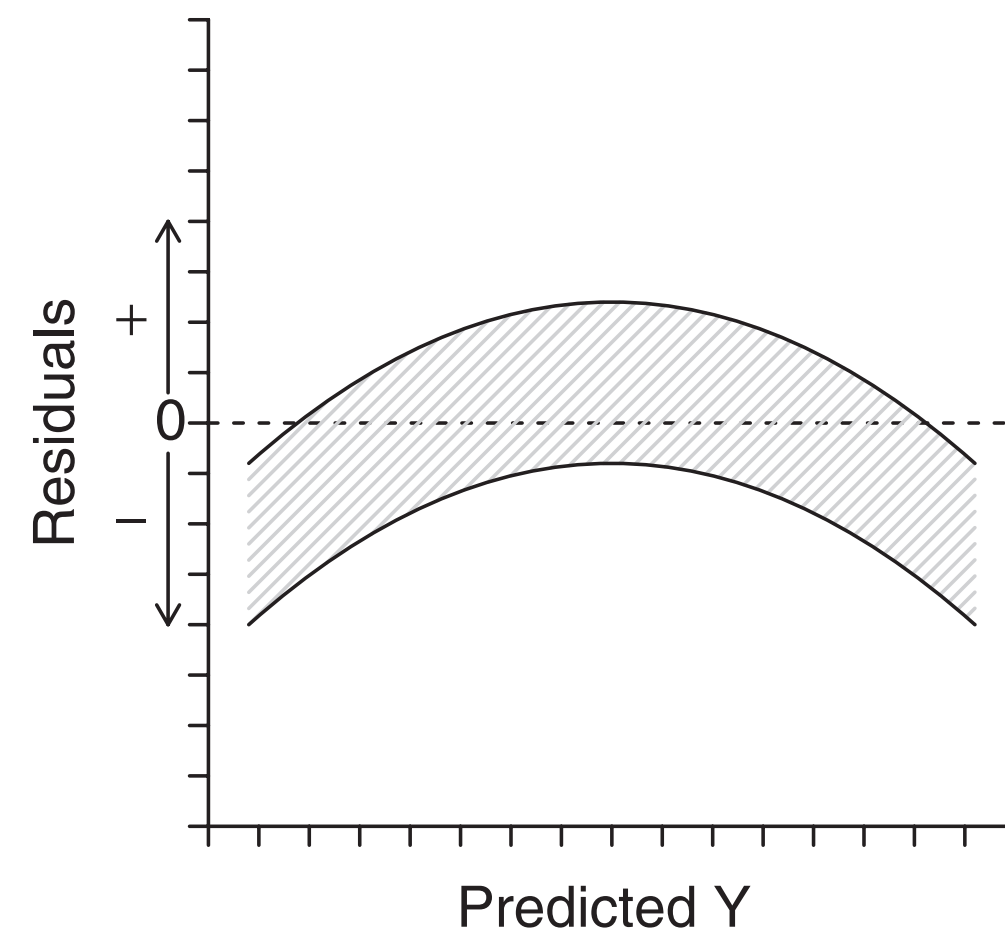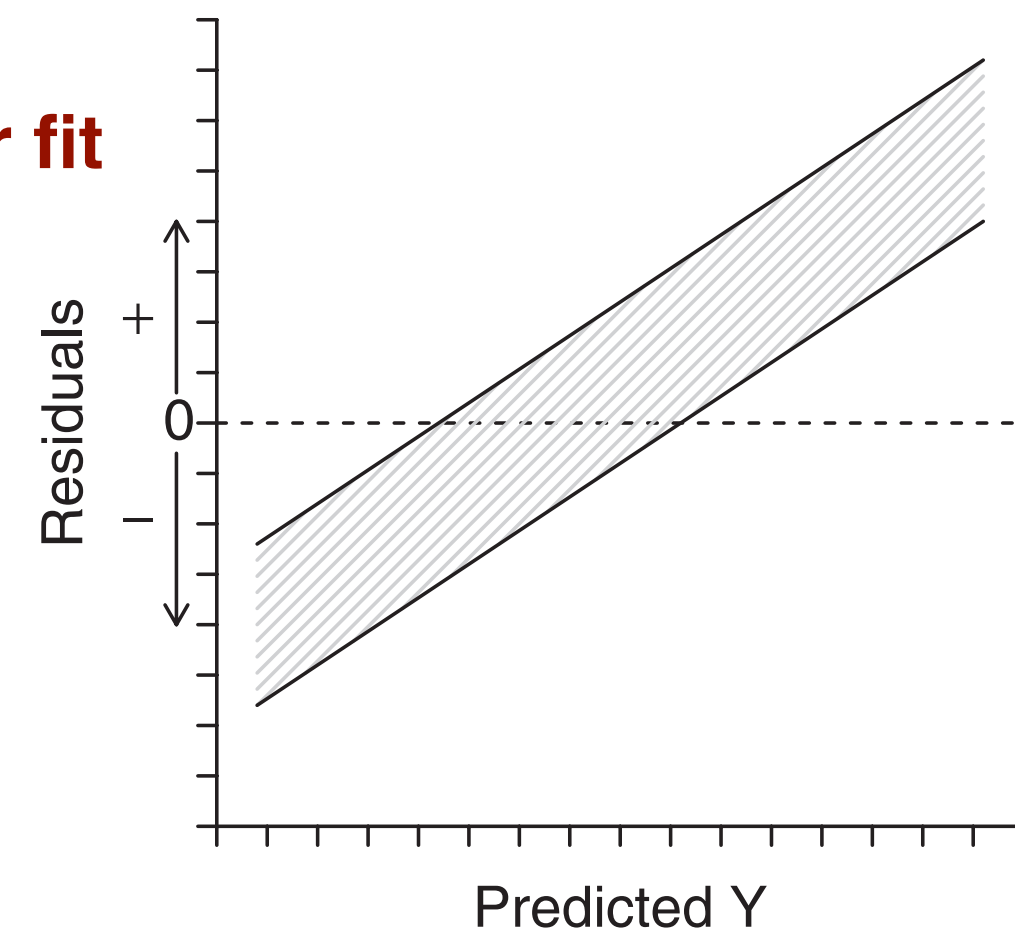# Assumptions of the *F*-ratio test for $\beta_1 = 0$
## (Using <u>residual plots</u>)

**uniform variance across pred. *Y***

**increasing variance as pred. *Y* increases**

**model is systematically poor fit (may be an outlier)**

**original relationship is curvilinear, not linear**
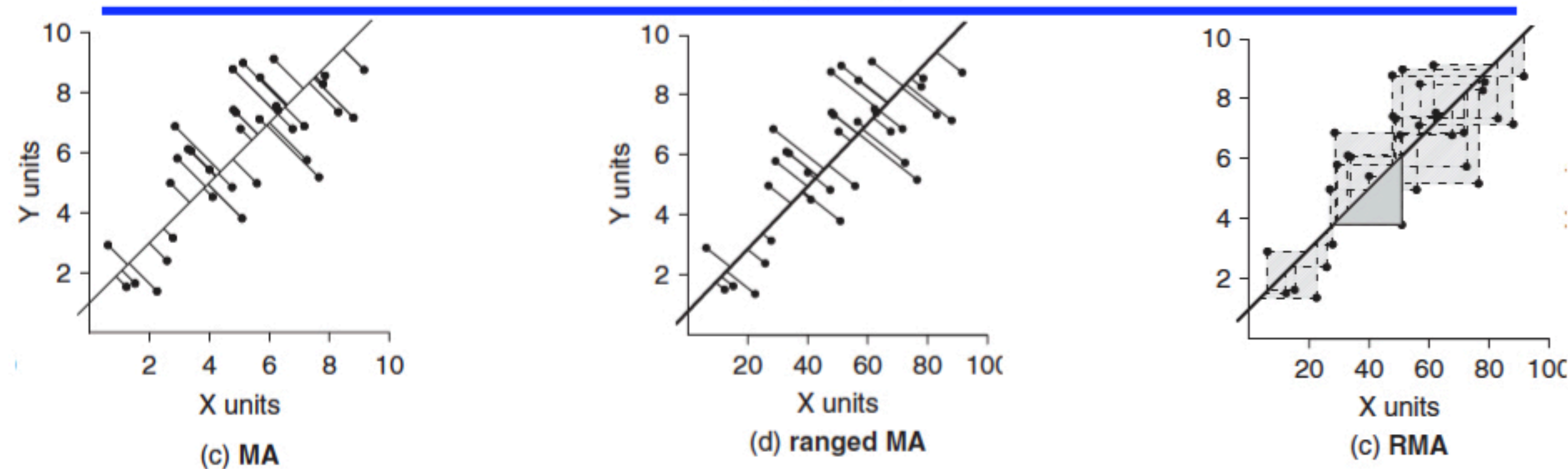
# The coefficient of determination ($r^2$)

$$r^2 = \frac{SS_{regression}}{SS_{total}} = 1 - \frac{SS_{residual}}{SS_{total}}$$

$$r^2 = 1 - \frac{SS_{residual(full)}}{SS_{total(reduced)}}$$

- proportion of variance in Y that is explained by X

# Model II regression: when *x* and *y* are both measured with error



(Type II)

(c) MA    (d) ranged MA    (c) RMA

**MA (Major Axis):** *x* and *y* have similar error and have same units

**Ranged MA:** *x* and *y* in different units or on different scales.

Assumes no outliers

**RMA or SMA (Reduced Major Axis):** *x* and *y* in different units or on different scales. Robust to outliers