# Linear Models

**Fernando Racimo, 2023**

# The two sides of statistics

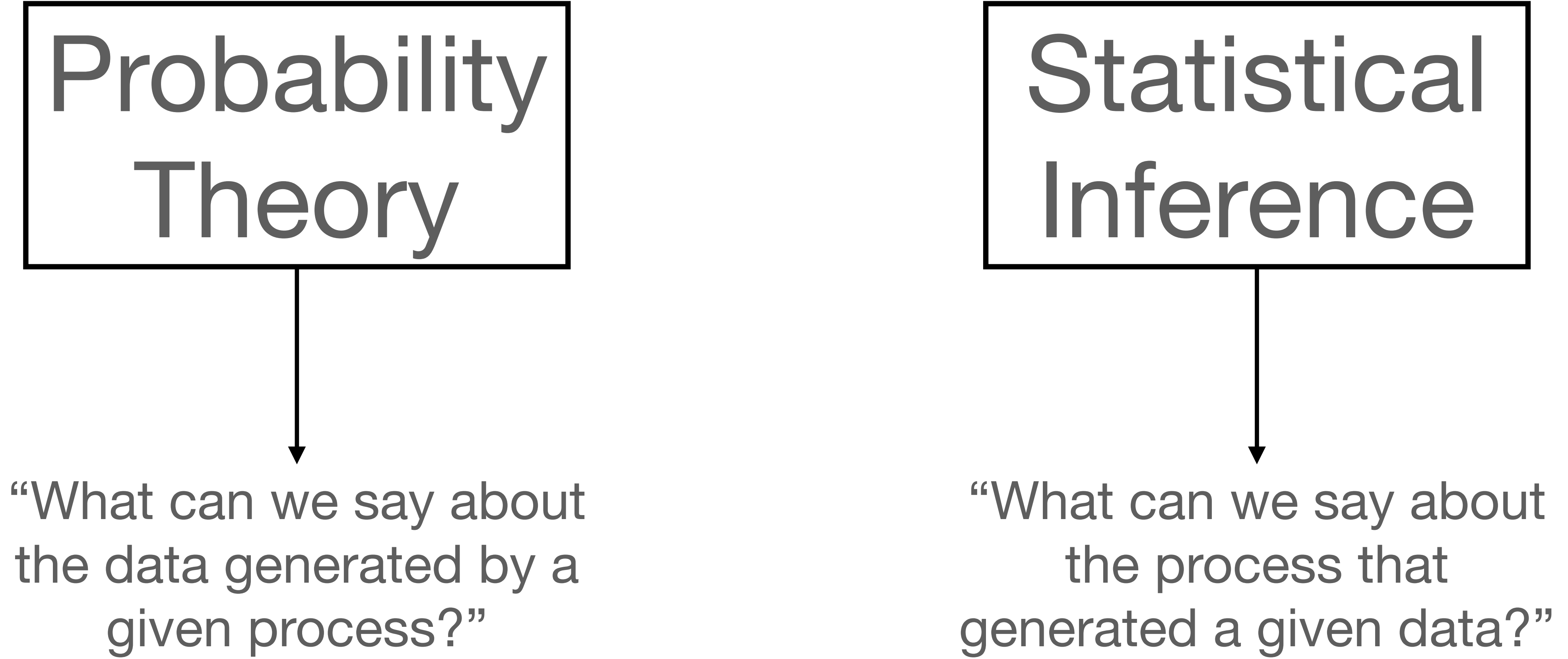Probability Theory

"What can we say about the data generated by a given process?"

# The two sides of statistics

| Probability Theory |
|:---:|

"What can we say about the data generated by a given process?"

| Statistical Inference |
|:---:|

"What can we say about the process that generated a given data?"

# The two sides of statistics
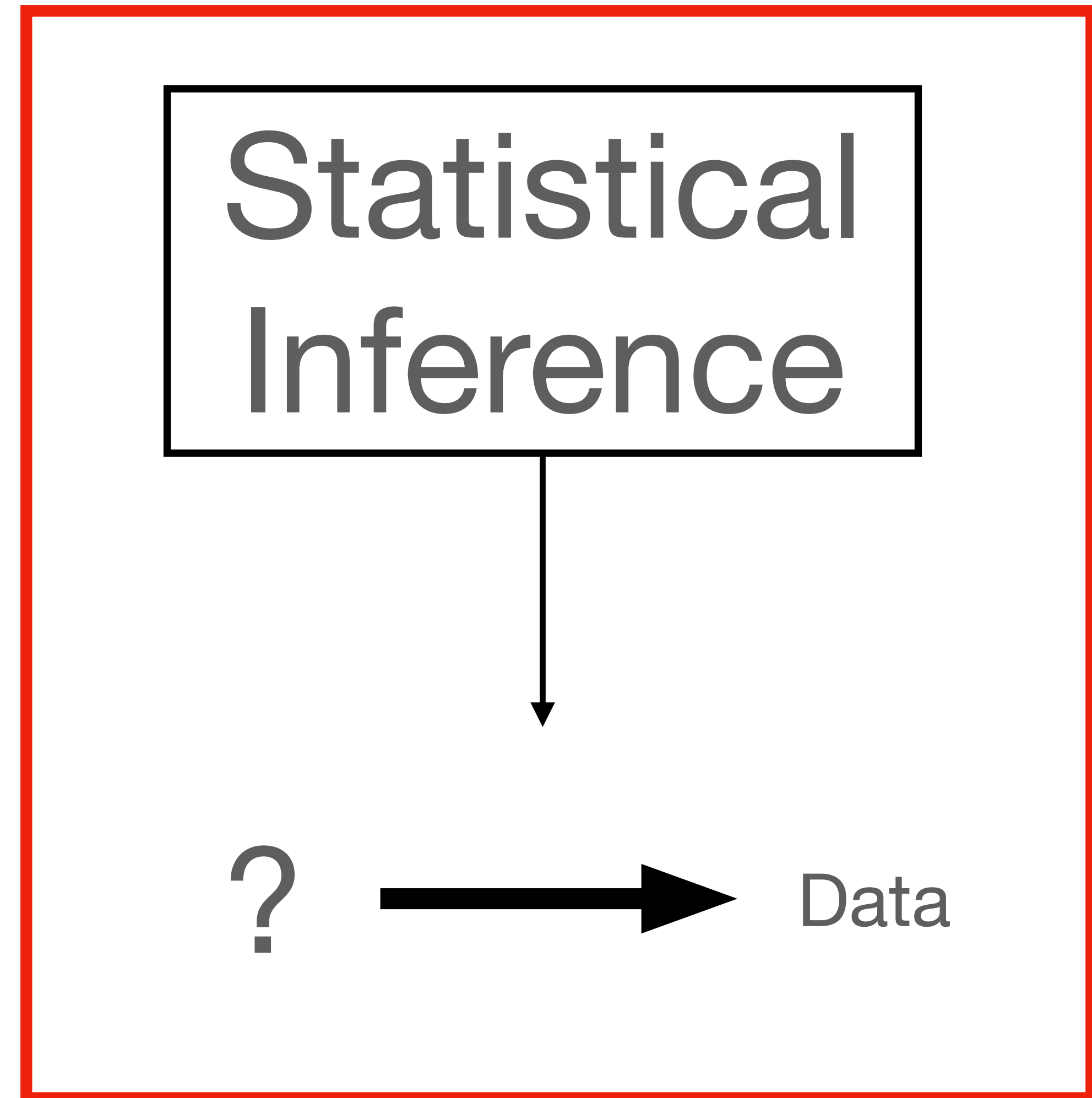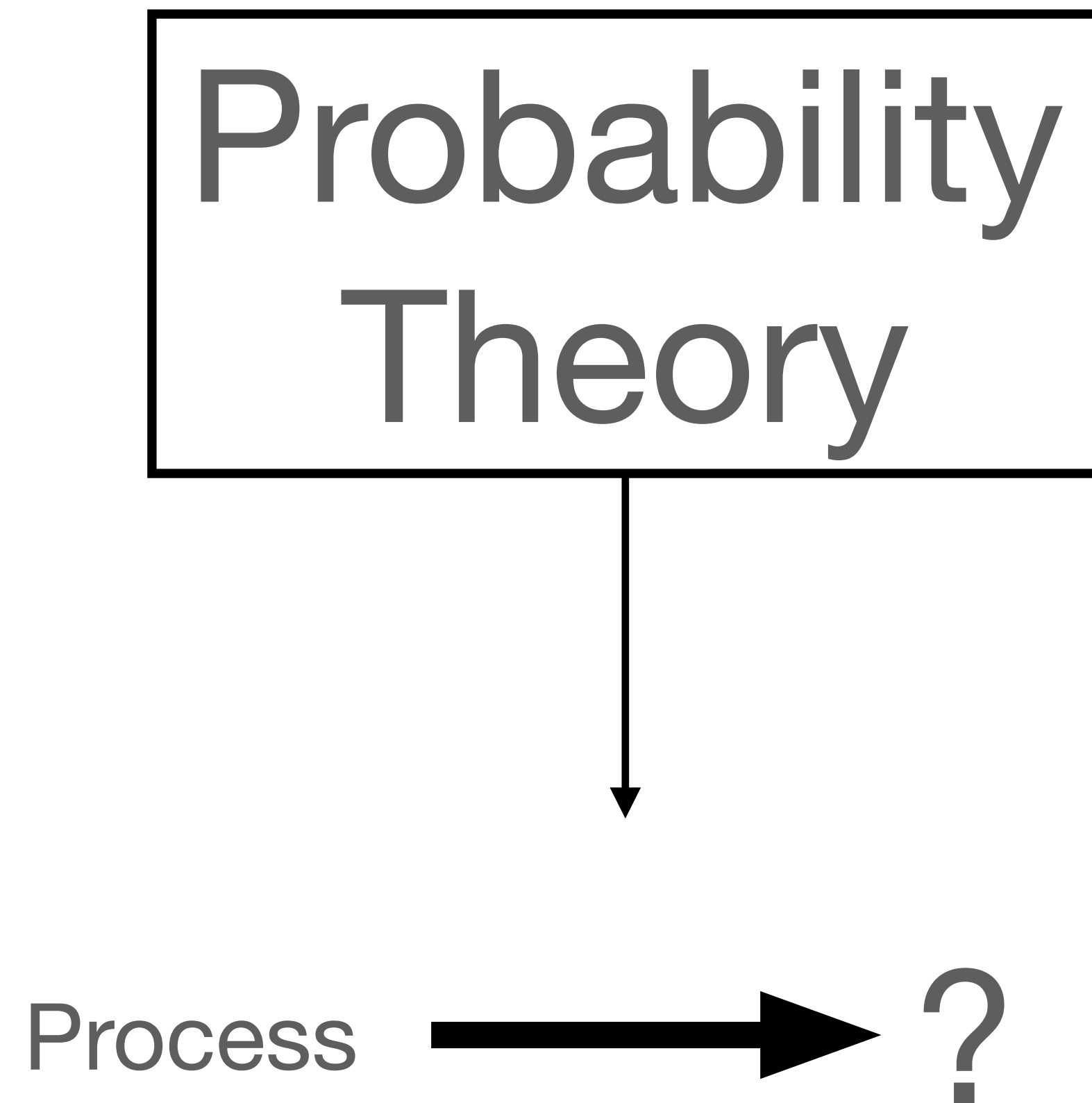
# The two sides of statistics

Probability
Theory

Process ➡ ?

Statistical
Inference

? ➡ Data

# Statistical inference: two "flavors"

Supervised learning: **today**

Predictor variables

Response variables

X $\xrightarrow{\quad ? \quad}$ Y

# Statistical inference: two "flavors"

Supervised learning: **today**

Predictor variables                    Response variables

$$X \xrightarrow{\quad ? \quad} Y$$

Unsupervised learning: **tomorrow**
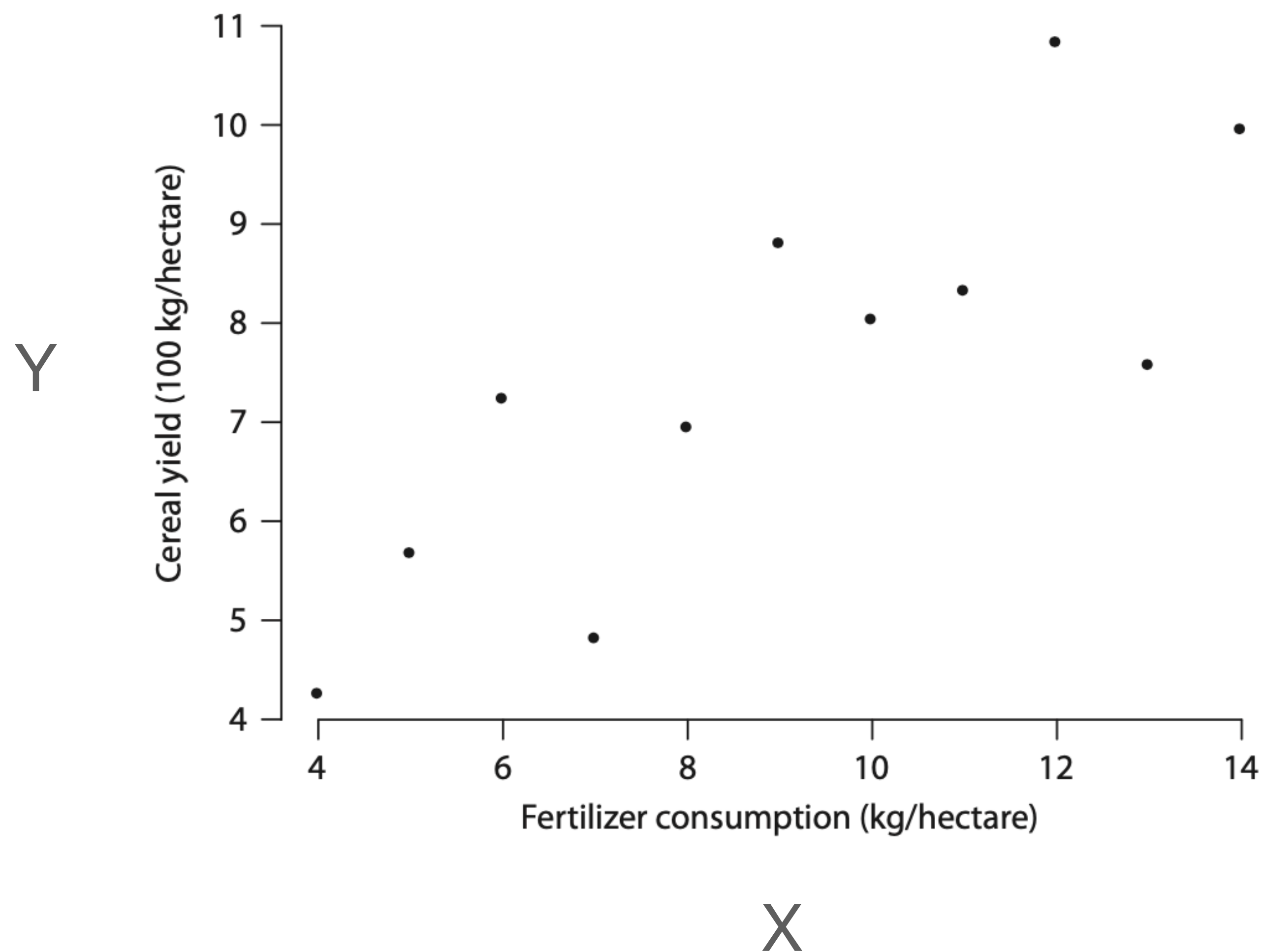
Variables

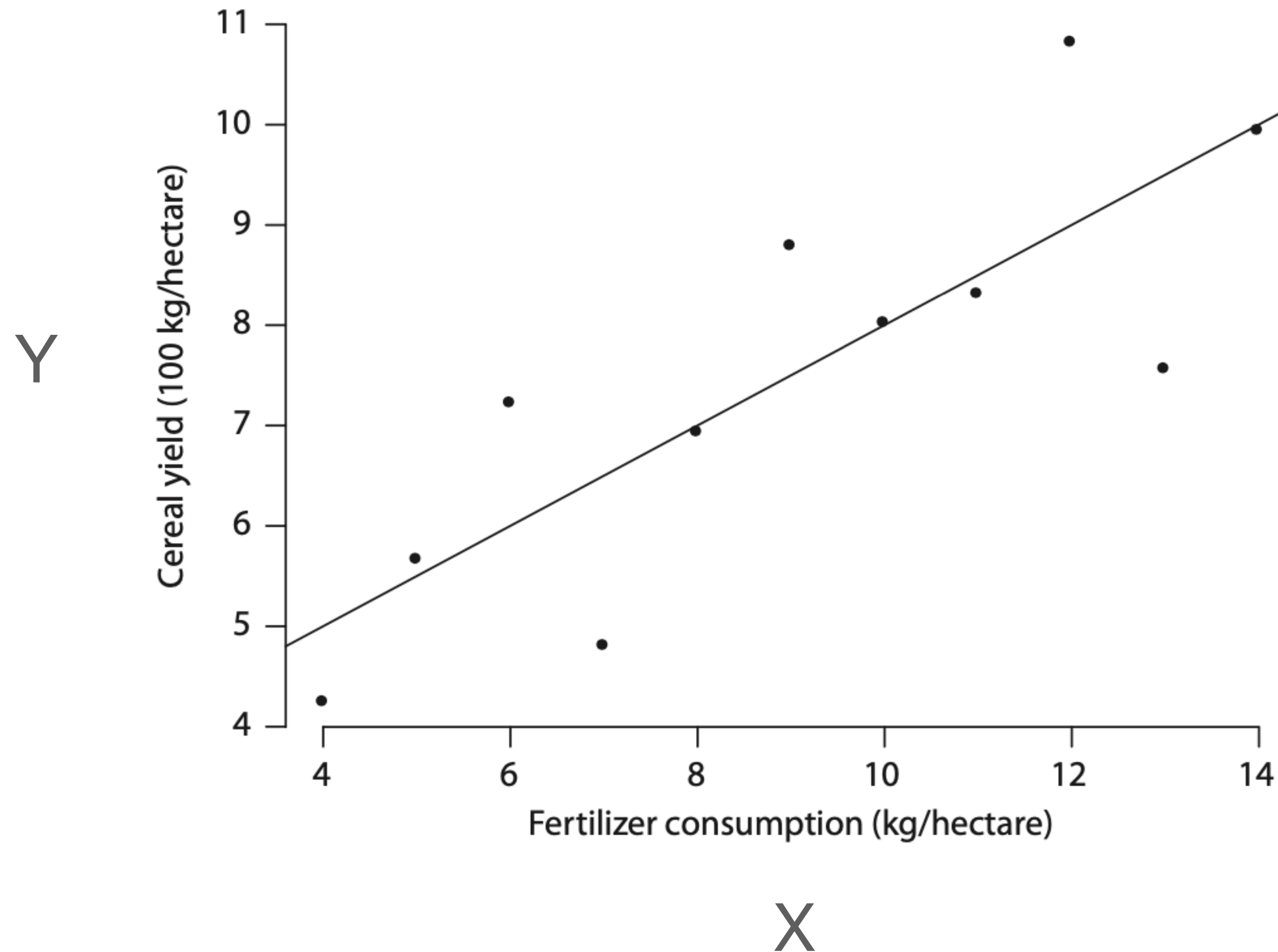$$? \xrightarrow{\quad ? \quad} Y$$

# Linear Regression

- Simplest model in supervised learning

- Jumping-off point for more complex models

- Many statistical models are extensions or generalizations of the linear model
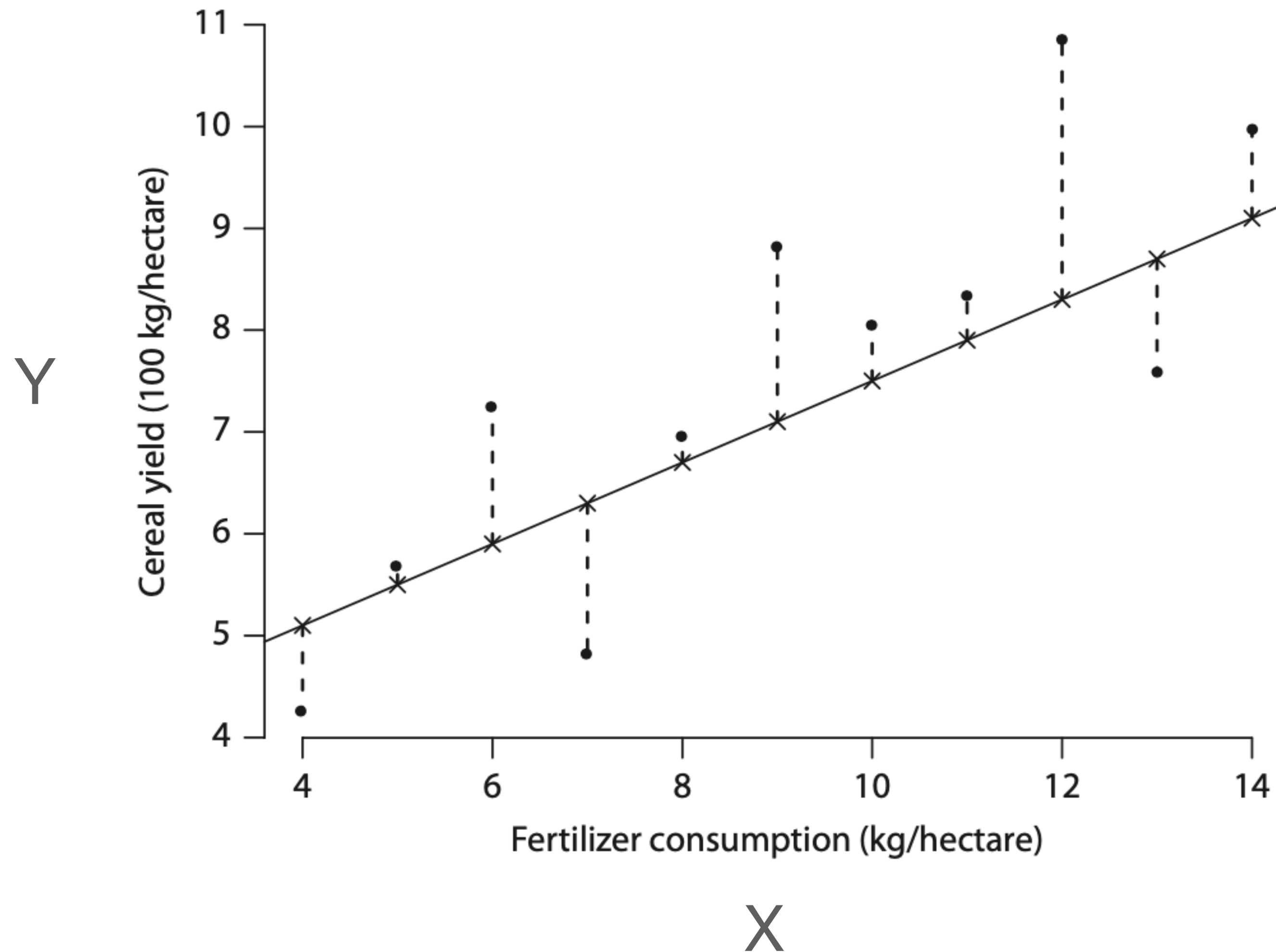
# Linear Regression

# Linear Regression: "a line of best fit"

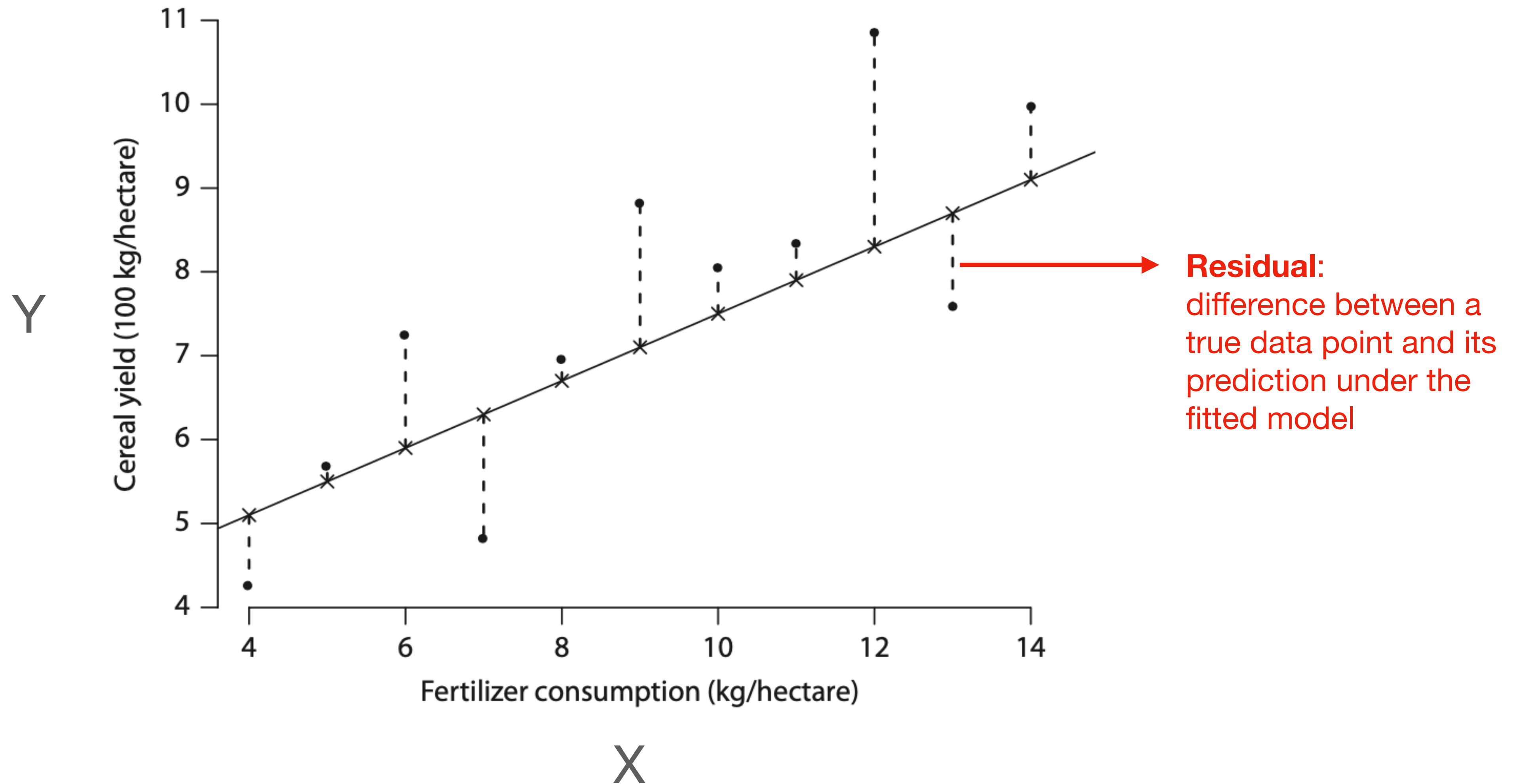# Linear Regression: types of questions

- Is there a relation between variable X and variable Y?

- How strong is the relation between variable X and variable Y?

- How well can we predict variable Y from the values of variable X?

- Is the relationship between variable X and Y linear?

- Which variables contribute to variable Y?

# Linear Regression: "a line of best fit"

# Linear Regression: "a line of best fit"



Y

Cereal yield (100 kg/hectare)

Fertilizer consumption (kg/hectare)

X

**Residual**: difference between a true data point and its prediction under the fitted model

# Simple linear regression

- **1** predictor variable (x)

- **1** response variable (y)

# Simple linear regression

- **1** predictor variable (x)

- **1** response variable (y)

$$\mathbf{y} = \mathbf{f(x)}$$

Variable **y** is a **function** of **x**

# Simple linear regression

- **1** predictor variable (x)

- **1** response variable (y)

$$\mathbf{y} = \mathbf{f}(\mathbf{x})$$

Variable **y** is a **function** of **x**

$$\mathbf{y} \approx \beta_0 + \beta_1 \mathbf{x}$$

Variable **y** is a **linear** function of variable **x**

# Simple linear regression

- **1** predictor variable (x)

- **1** response variable (y)

$$\mathbf{y} = \mathbf{f}(\mathbf{x})$$

Variable **y** is a **function** of **x**

$$\mathbf{y} \approx \beta_0 + \beta_1 \mathbf{x}$$

Variable **y** is a **linear** function of variable **x**

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$$

Variable **y** is a linear function of **x**, plus some **noise**

# Simple linear regression

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$$

Variable **y** is a linear function of **x**, plus some **noise**

# Simple linear regression

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$$

Variable **y** is a linear function of **x**, plus some **noise**

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

Each value $y_i$ has a specific predictor value $x_i$ and noise value $\epsilon_i$

# Simple linear regression

$$\mathbf{y} = \beta_0 + \beta_1\mathbf{x} + \epsilon$$

Variable **y** is a linear function of **x**, plus some **noise**

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

Each value $y_i$ has a specific predictor value $x_i$ and noise value $\epsilon_i$

**Individual values are represented with subscripts**

# Simple linear regression

**Vectors of values (variables) are represented in bold**

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$$

Variable **y** is a linear function of **x**, plus some **noise**

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

Each value $y_i$ has a specific predictor value $x_i$ and noise value $\epsilon_i$

**Individual values are represented with subscripts**

# Parameters vs. Estimates

$$y = \boxed{\beta_0} + \boxed{\beta_1}x + \epsilon$$

$\beta_0$ and $\beta_1$ are unknown parameters in our model: **we do not know their value**

$$y = \boxed{\hat{\beta}_0} + \boxed{\hat{\beta}_1}x + \epsilon$$

$\hat{\beta}_0$ and $\hat{\beta}_1$ are our **best estimates** of the above parameters

# Model vs. Inference Method

Model

Inference method

Parameter estimates



$$\hat{\alpha} = 0.56$$

$$\hat{\beta} = 3.2$$

$$\hat{\mu} = -2$$

$$\cdots$$

$v$

$\eta$

$\gamma$

$\mu$

$\beta$

Data

# Model vs. Inference Method

Model                          Inference method

Linear regression   ⟶   Ordinary least squares

Weighted least squares

Maximum likelihood

Ridge regression

…

# Model vs. Inference Method

Model

Inference method

Linear regression

Ordinary least squares

Weighted least squares

Maximum likelihood

Ridge regression

…

# Ordinary least squares

"Find estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ of the parameters $\beta_0$ and $\beta_1$ by **minimizing the Sum of Squared Residuals**"

# Ordinary least squares

"Find estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ of the parameters $\beta_0$ and $\beta_1$ by **minimizing the Sum of Squared Residuals**"

?

# Ordinary least squares: what does this mean?



$$\text{res} = y_i - \hat{y}$$

# Ordinary least squares: what does this mean?



Residuals

$$\text{res} = y_i - \hat{y}$$

$$(y_i - \hat{y})^2$$

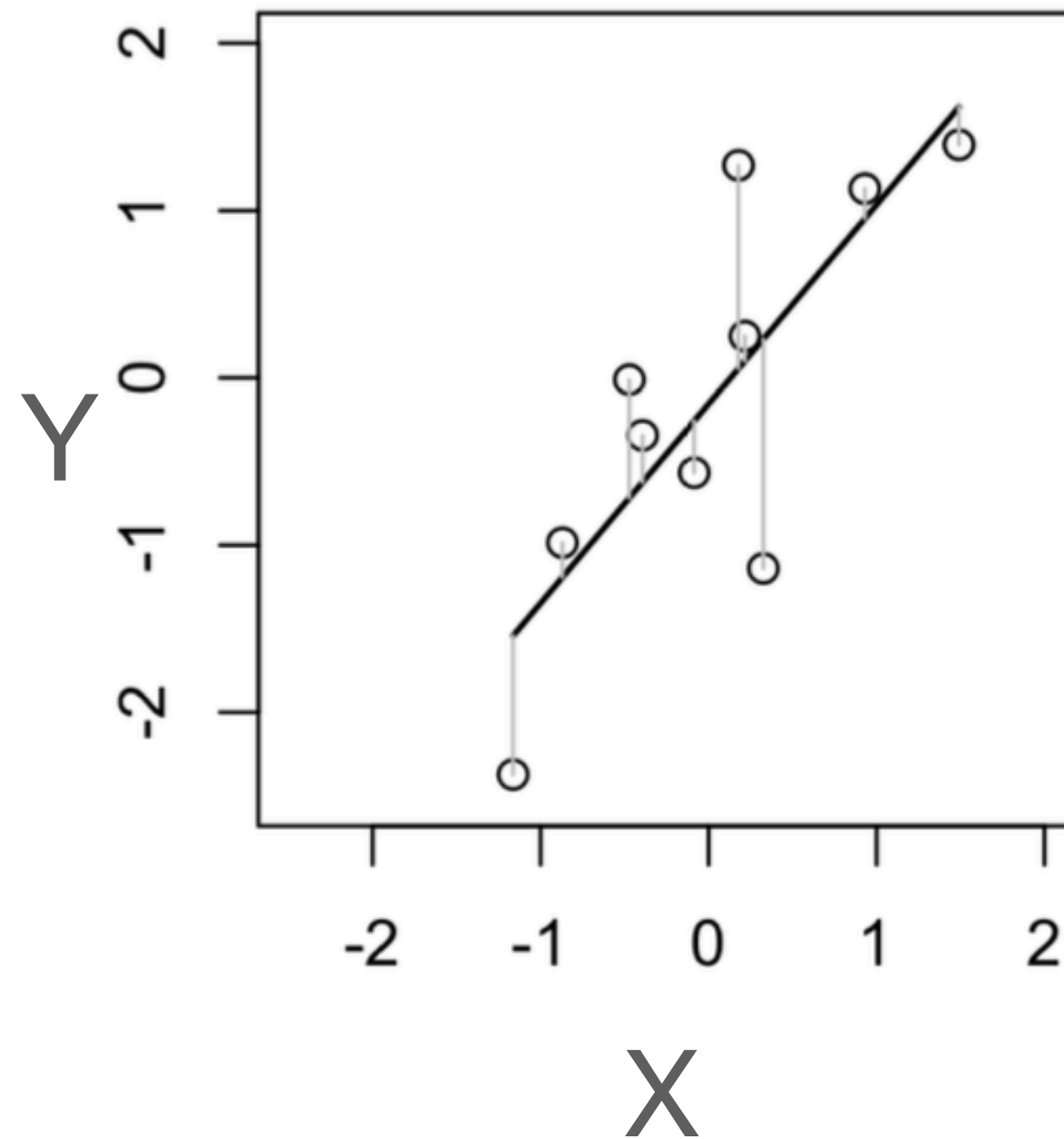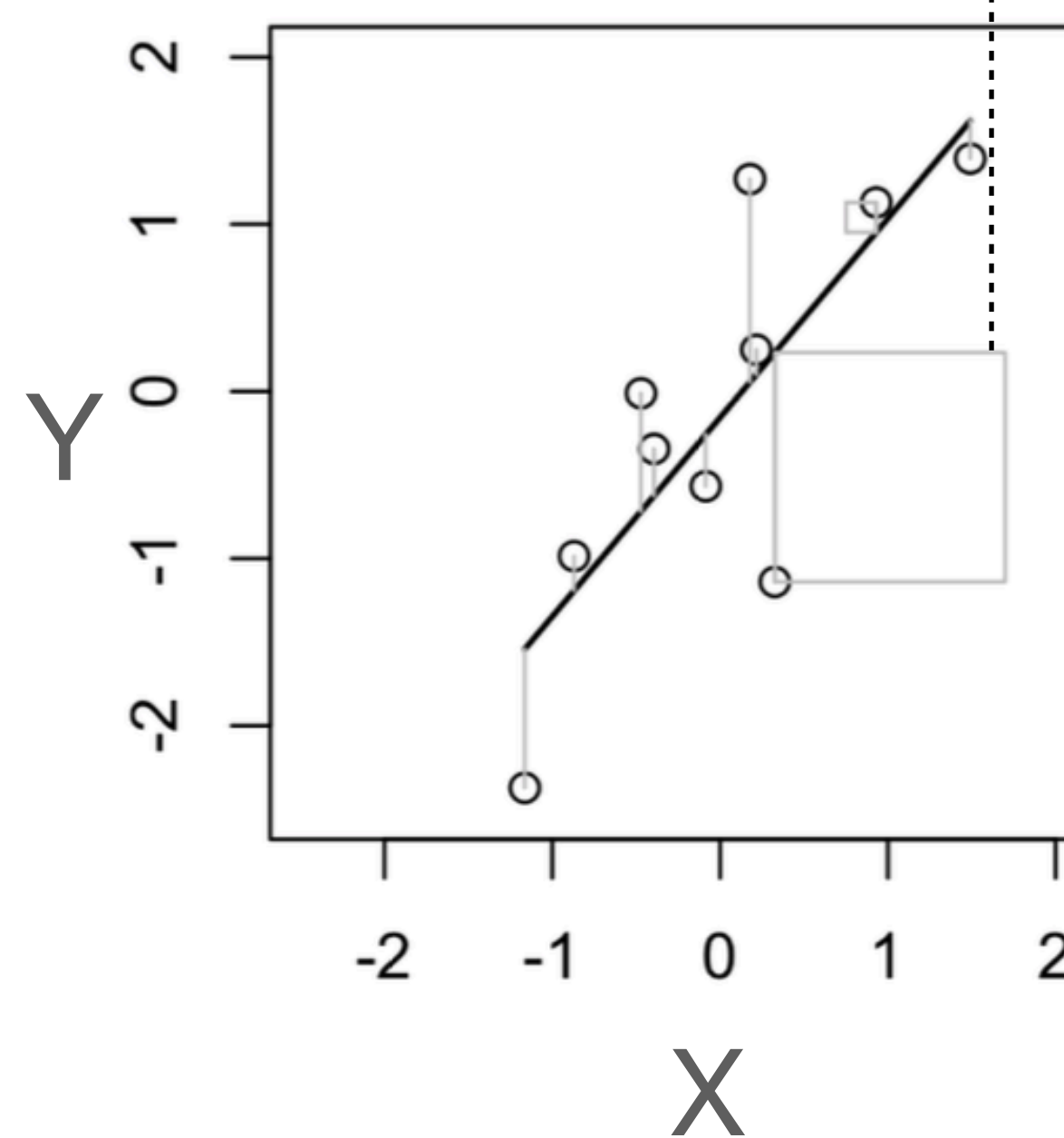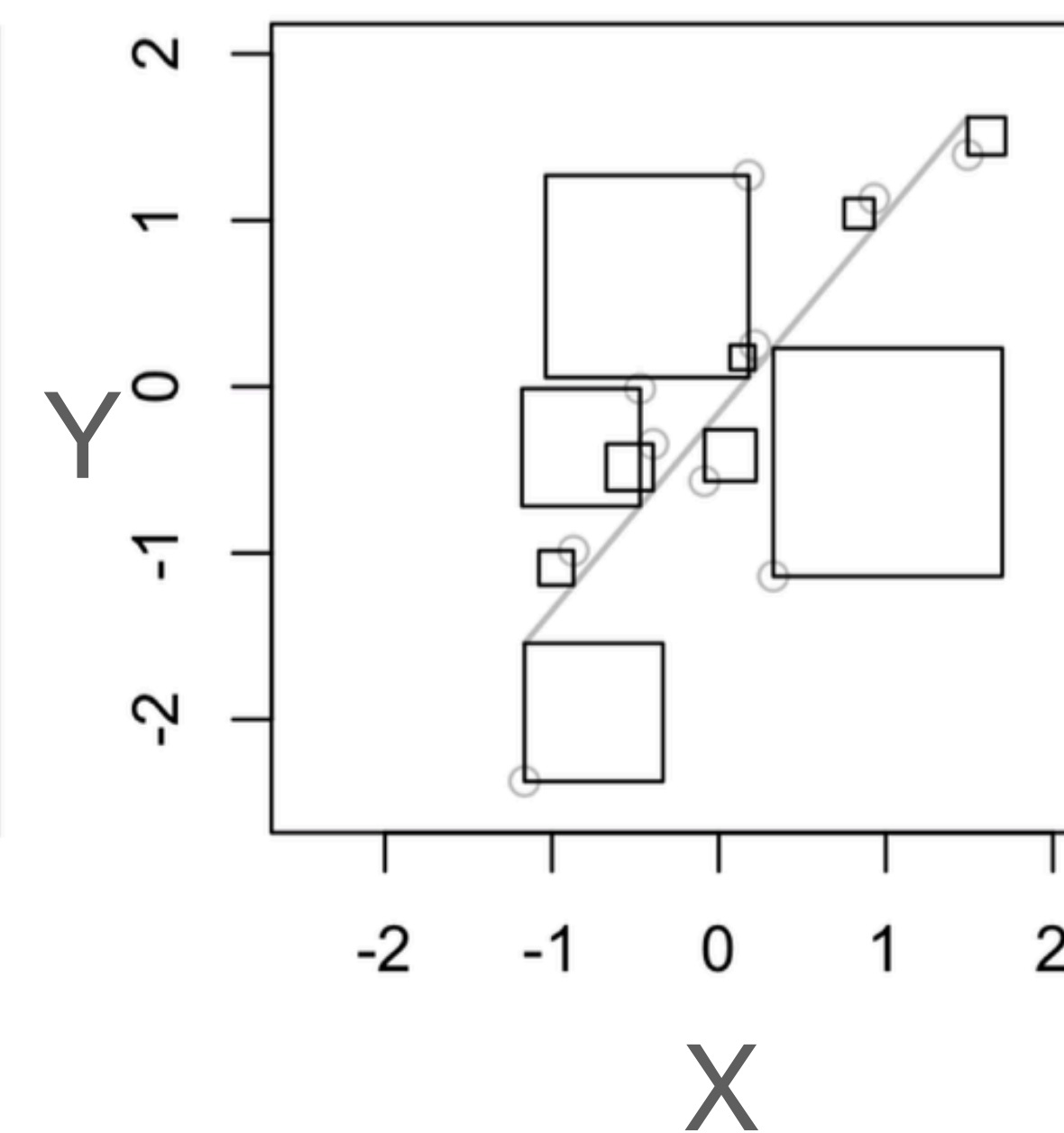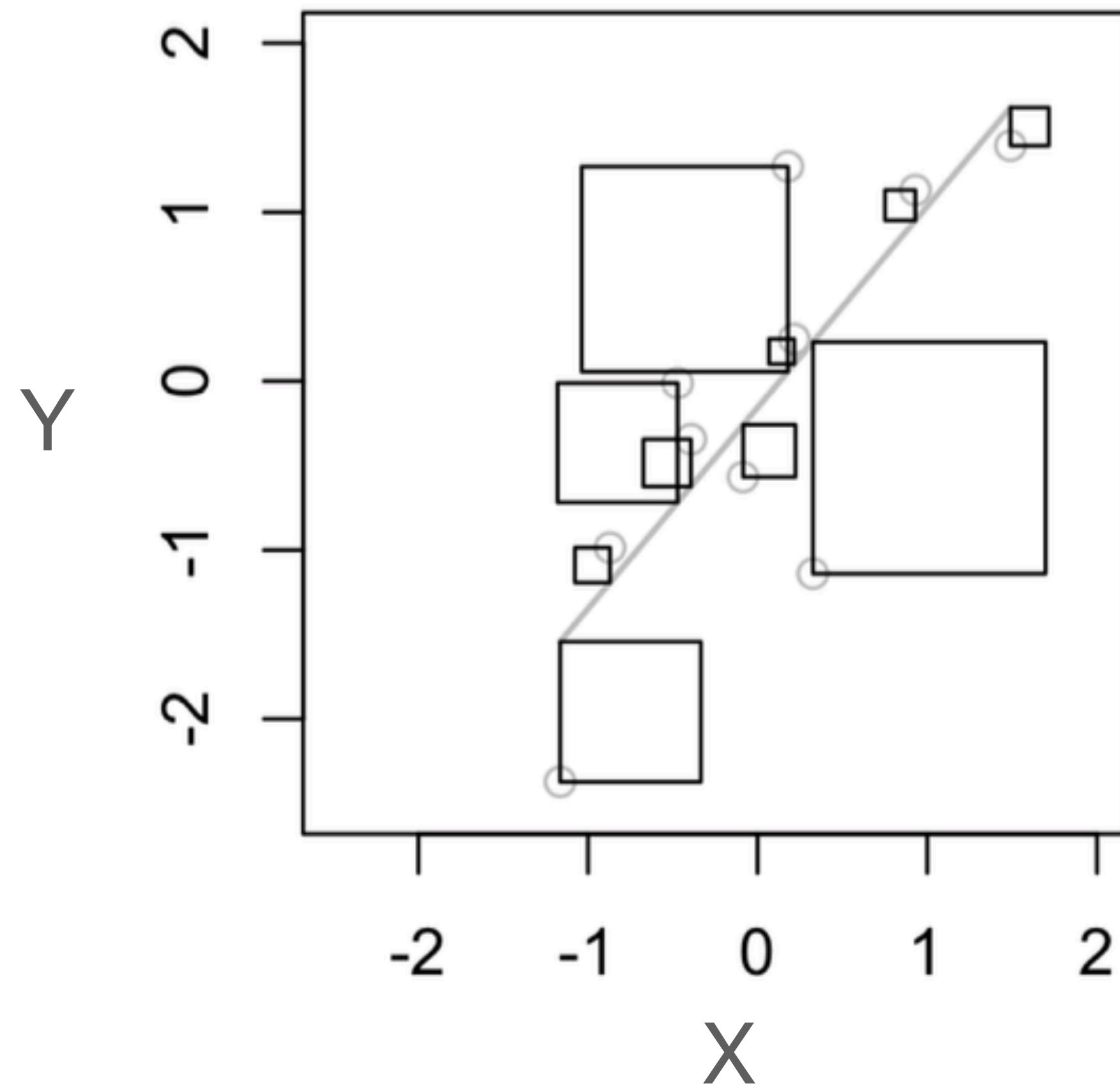# Ordinary least squares: what does this mean?

Residuals

A squared residual



res = $y_i - \hat{y}$

$(y_i - \hat{y})^2$

# Ordinary least squares: what does this mean?



**Residuals**

**A squared residual**

**Sum of Squared Residuals**

$$\text{res} = y_i - \hat{y}$$

$$(y_i - \hat{y})^2$$

$$SS_{res} = \sum_{i}^{n}(y_i - \hat{y})^2$$

# Ordinary least squares: what does this mean?

**Sum of Squared Residuals**



$$SS_{res} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

# Ordinary least squares: what does this mean?
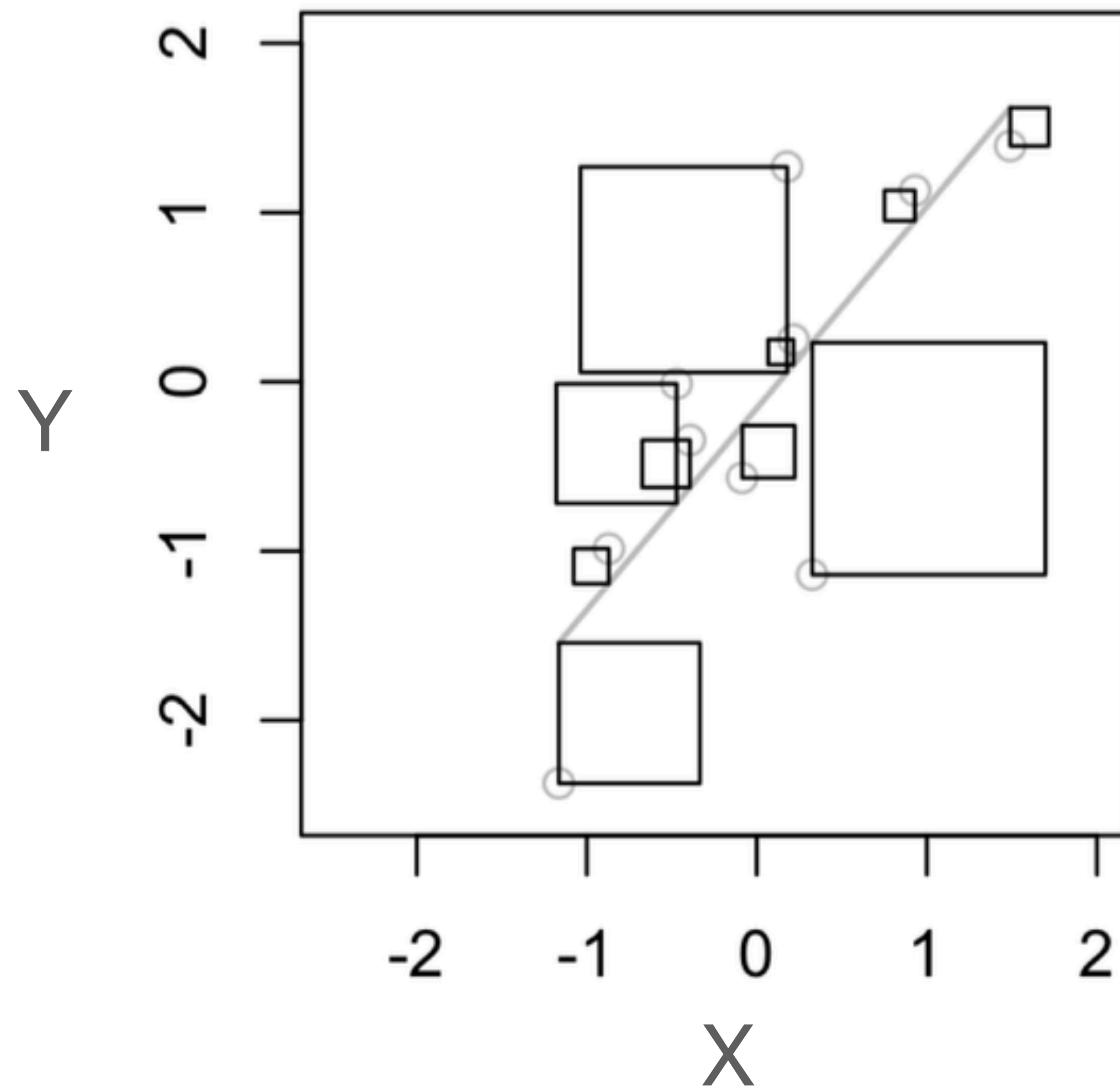
**Sum of Squared Residuals**



$$SS_{res} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

but our estimate $\hat{y}_i$ is simply a linear function of $x_i$:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

# Ordinary least squares: what does this mean?

$$SS_{res} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

**Sum of Squared Residuals**



but our estimate $\hat{y}_i$ is simply a linear function of $x_i$:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$SS_{res} = \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

# Ordinary least squares: what does this mean?

$$SS_{res} = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$$

The values of $\beta_0$ and $\beta_1$ that **minimize SSres** are:

# Ordinary least squares: what does this mean?

$$SS_{res} = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$$

The values of $\beta_0$ and $\beta_1$ that **minimize SSres** are:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

# Ordinary least squares: what does this mean?

$$SS_{res} = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$$

The values of $\beta_0$ and $\beta_1$ that **minimize SSres** are:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Average of y      Average of x

# Ordinary least squares: what does this mean?

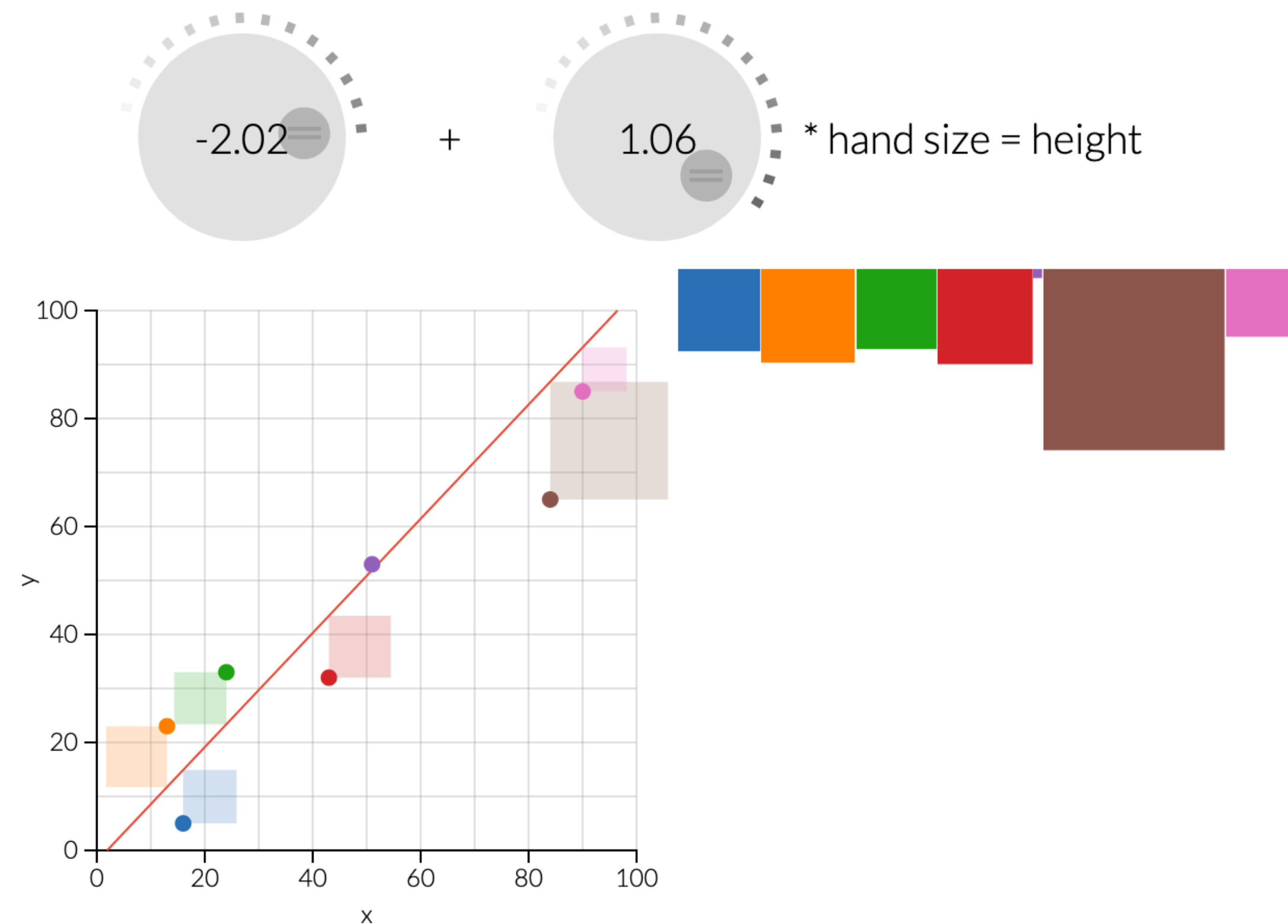$$SS_{res} = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$$

The values of $\beta_0$ and $\beta_1$ that **minimize SSres** are:

$$\hat{\beta}_0 = \overline{y} - \hat{\beta}_1 \overline{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^{n} [(x_i - \overline{x})^2]}$$

Average of y     Average of x

# Ordinary least squares: what does this mean?

$$SS_{res} = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$$

The values of $\beta_0$ and $\beta_1$ that **minimize SSres** are:

$$\hat{\beta}_0 = \overline{y} - \hat{\beta}_1 \overline{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^{n} [(x_i - \overline{x})^2]}$$

Average of y

Average of x

**You'll have to trust me here… or not:**
for proof, try taking the derivative of SSres and equalizing it to 0.

# Ordinary least squares: interactive session



-2.02 + 1.06 * hand size = height

https://setosa.io/ev/ordinary-least-squares-regression/

# Ordinary least squares

"Find estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ of the parameters $\beta_0$ and $\beta_1$ by **minimizing the Sum of Squared Residuals**"
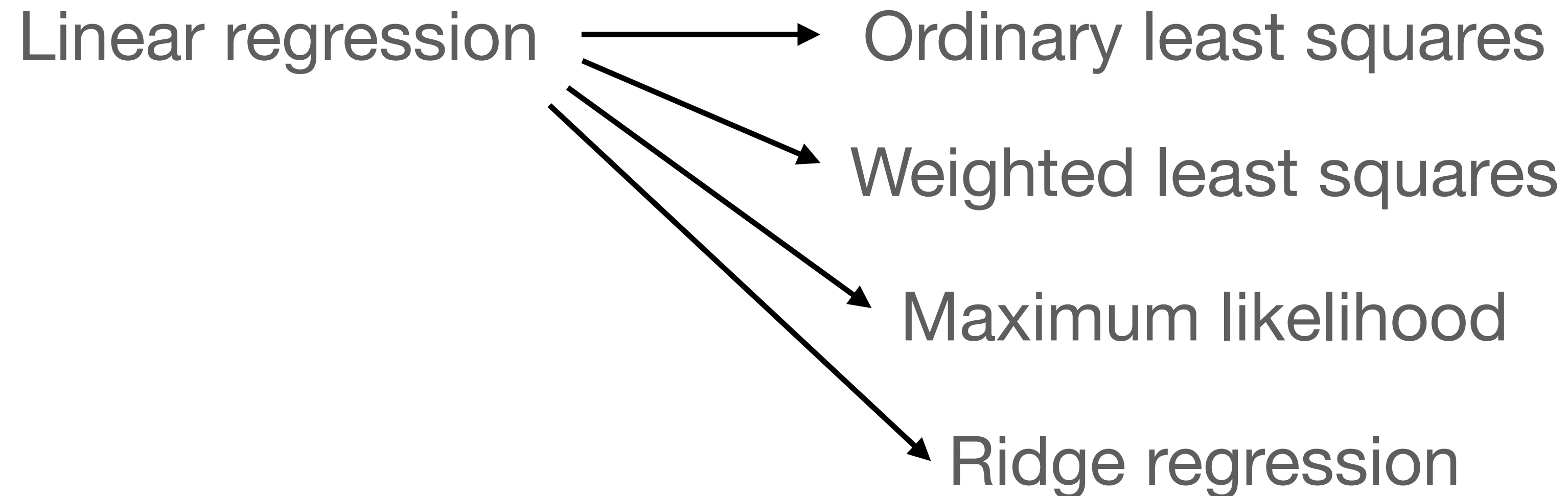
# Ordinary least squares

"Find estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ of the parameters $\beta_0$ and $\beta_1$ by **minimizing the Sum of Squared Residuals**"

Loss function

# Model vs. Inference Method

Model

Inference method

Linear regression → Ordinary least squares

Weighted least squares

Maximum likelihood

Ridge regression

…

Each of these methods has a different loss function!

# What does the slope mean?

$$\hat{\beta}_1 = ?$$

# What does the slope mean?

$$\hat{\beta}_1 = ?$$

The **slope** measures the **covariance between X and Y**, as a proportion of the **variance of X**

# What does the slope mean?

$$\hat{\beta}_1 = ?$$

The **slope** measures the **covariance between X and Y**, as a proportion of the **variance of X**

Let's unpack this a bit…

# Measuring variation

# Measuring variation

**Sum of squares of X**

$$SS_X = \sum_{i=1}^{n} (x_i - \bar{x})^2$$

# Measuring variation

**Sum of squares of X**

$$SS_X = \sum_{i=1}^{n} (x_i - \bar{x})^2$$

**Sample variance of X:**

$$s_X = \frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n - 1}$$

# Measuring variation

**Sum of squares of X**

$$SS_X = \sum_{i=1}^{n} (x_i - \bar{x})^2$$

**Sample variance of X:**

$$s_X = \frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n - 1}$$

This constant is applied to the sum of squares to obtain an **unbiased estimate of the true variance** when we only have **finite samples**

# Measuring variation

**Sample variance of X:**

$$s_X = \frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n - 1}$$

# Measuring variation

## Sample variance of X:

$$s_X = \frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n-1}$$



"Sum of squares of X"



"Sum of squares of Y"

# Measuring variation

## Sample variance of X:

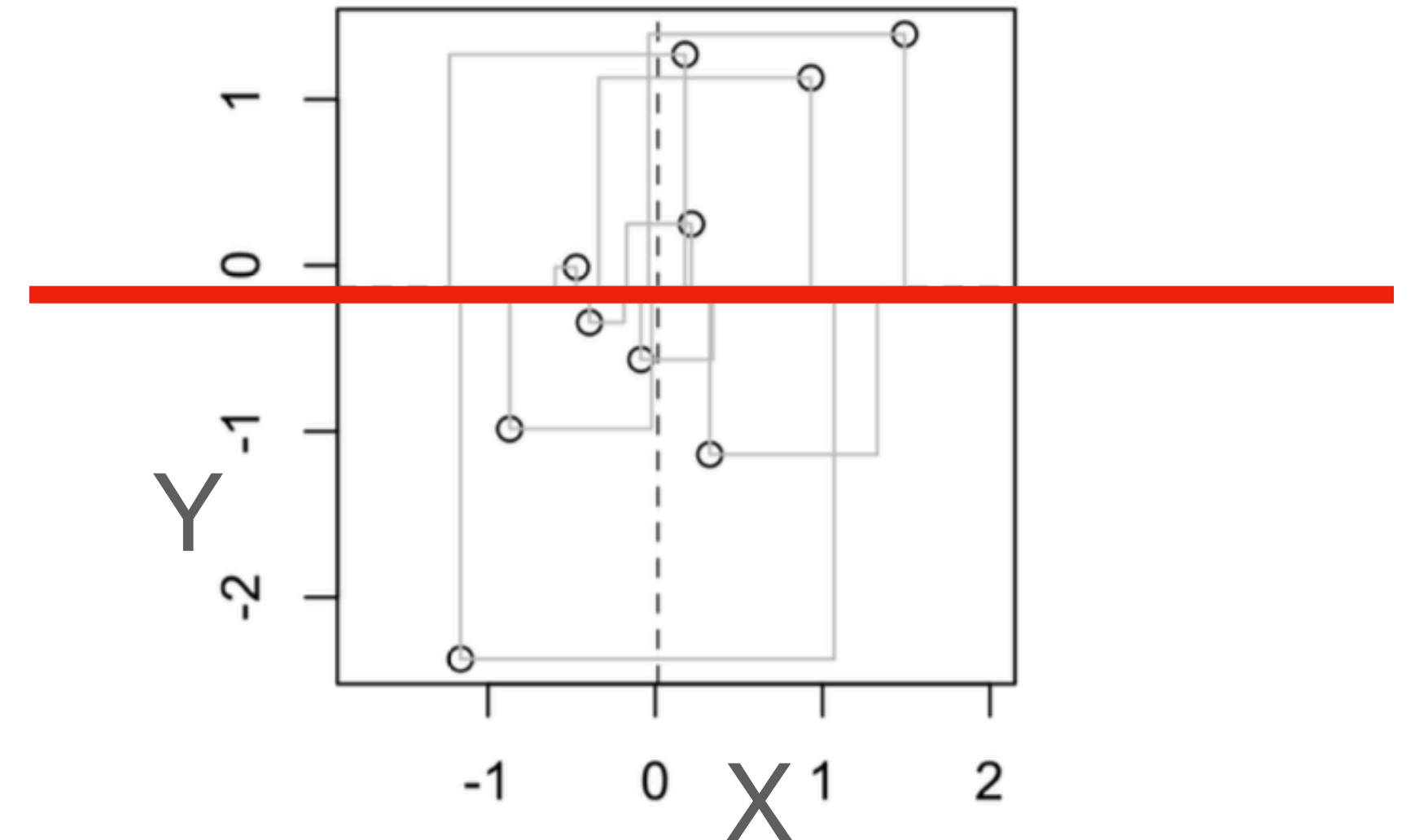$$s_X = \frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n - 1}$$

## Sample variance of Y:

$$s_Y = \frac{\sum_{i=1}^{n} (y_i - \bar{y})^2}{n - 1}$$



"Sum of squares of X"

"Sum of squares of Y"

# Measuring co-variation

**Sample covariance of X and Y:**

$$s_{X,Y} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

"Sum of XY rectangles"

# The Slope

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n} [(x_i - \bar{x})^2]}$$

Sample variance of X

# The Slope

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n} [(x_i - \bar{x})^2]}$$
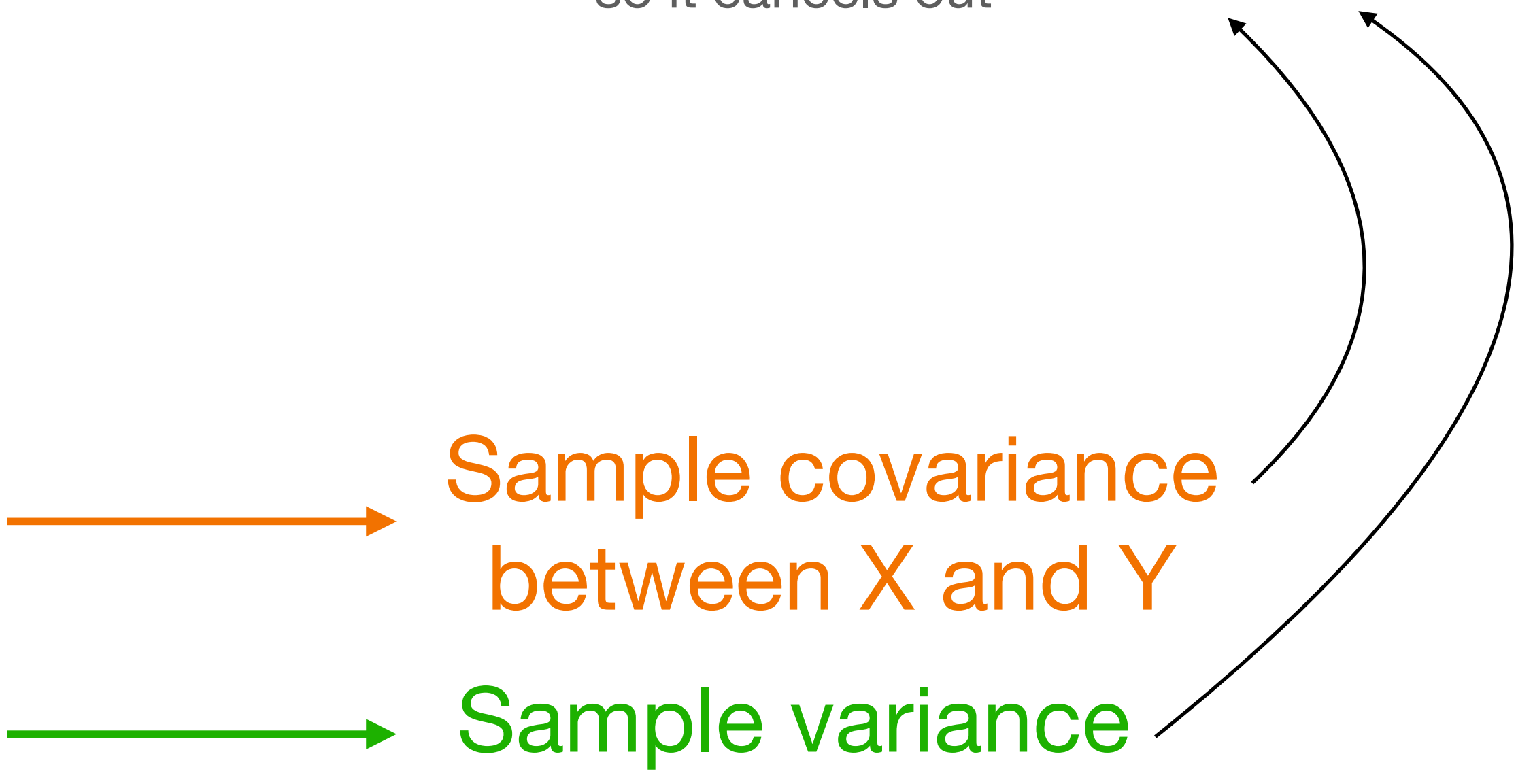
Sample covariance between X and Y
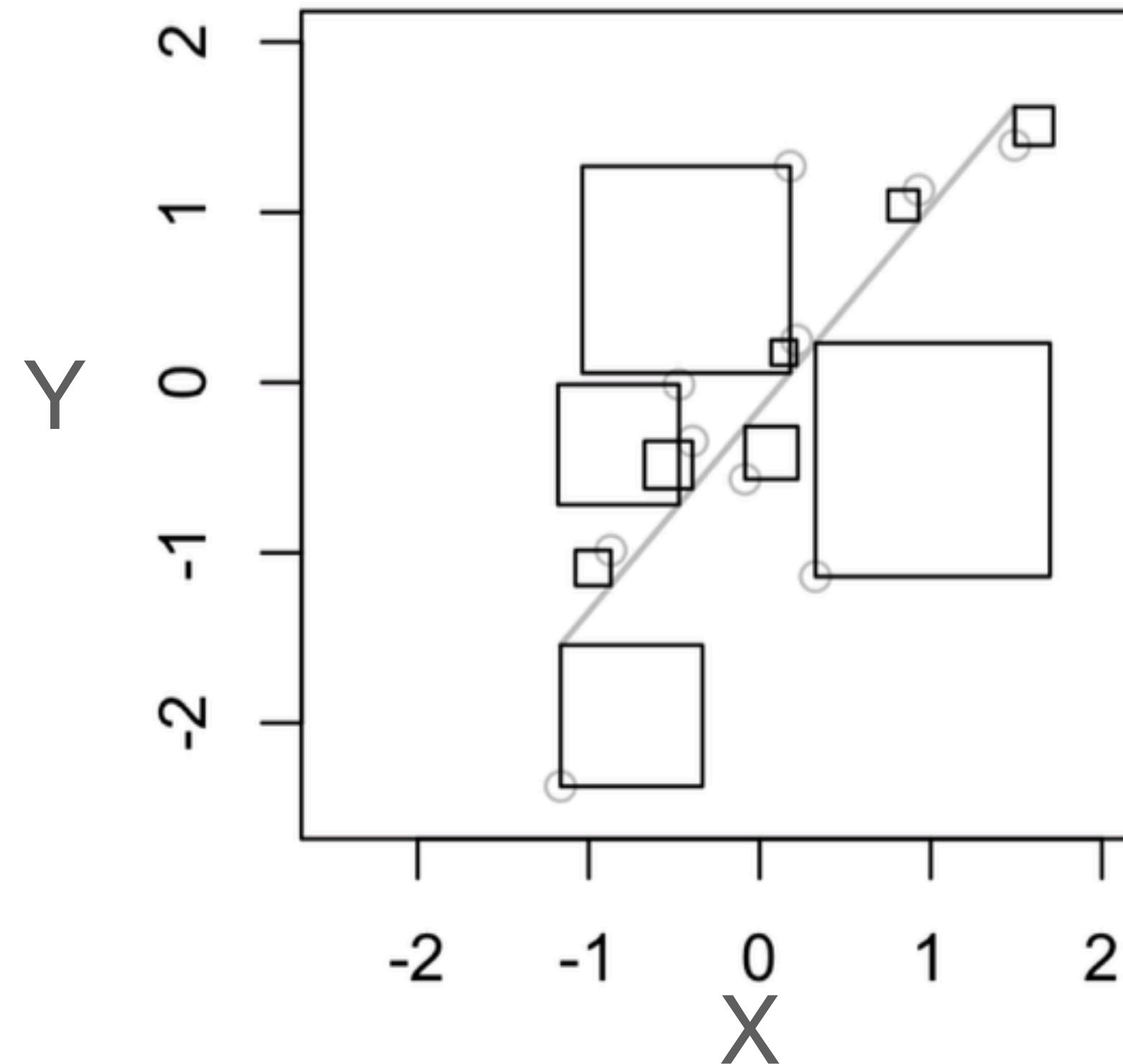
Sample variance of X

# The Slope

The constant n-1 is in both the numerator and denominator, so it cancels out

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^{n}[(x_i - \overline{x})^2]}$$

Sample covariance between X and Y

Sample variance of X

# The Slope

The constant n-1 is in both the numerator and denominator, so it cancels out

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}[(x_i - \bar{x})^2]}$$

→ Sample covariance between X and Y

→ Sample variance of X

The **slope** measures the **covariance between X and Y**, as a proportion of the **variance of X**

# How good is our model?

- No model is perfect, but some are more useful than others

- After fitting our model, we still have **unexplained variation in the dependent variable**: the sums of squares of the residuals (SS_res).
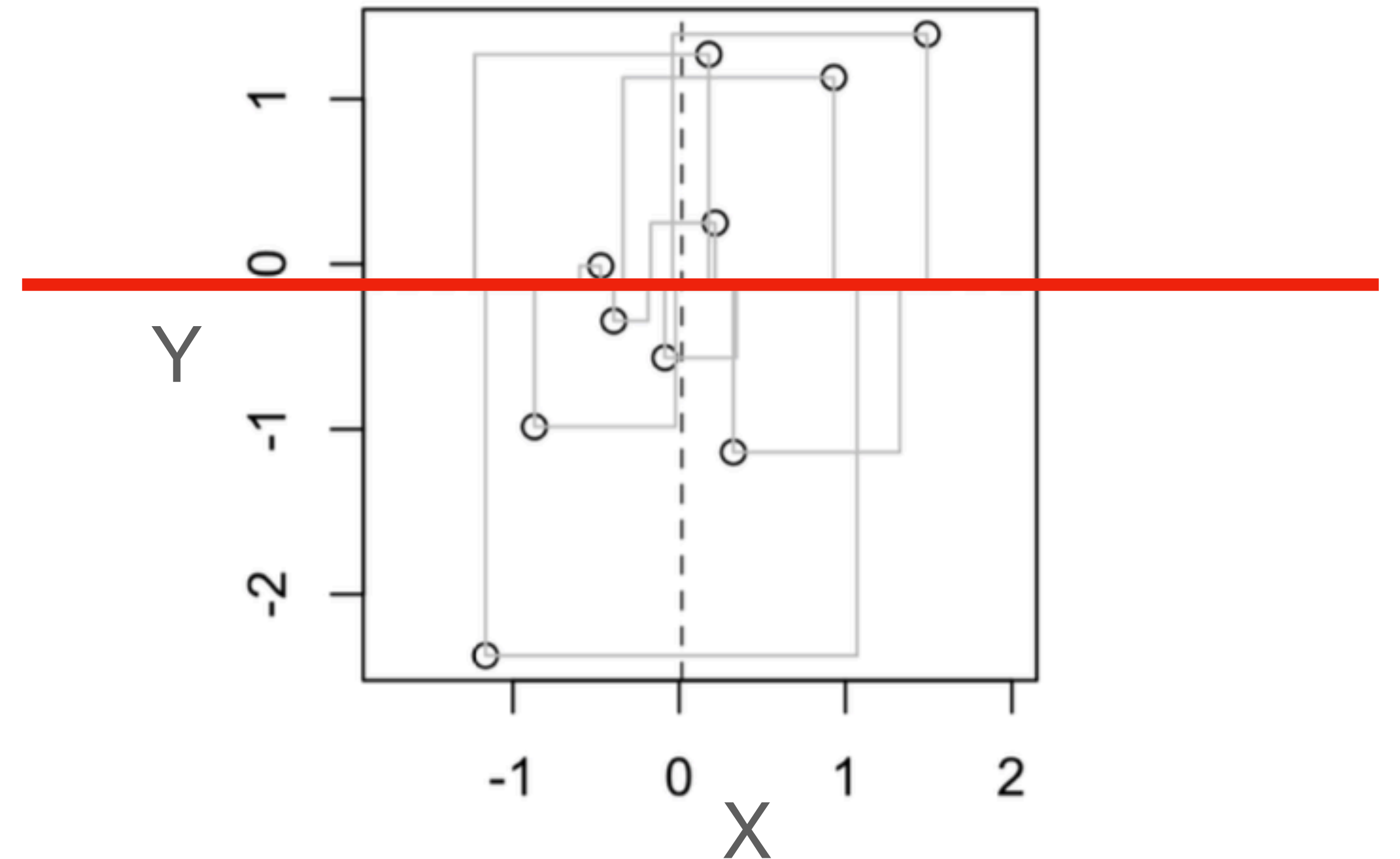
$$SS_{res} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

# How good is our model?

- We also have a measure of the **total variation in the dependent variable**: the sum of the squares of Y ($SS_Y$)

$$SS_Y = \sum_{i=1}^{n} (y_i - \bar{y})^2$$

# How good is our model?

- A natural measure of fit: the **coefficient of determination (R²)**

$$R^2 = 1 - \frac{SS_{res}}{SS_Y}$$

# How good is our model?

- A natural measure of fit: the **coefficient of determination (R²)**

$$R^2 = 1 - \frac{SS_{res}}{SS_Y}$$

How **bad** we are at explaining variation in Y

# How good is our model?

- A natural measure of fit: the **coefficient of determination ($R^2$)**

$$R^2 = 1 - \frac{SS_{res}}{SS_Y}$$

How **good** we are at explaining variation in Y
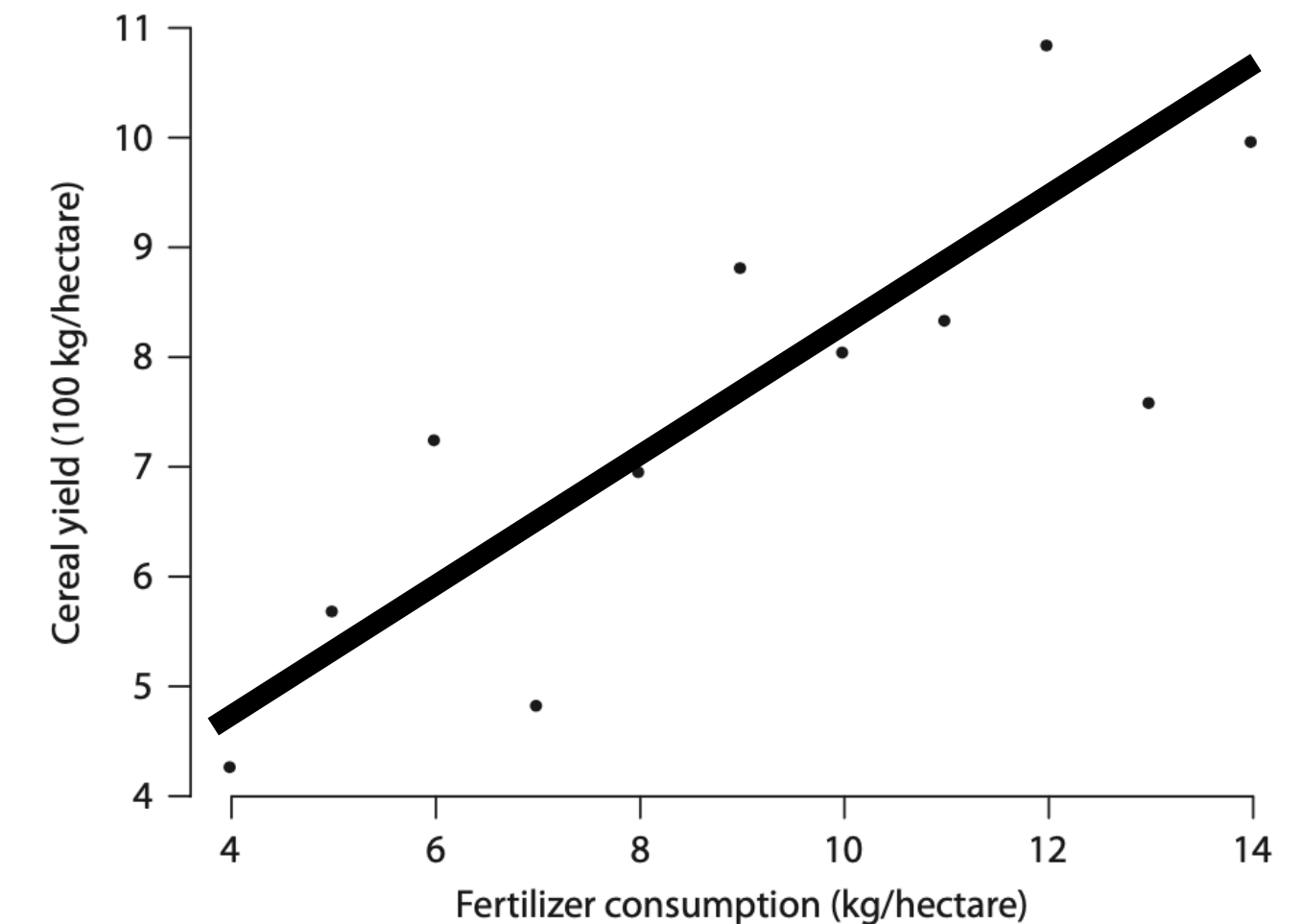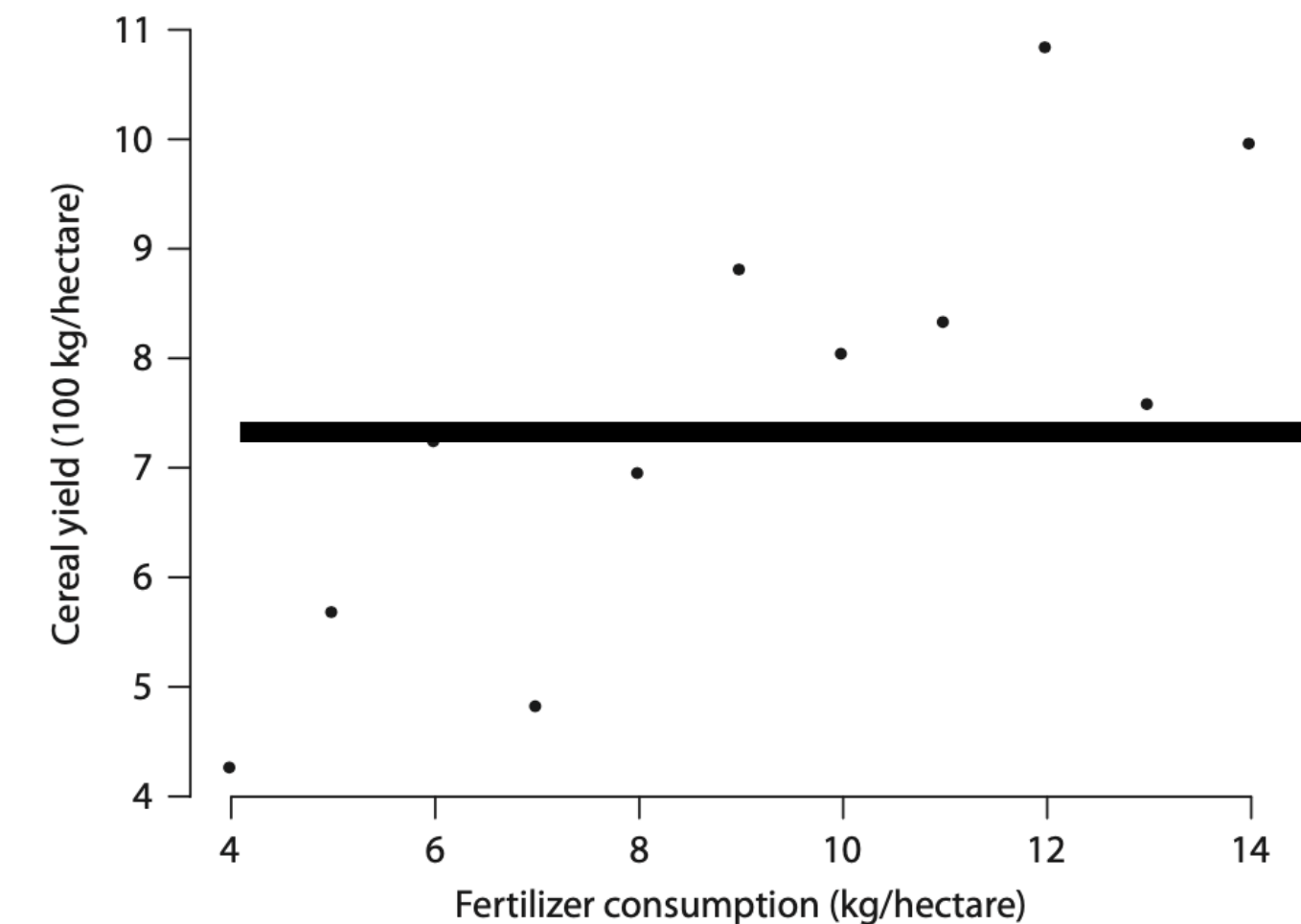
# Hypothesis testing via a linear model

# Hypothesis testing

$H_0$: there is **no** relationship between fertilizer consumption and cereal yield

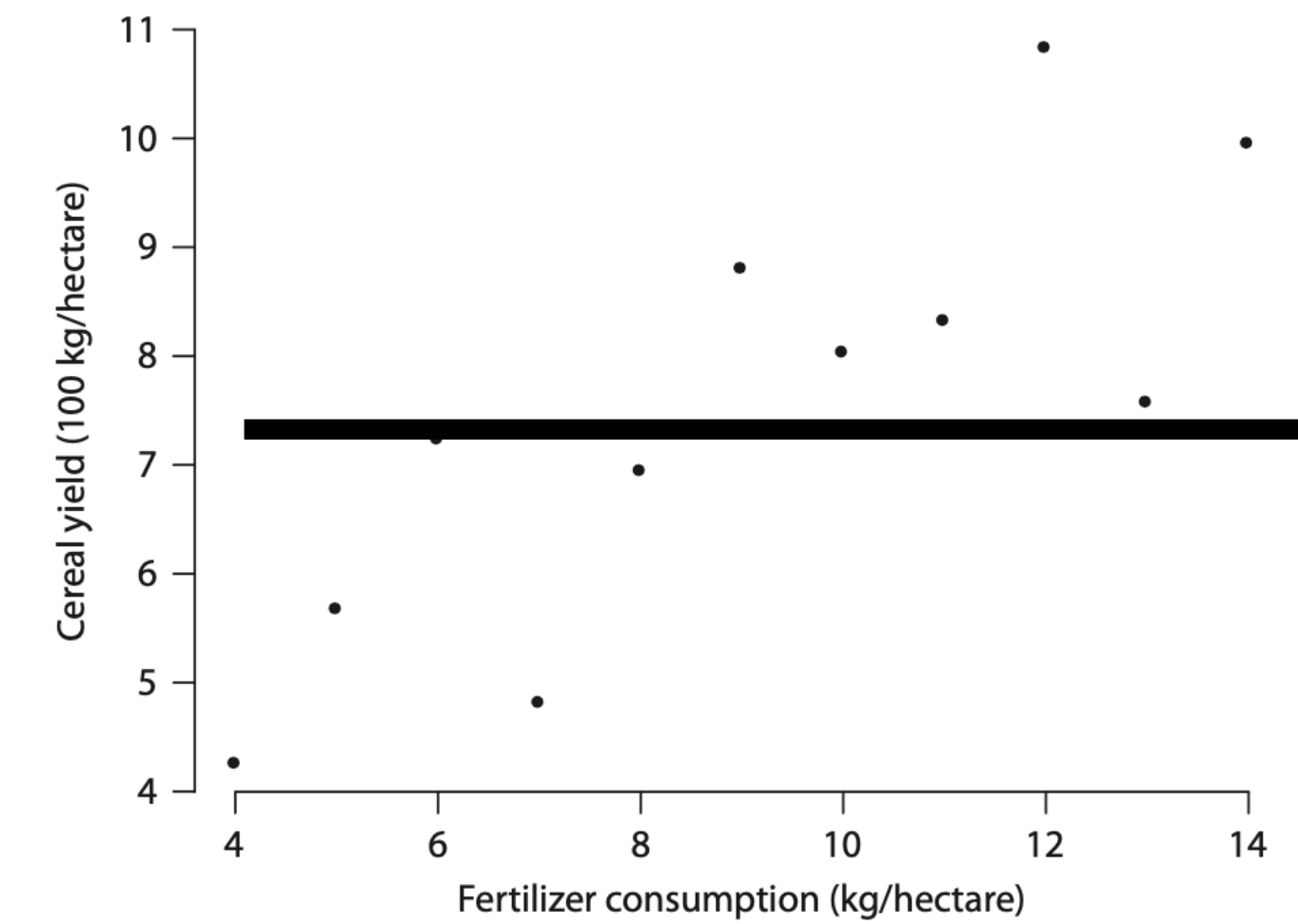Null (simpler) hypothesis ⟶ Can we reject it?

$H_1$: there is **some** relationship between fertilizer consumption and cereal yield
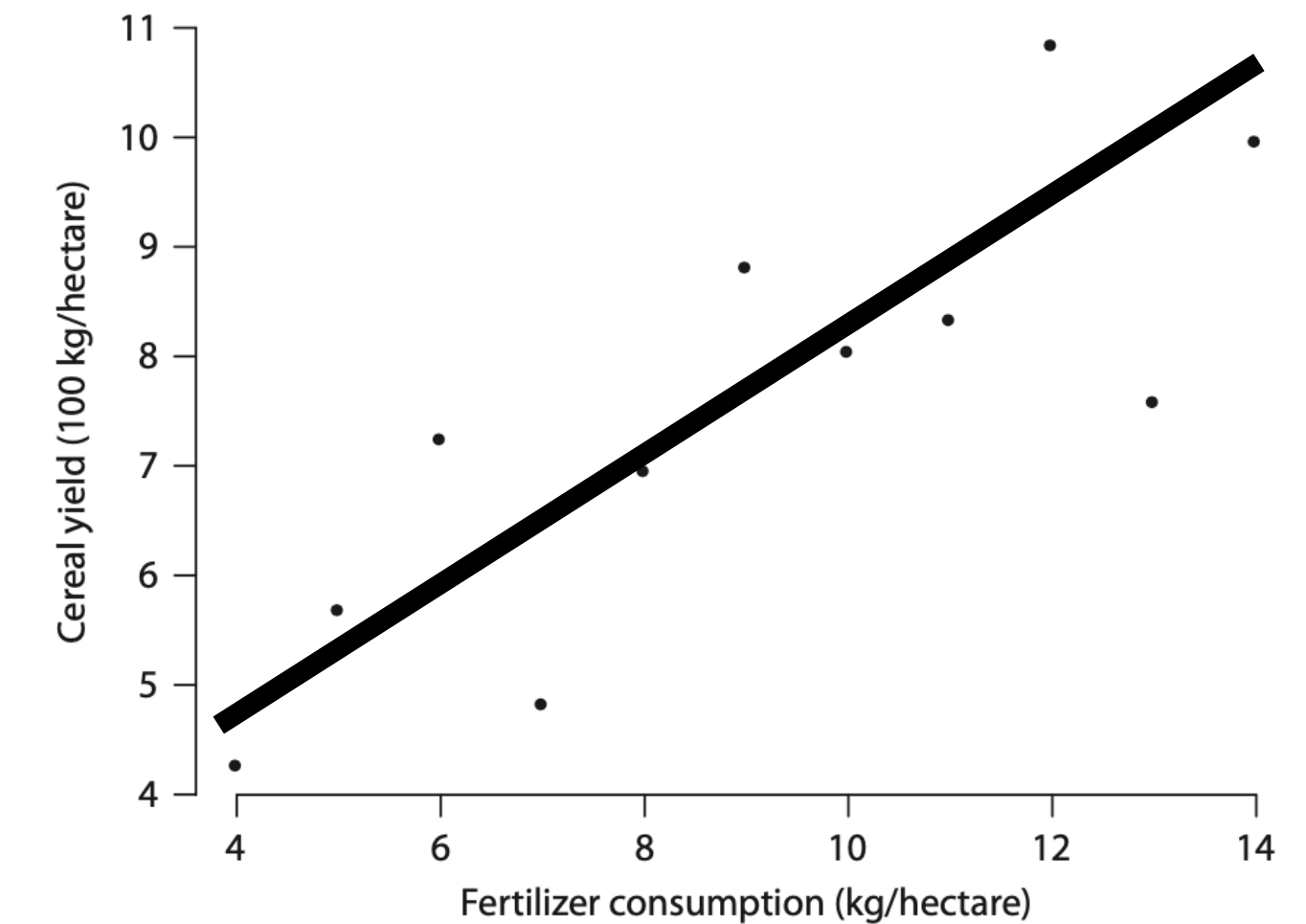
Alternative (more complicated) hypothesis

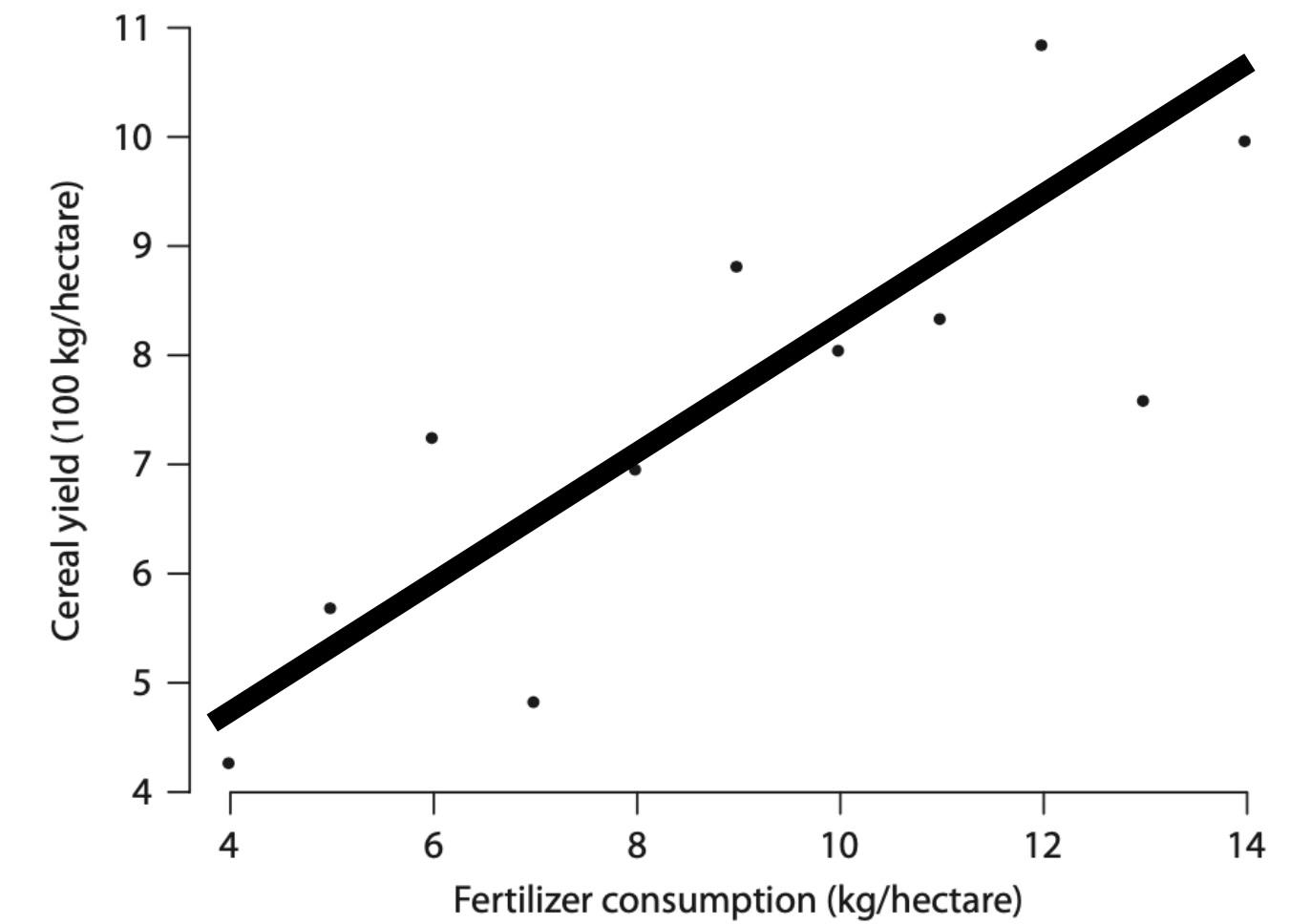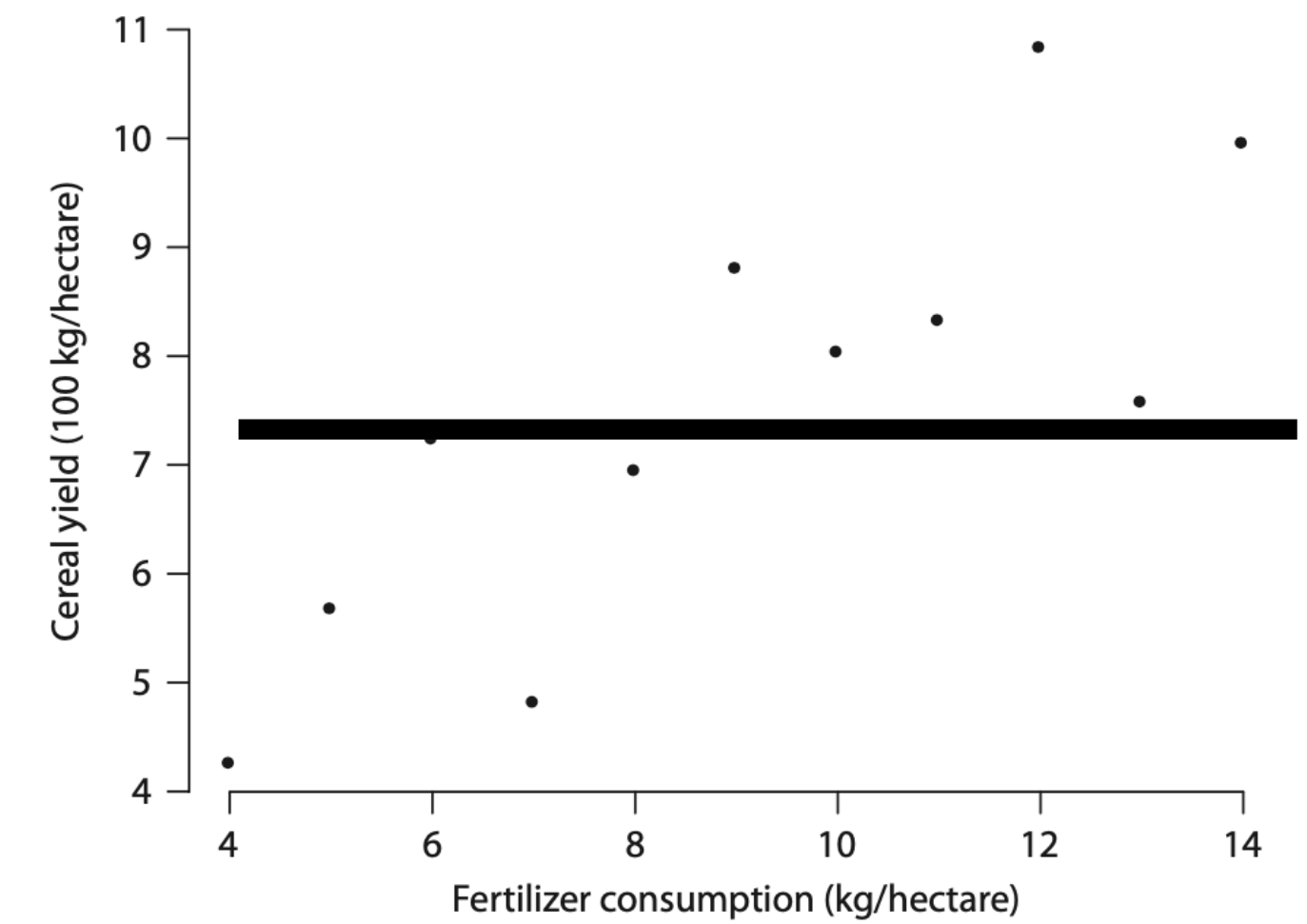# Hypothesis testing

$$\mathbf{y} = \beta_0 + \epsilon$$

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$$

# Hypothesis testing

$$\beta_1 = 0$$



$$\beta_1 \neq 0$$

# Test statistic: *t*

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

# Test statistic: *t*

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} \longrightarrow \text{Difference between our estimate } \hat{\beta}_1 \text{ and 0}$$

# Test statistic: *t*

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

Difference between our estimate $\hat{\beta}_1$ and 0

# Test statistic: *t*

The bigger the difference, the bigger is *t*
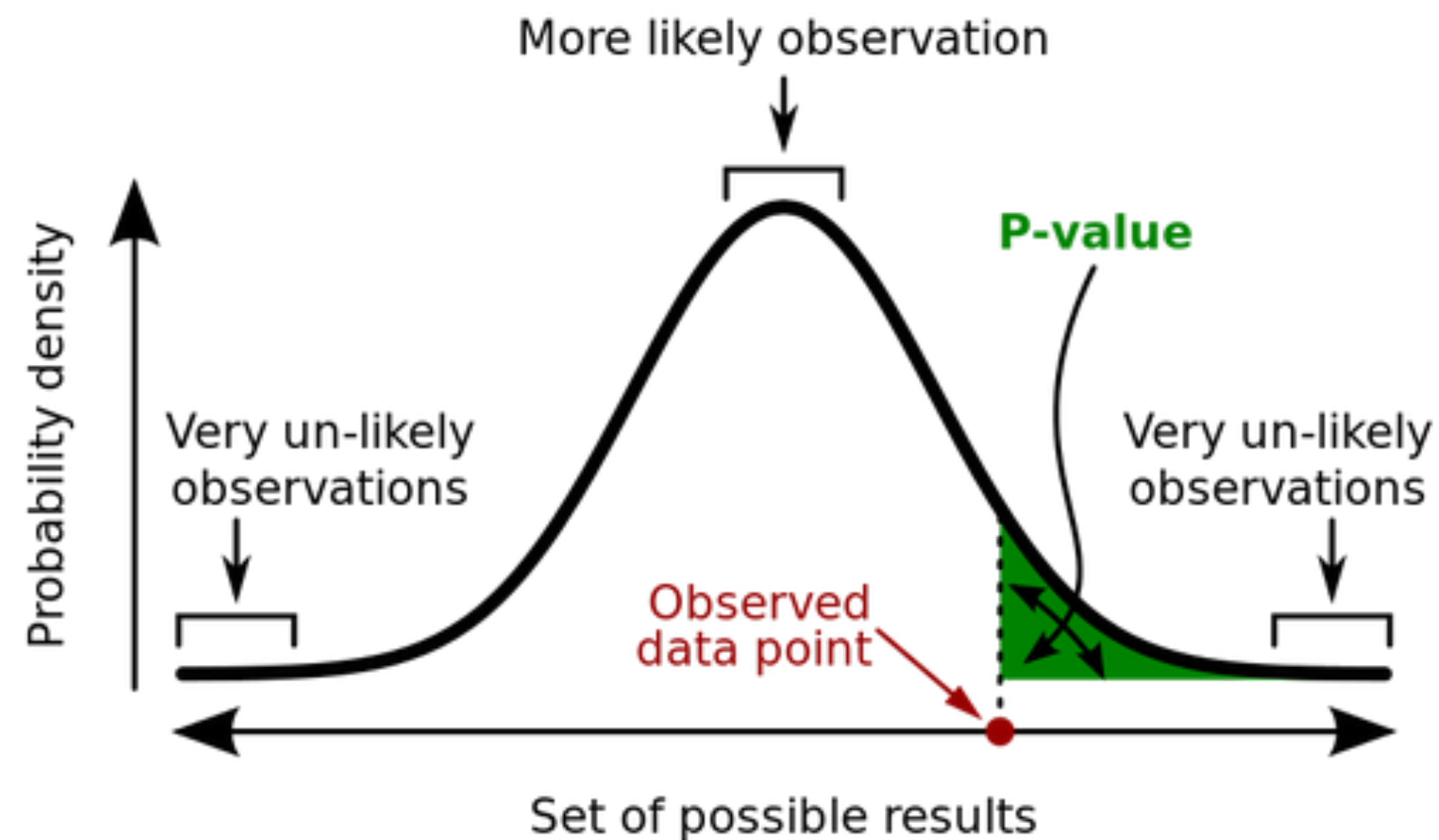
$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

Difference between our estimate $\hat{\beta}_1$ and 0

Standard Error: a measure of how inaccurate our estimate $\hat{\beta}_1$ is at estimating $\beta_1$

# Test statistic: *t*

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

<span style="color:red">The bigger the difference, the bigger is *t*</span>

→ Difference between our estimate $\hat{\beta}_1$ and 0

→ Standard Error: a measure of how inaccurate our estimate $\hat{\beta}_1$ is at estimating $\beta_1$

<span style="color:red">The bigger the error, the smaller is *t*</span>
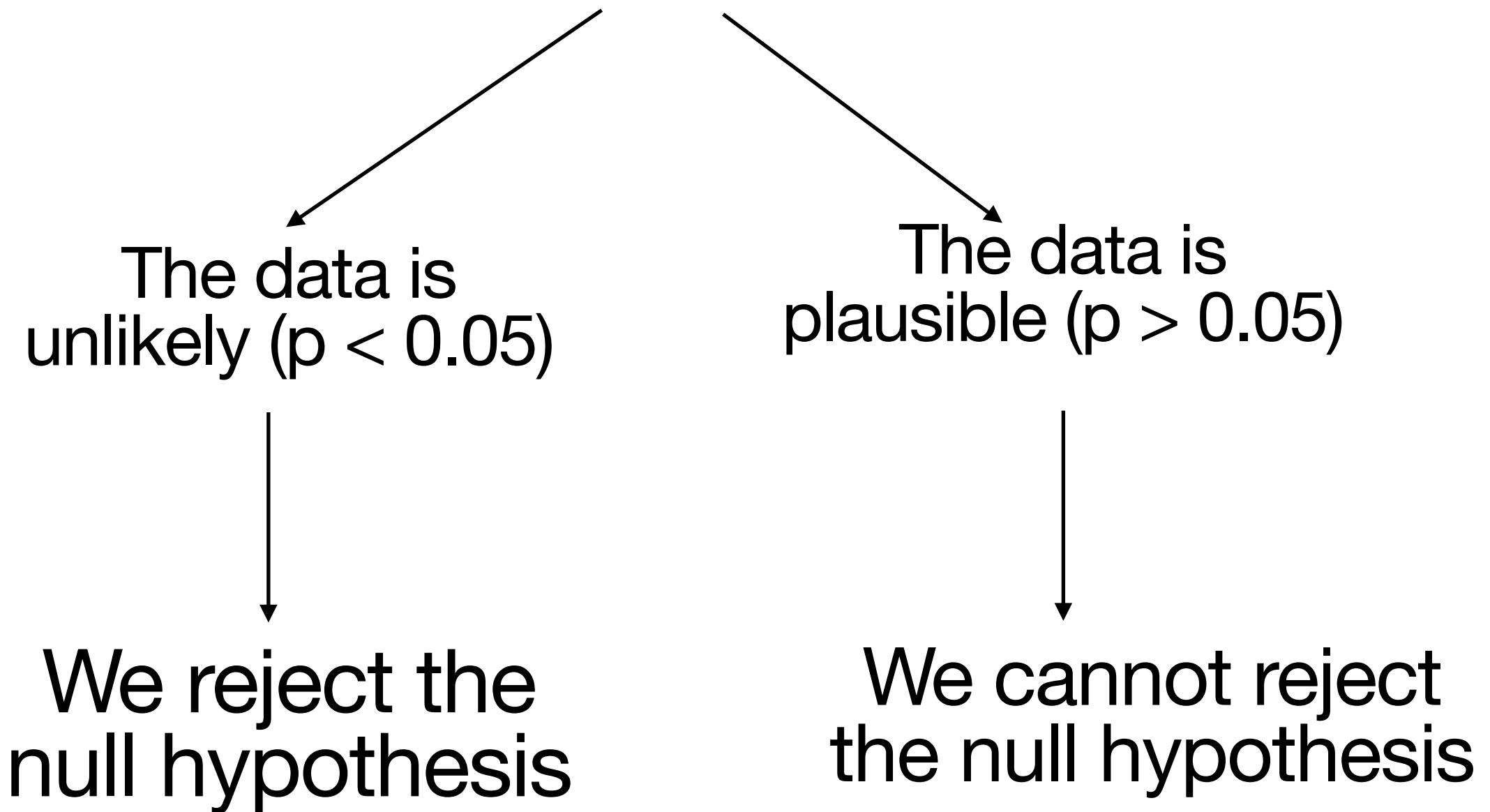
# Test statistic: *t*

Assuming $\beta_1 = 0$, then the t-statistic should follow a well-known distribution: the t-distribution.



If our estimate is **too unlikely** under this distribution, then we can **reject the hypothesis** $\beta_1 = 0$

# Frequentist hypothesis testing

Assuming the null hypothesis is true…

The data is unlikely ($p < 0.05$)

The data is plausible ($p > 0.05$)

We reject the null hypothesis
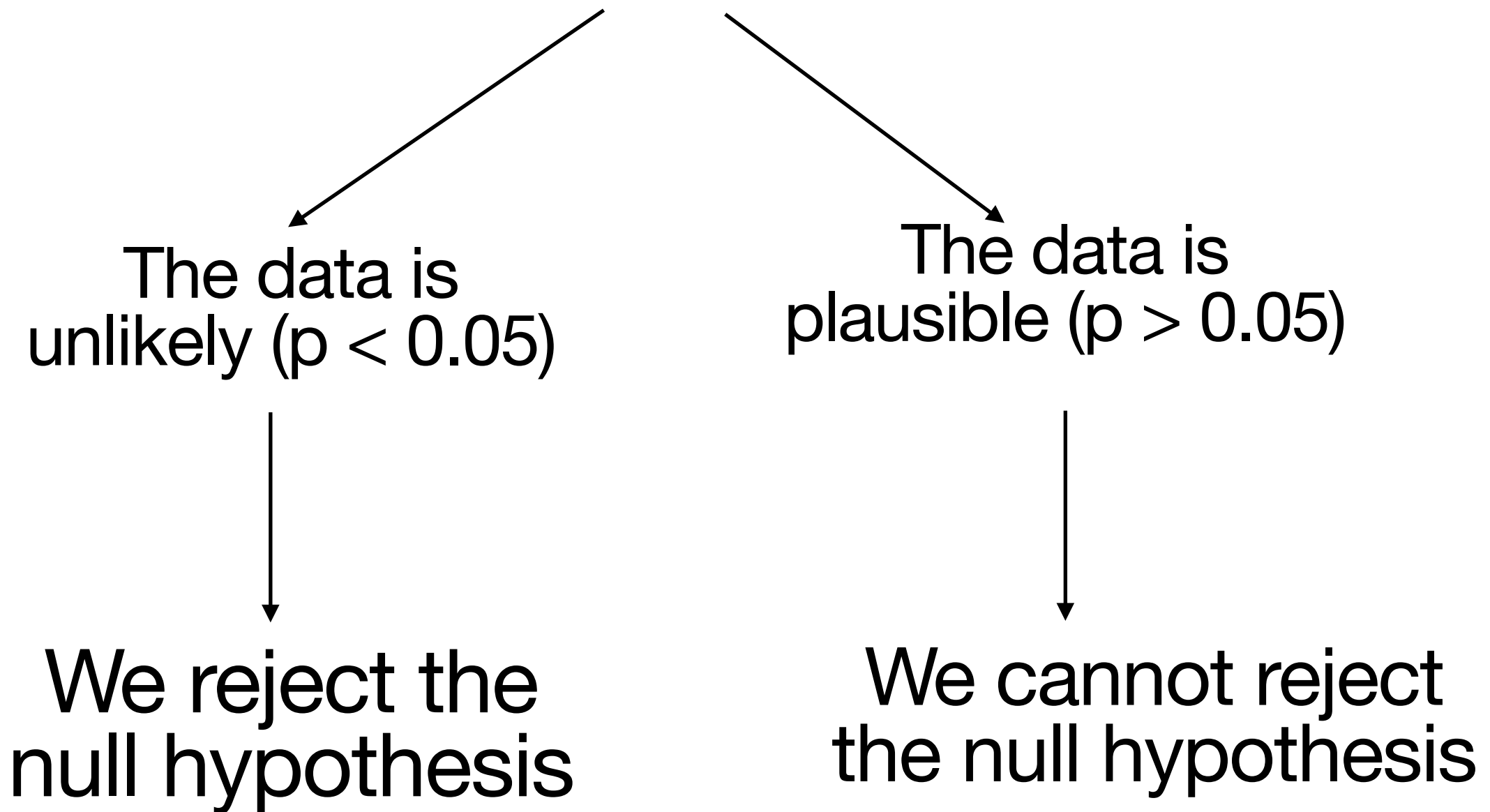
We cannot reject the null hypothesis

# Problems with hypothesis testing

- Who decides what is "too unlikely"? ($p < 0.05$ is an arbitrary cutoff)

- Could a "rejected" model still be a plausible explanation of the data?

- If we "reject" a null model, are we necessarily "accepting" the alternative?
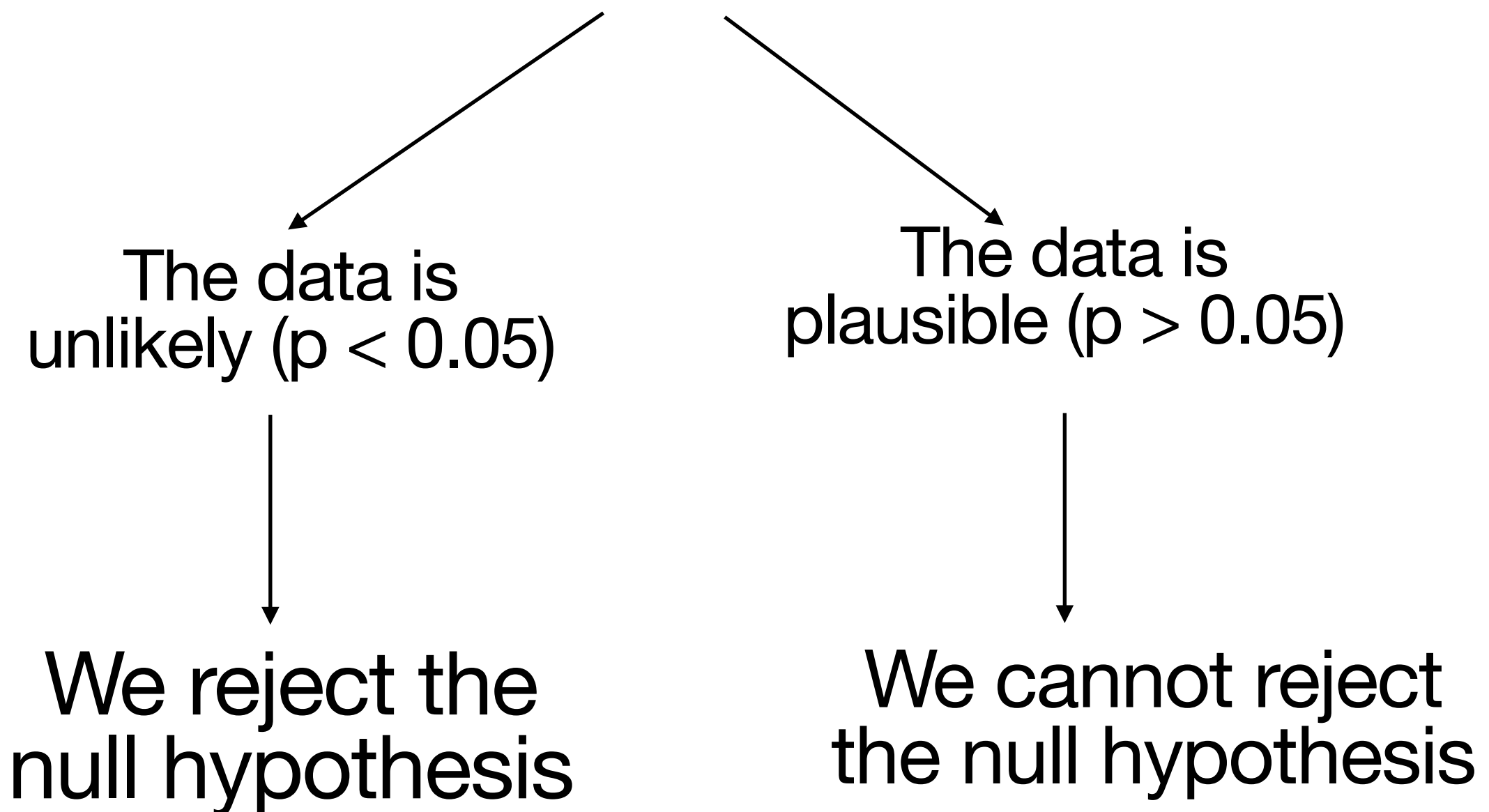
- What if there are multiple alternative models?

# Problems with hypothesis testing

- Who decides what is "too unlikely"? ($p < 0.05$ is an arbitrary cutoff)

- Could a "rejected" model still be a plausible explanation of the data?

- If we "reject" a null model, are we necessarily "accepting" the alternative?

- What if there are multiple alternative models?

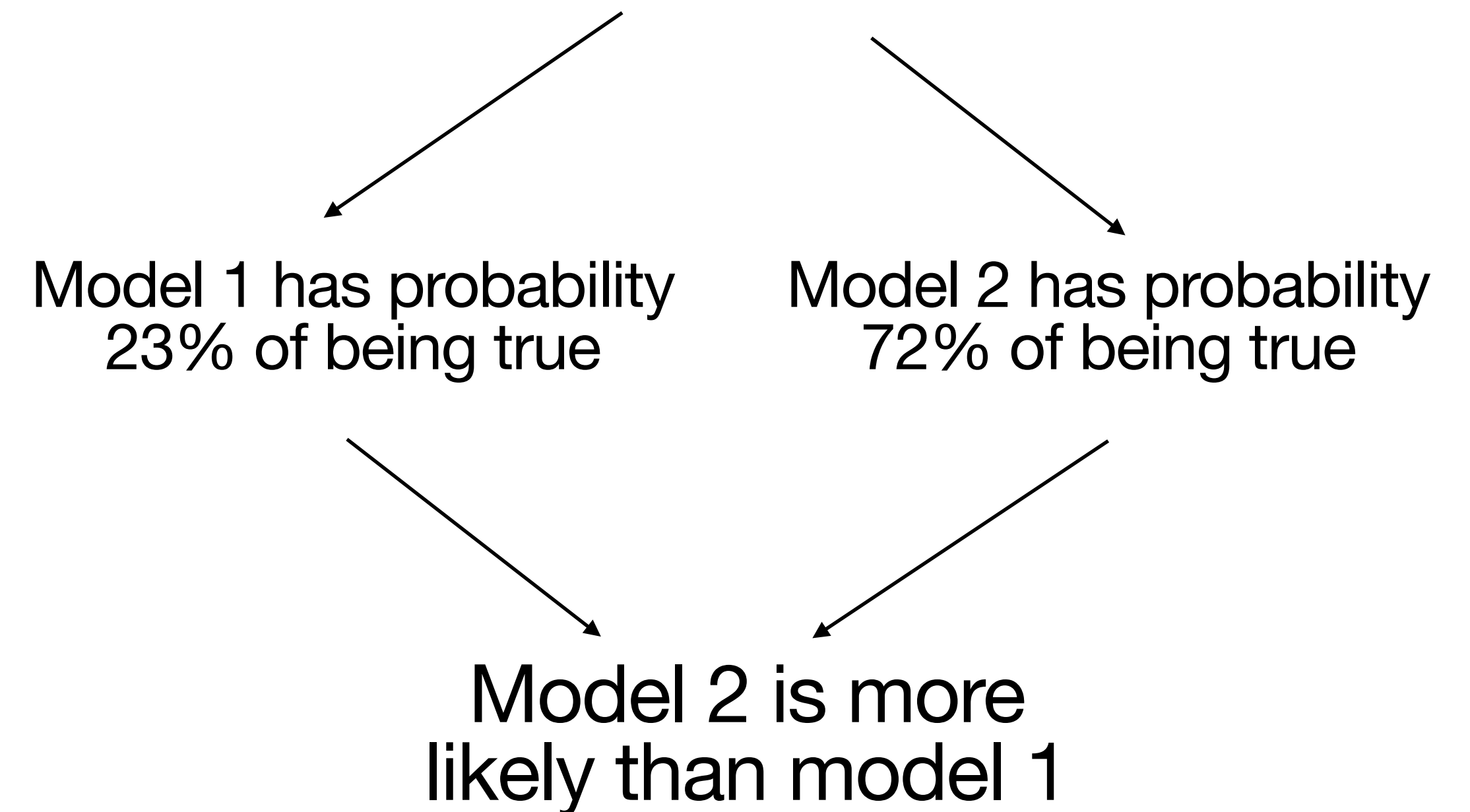# Frequentist hypothesis testing

Assuming the null hypothesis is true…

The data is unlikely ($p < 0.05$)

The data is plausible ($p > 0.05$)

We reject the null hypothesis

We cannot reject the null hypothesis

**Frequentist hypothesis testing**

Assuming the null hypothesis is true…

The data is unlikely ($p < 0.05$)

The data is plausible ($p > 0.05$)

We reject the null hypothesis

We cannot reject the null hypothesis

**Bayesian model choice**

Given the data…

Model 1 has probability 23% of being true

Model 2 has probability 72% of being true

Model 2 is more likely than model 1

# Frequentist hypothesis testing

Assuming the null hypothesis is true…

The data is unlikely ($p < 0.05$)

The data is plausible ($p > 0.05$)

We reject the null hypothesis

We cannot reject the null hypothesis

# Bayesian model choice

Given the data…

Model 1 has probability 23% of being true

Model 2 has probability 72% of being true

Model 2 is more likely than model 1

To find out more: take the course "Advanced Topics in Data Analysis"!

- **Fitting a simple linear model in R**
- **Interpreting a linear model in R**

# The two sides of statistics

# The two sides of statistics

Probability Theory → ← Statistical Inference

Probability Theory

"What can we say about the data generated by a given process?"

Statistical Inference

"What can we say about the process that generated a given data?"

# Linear Regression: a probabilistic interpretation

**Question:** What can we say about Y after we've fitted our model?

?

# Linear Regression: a probabilistic interpretation

**Question:** What can we say about Y after we've fitted our model?

!

**Answer:** It depends on what we're willing to assume…

# Linear Regression: a probabilistic interpretation

Let's treat the variable Y as a random variable with some (possibly unknown) distribution

# "Linearity" assumption

$$E[\epsilon \mid X = x] = 0$$

**"The expected value of the error term is zero, regardless of the value of X"**

# "Linearity" assumption

$$E[\epsilon \mid X = x] = 0$$

"The expected value of the error term is zero, regardless of the value of X"

**Then:**

$$E[Y \mid X = x] = \beta_0 + \beta_1 x$$

"Given a value of X, the expected value of Y is a linear function of that value of X"

# E[Y|X=x] is the best linear predictor for X=x

# E[Y|X=x] is the best linear predictor for X=x
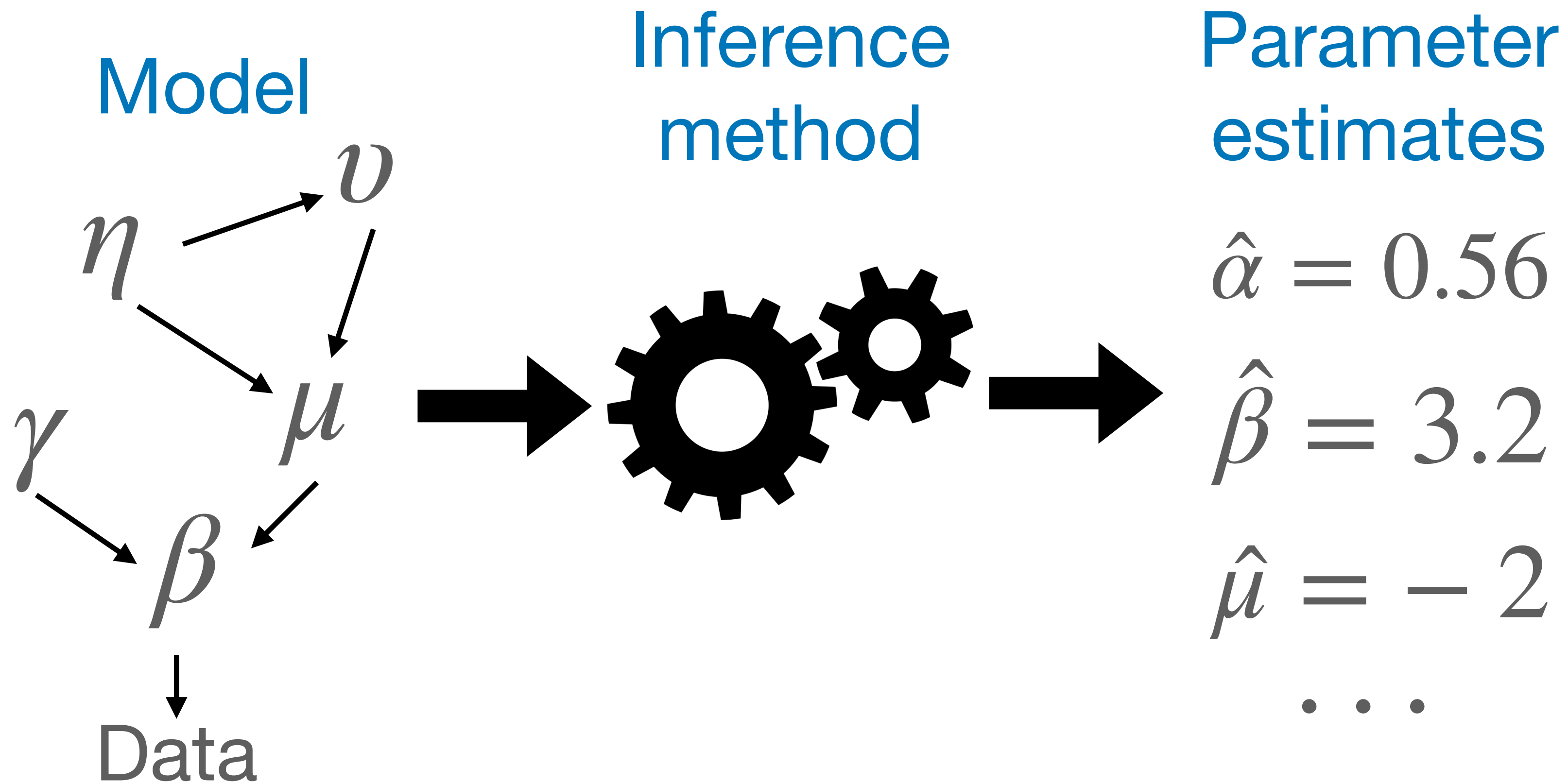
# E[Y|X=x] is the best linear predictor for X=x
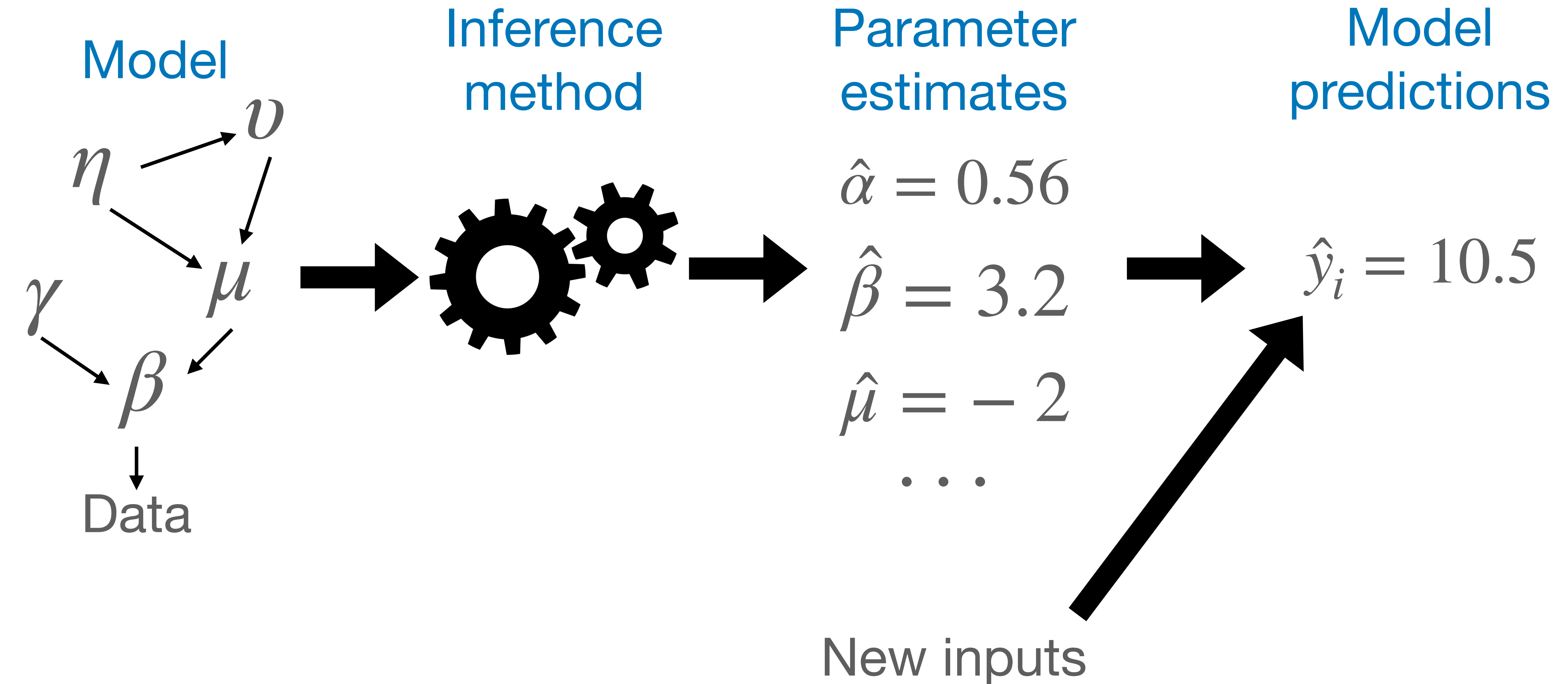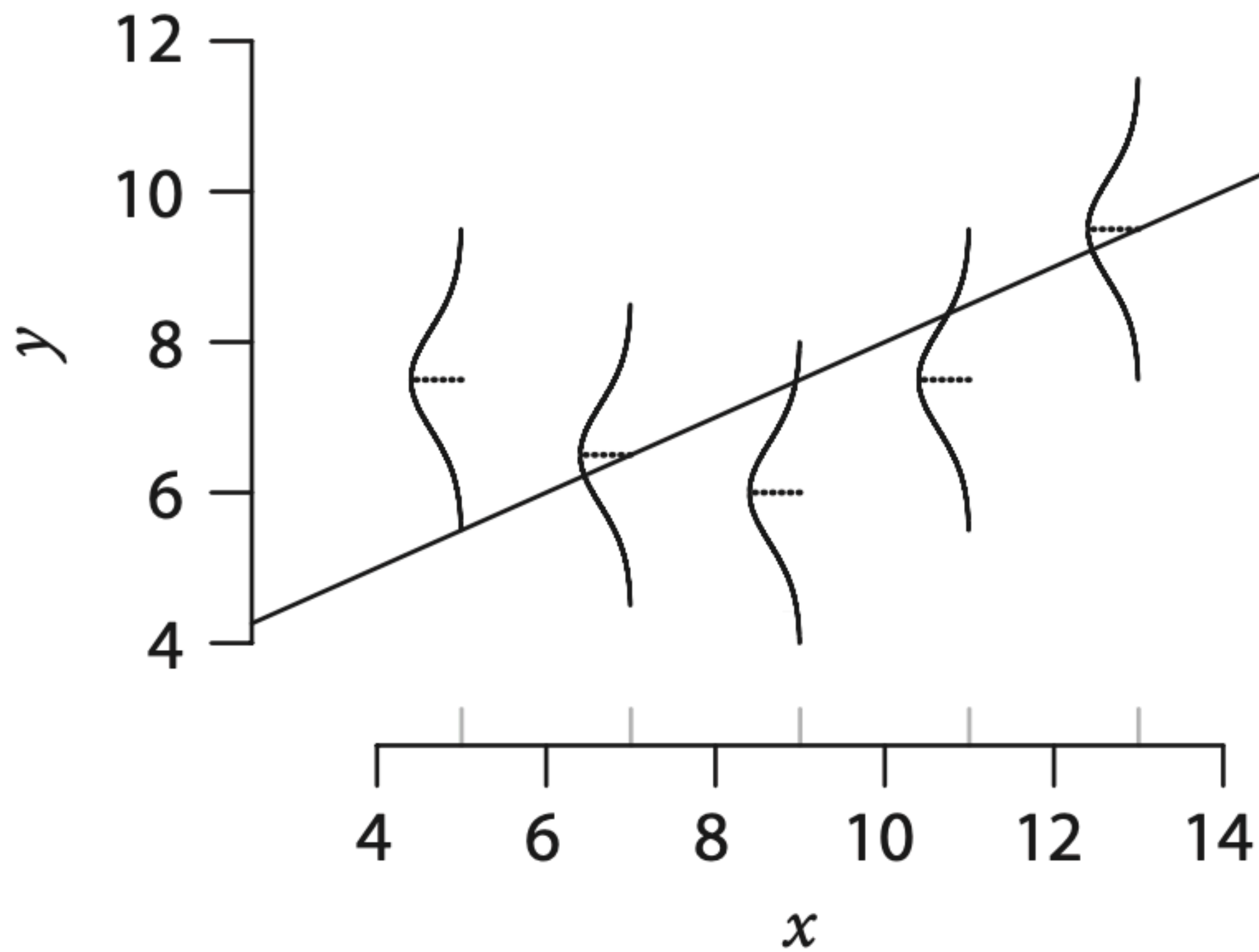
# E[Y|X=x] is the best linear predictor for X=x

# E[Y|X=x] is the best linear predictor for X=x

E[Y|X=x] is the best linear predictor for X=x

Edge 2019

# E[Y|X=x] is the best linear predictor for X=x

# Inference vs. Prediction

Model

$\upsilon$
$\eta$
$\gamma$
$\mu$
$\beta$

Data

Inference
method

Parameter
estimates

$\hat{\alpha} = 0.56$

$\hat{\beta} = 3.2$

$\hat{\mu} = -2$

$\ldots$

# Inference vs. Prediction

Model

Inference
method

Parameter
estimates

Model
predictions

$\upsilon$

$\eta$

$\gamma$

$\mu$

$\beta$

Data

$\hat{\alpha} = 0.56$

$\hat{\beta} = 3.2$

$\hat{\mu} = -2$

$\dots$

$\hat{y}_i = 10.5$

New inputs

# Linearity violated

# "Homoscedasticity" assumption

$$Var[\epsilon \mid X = x] = \sigma_\epsilon^2$$

**"The variance of the error is a constant, regardless of the value of X"**

# "Homoscedasticity" assumption

$$Var[\epsilon \mid X = x] = \sigma_\epsilon^2$$

**"The variance of the error is a constant, regardless of the value of X"**

## Then:

$$Var[Y \mid X = x] = \sigma_\epsilon^2$$

**"The variance of Y is a constant (equal to the variance of the error), regardless of the value of X"**

# Homoscedasticity violated

# "Normality" assumption

$$\epsilon_i \,|\, (X = x) \sim Normal(0, \sigma_\epsilon^2)$$

and all $\epsilon_i$ are independent of each other

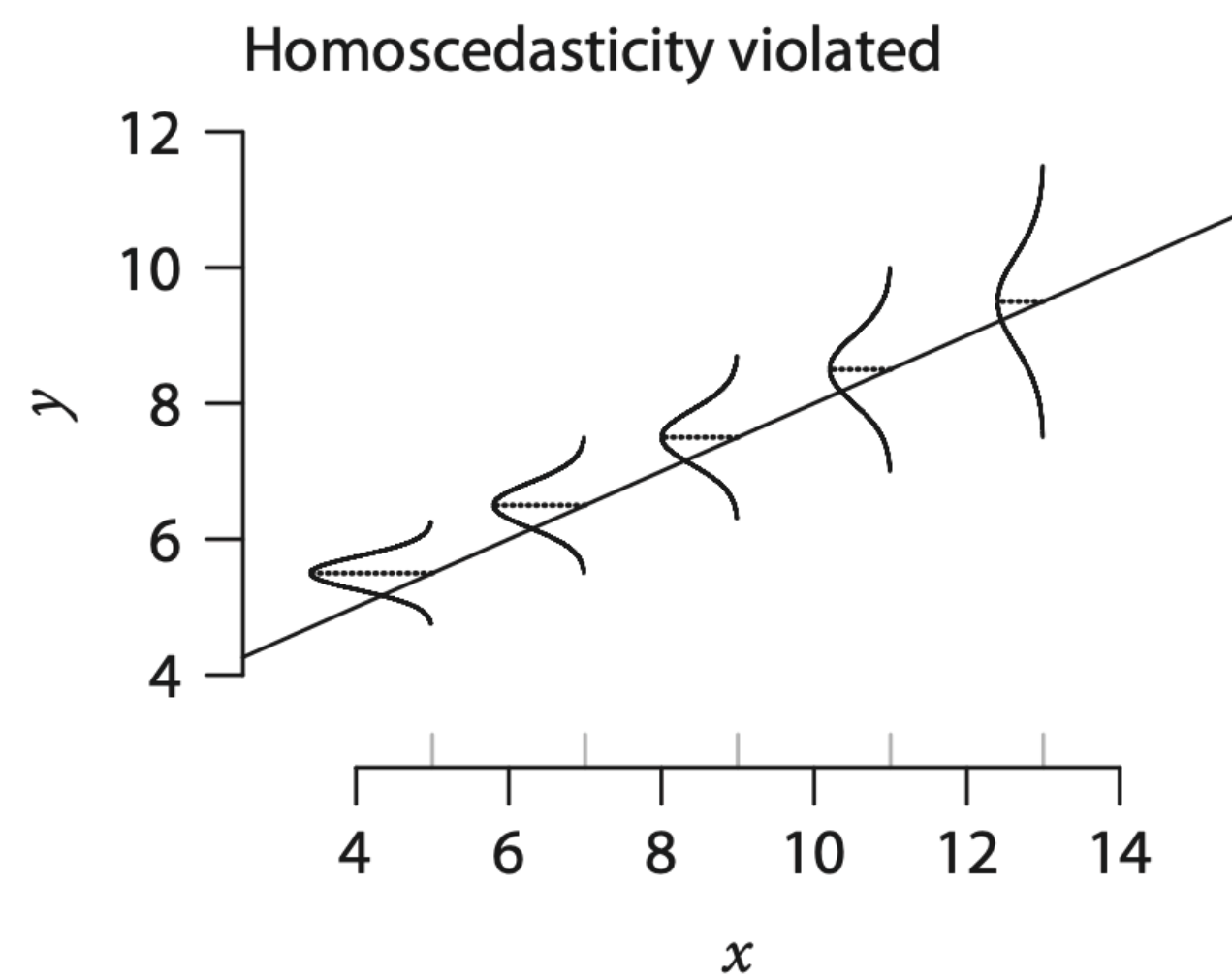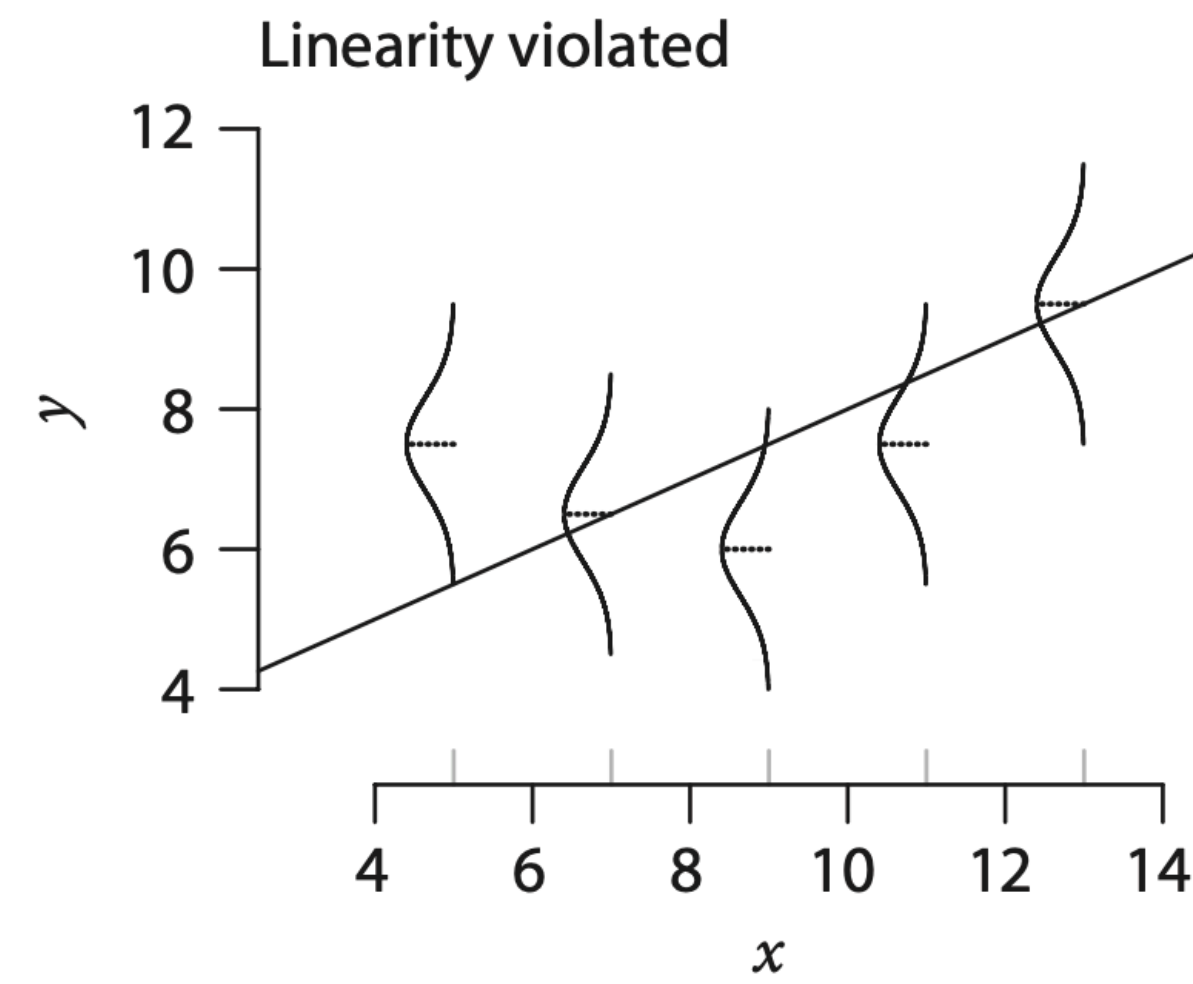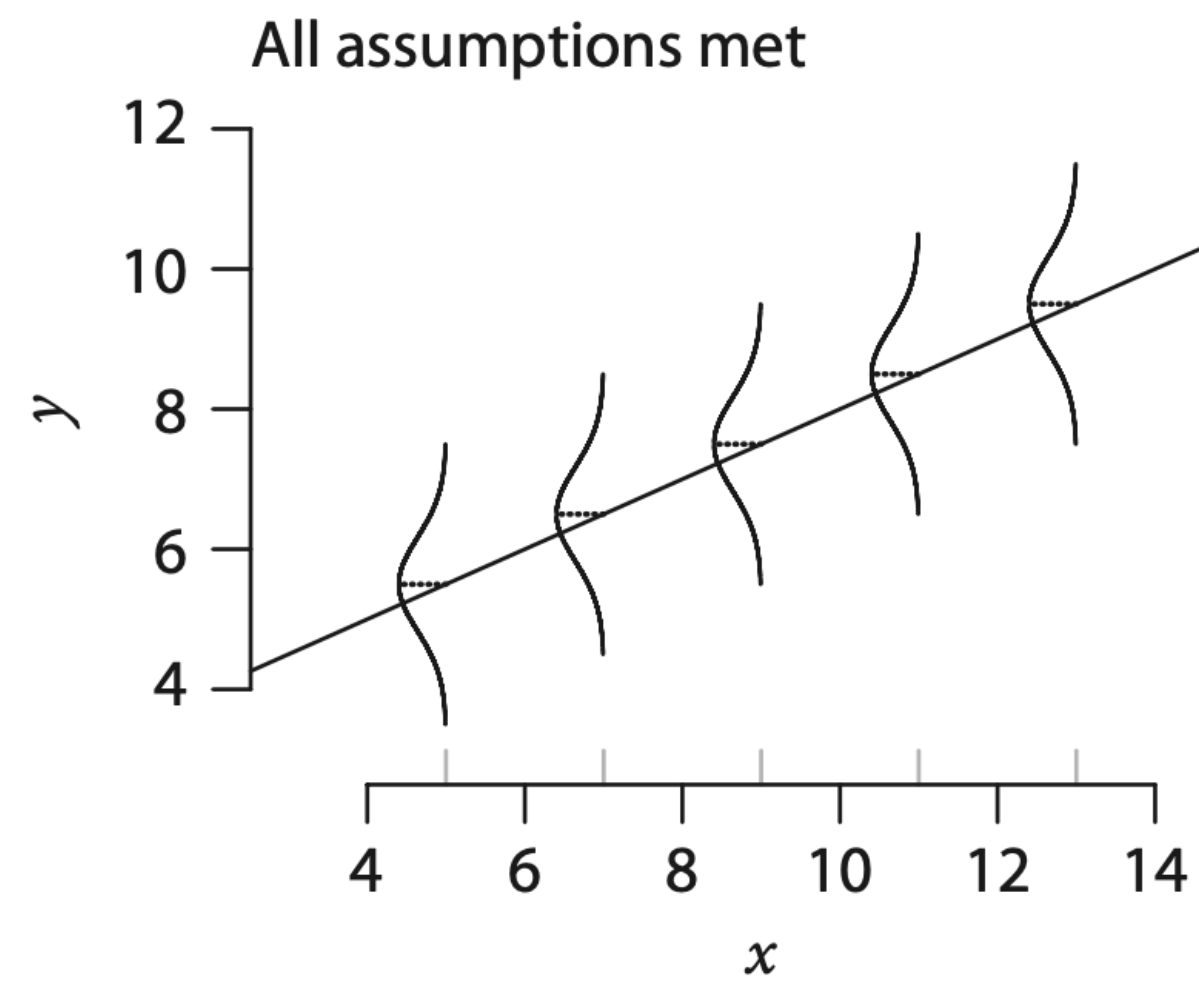**"The error is normally distributed, regardless of the value of X"**

# "Normality" assumption

$$\epsilon_i \,|\, (X = x) \sim Normal(0, \sigma_\epsilon^2)$$

and all $\epsilon_i$ are independent of each other

**"The error is normally distributed, regardless of the value of X"**

## Then:

$$Y \,|\, (X = x) \sim Normal(\beta_0 + \beta_1 x, \sigma_\epsilon^2)$$

**"Y is normally distributed, with mean equal to a linear function of the value of X"**

# Normality violated

# Violations of assumptions

That's it for simple linear regression today, but:

# A zoo of models…

# Simple linear regression

- **1** predictor variable (x)

- **1** response variable (y)

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$$

Variable **y** is a linear function of **x**, plus some **noise**

James et al. 2014

# Simple linear regression

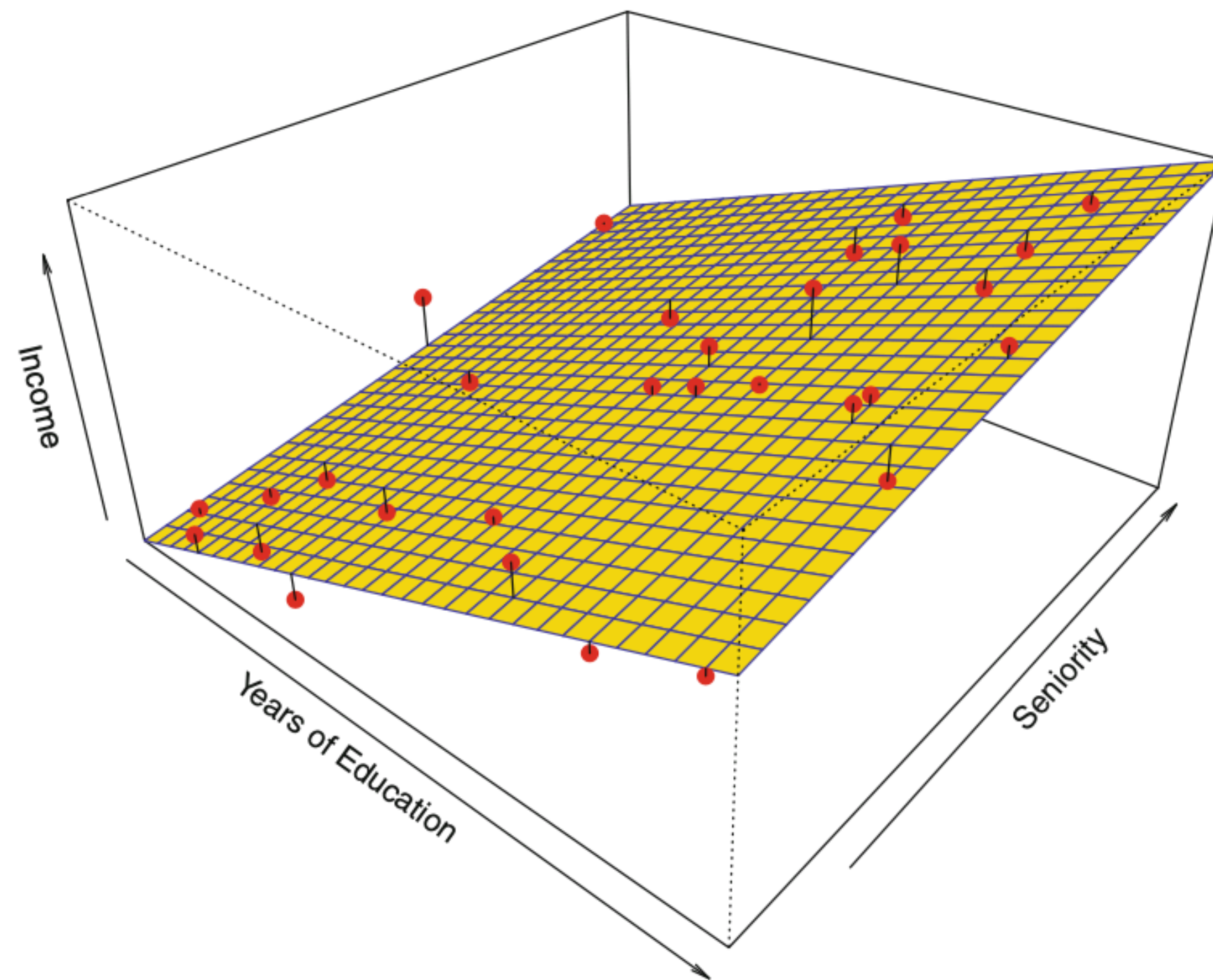- **1** predictor variable (x)

- **1** response variable (y)

# Multiple linear regression

- **Several** predictor variables

- **1** response variable

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \epsilon$$

Variable **y** is a linear function of **$x_1$, $x_2$, $x_3$, etc.** plus some **noise**
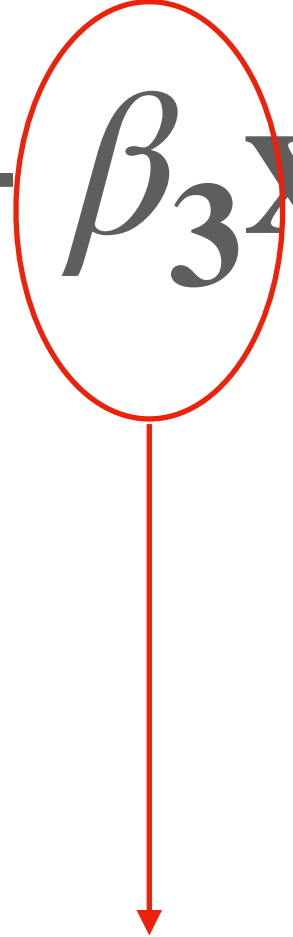
# Multiple linear regression

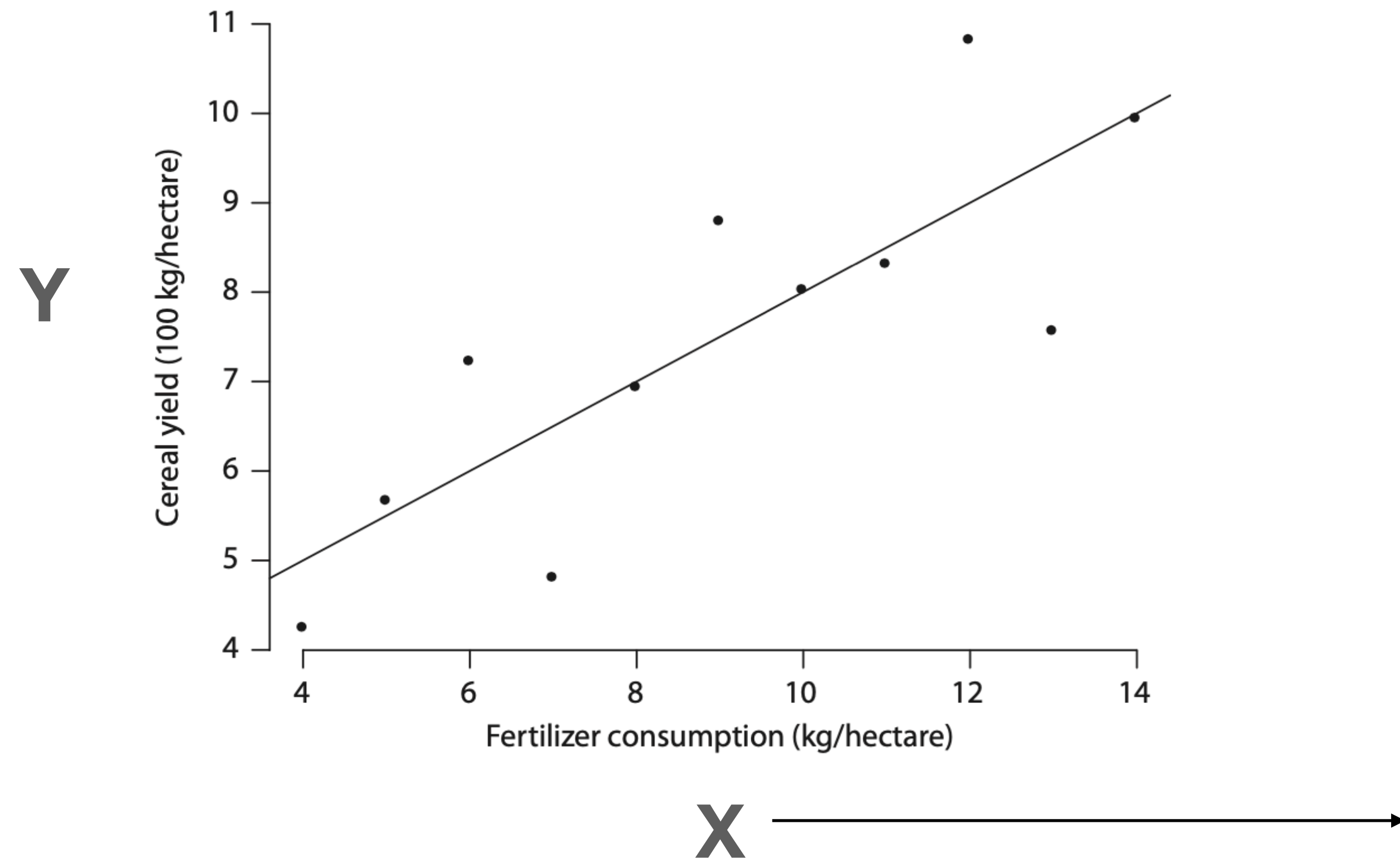- **Several** predictor variables

- **1** response variable



James et al. 2014

# Expanding the model: interaction terms

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \ldots + \epsilon$$

# Expanding the model: interaction terms

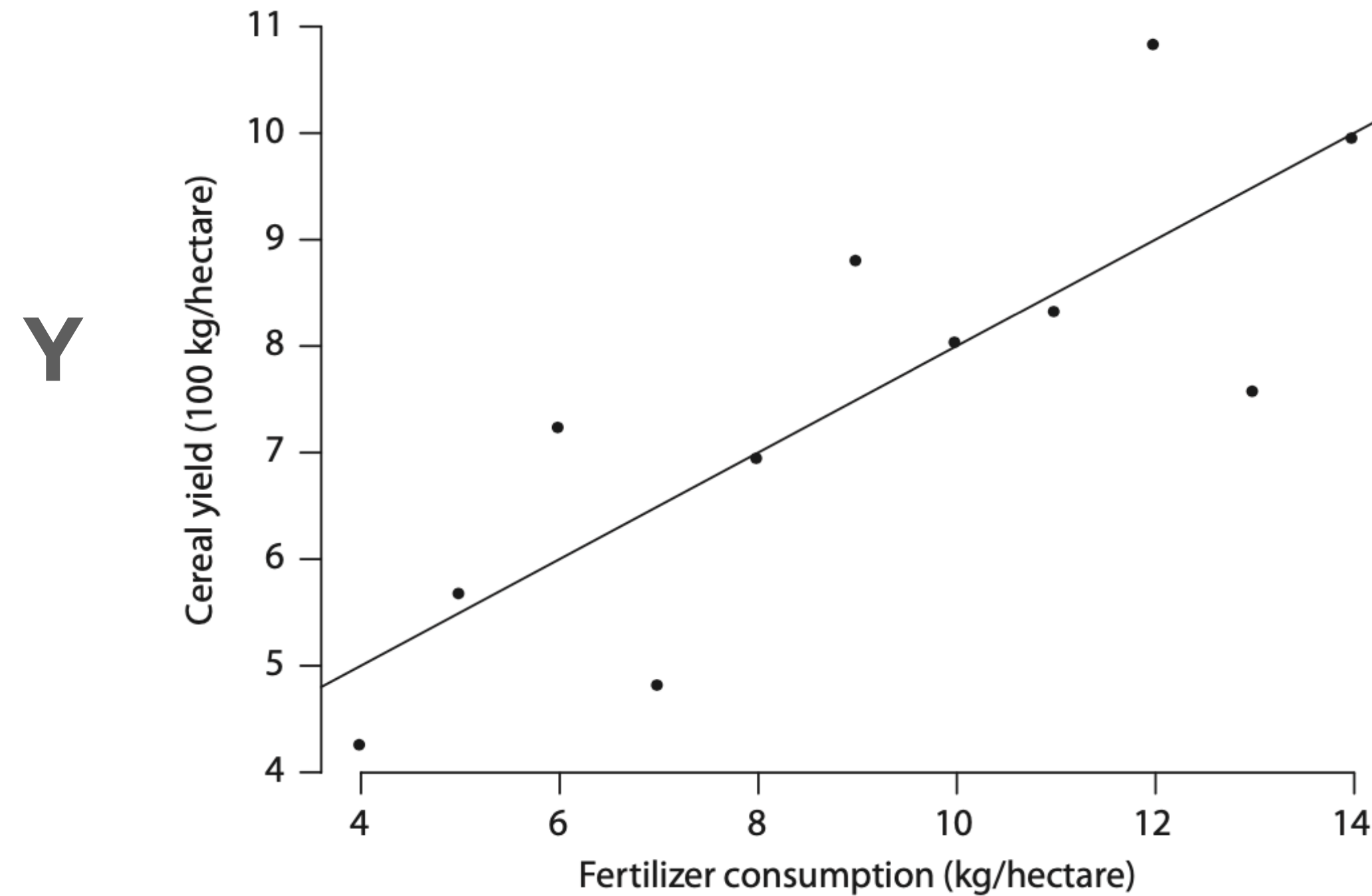$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \ldots + \epsilon$$

will be large when the combined behavior of $x_1$ and $x_2$ jointly influence the value of y

# Multiple linear regression with categorical predictors (ANOVA)



**Y**

Cereal yield (100 kg/hectare)

Fertilizer consumption (kg/hectare)

**X** →

So far, we've assumed that the value of **X** is a continuous number

# Multiple linear regression with categorical predictors (ANOVA)



**Y**

Cereal yield (100 kg/hectare)

Fertilizer consumption (kg/hectare)

**X** ⟶

What if it was a **category** instead?

So far, we've assumed that the value of **X** is a continuous number

# Multiple linear regression with categorical predictors (ANOVA)



**Height**

0      1
category
(e.g. male vs. female)

$$y = \beta_0 + \beta_1 x + \epsilon$$

We can modify our original linear model such that **x** is now a **dummy variable.**

# Multiple linear regression with categorical predictors (ANOVA)

$$y = \beta_0 + \beta_1 x + \epsilon$$

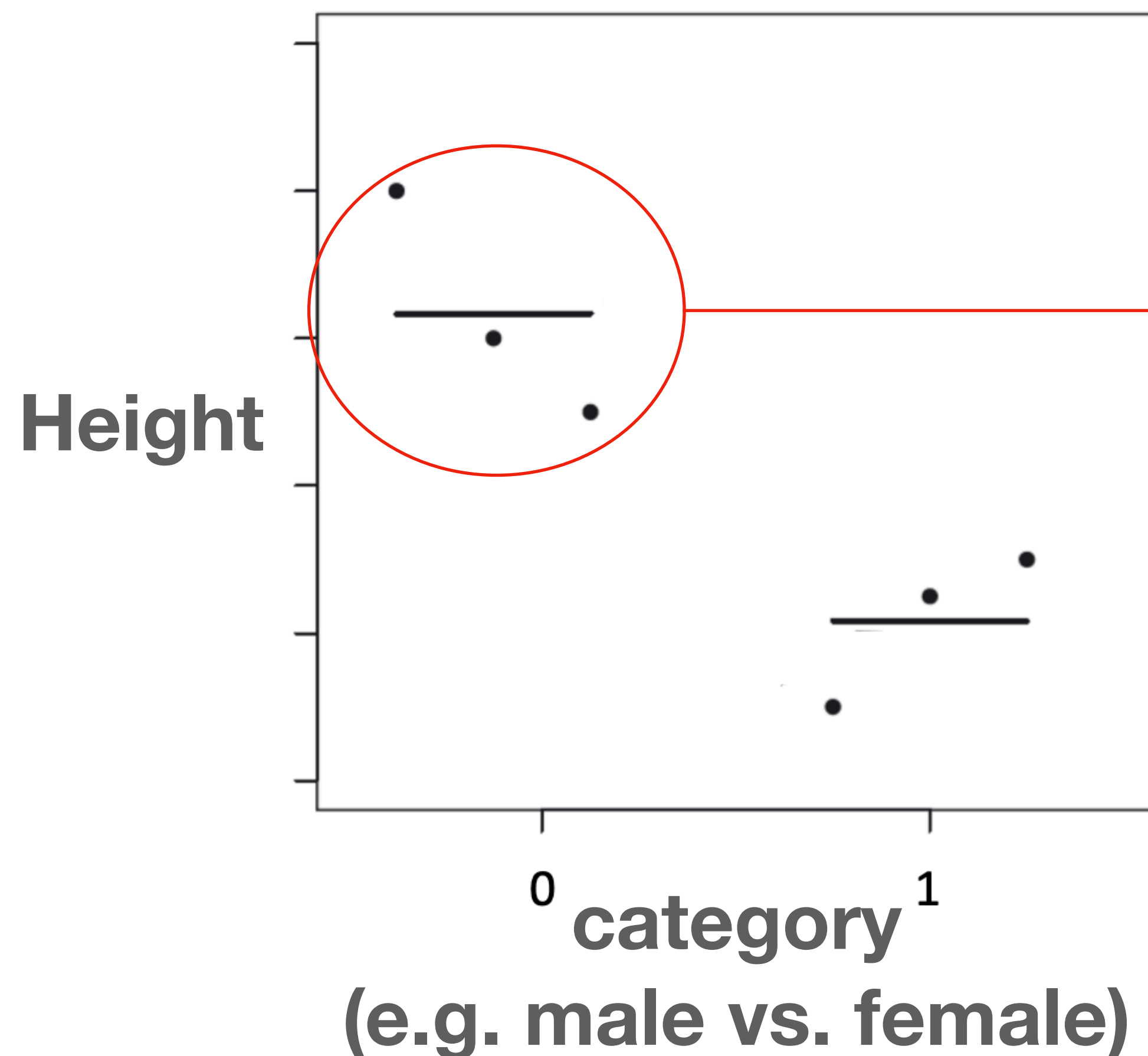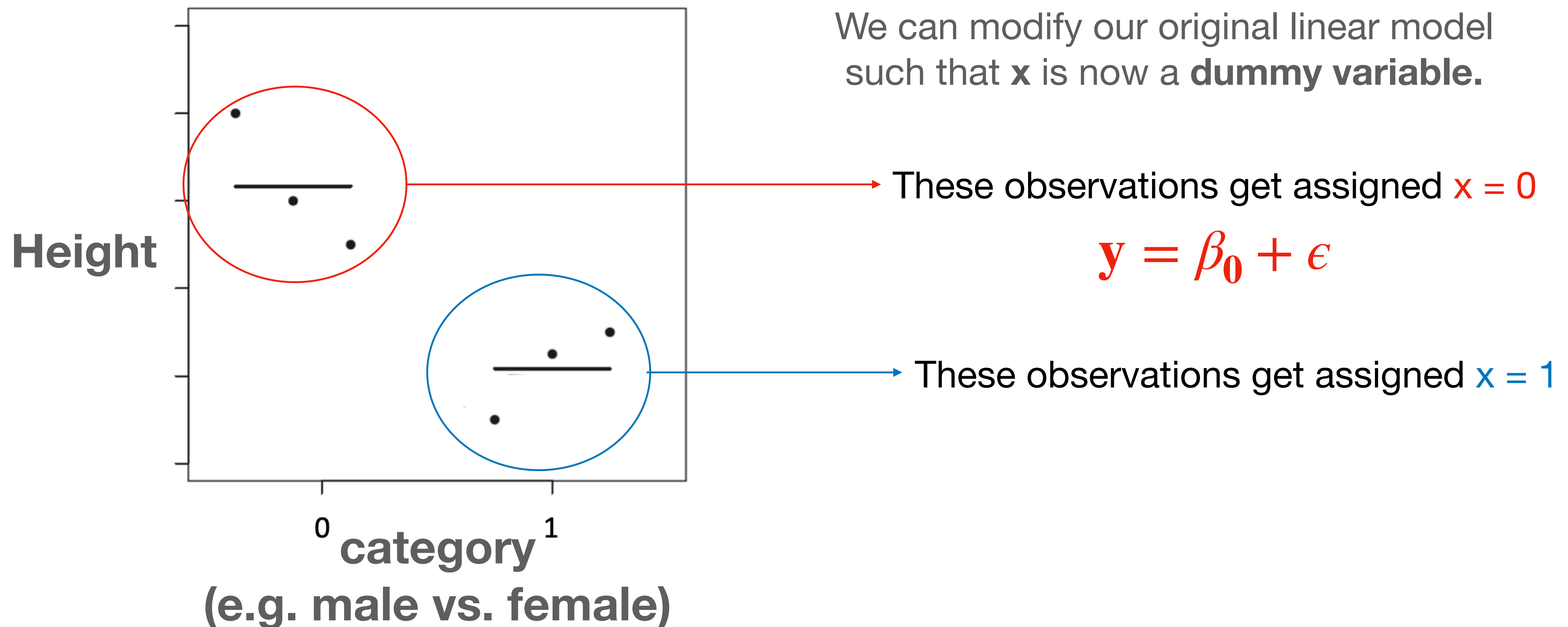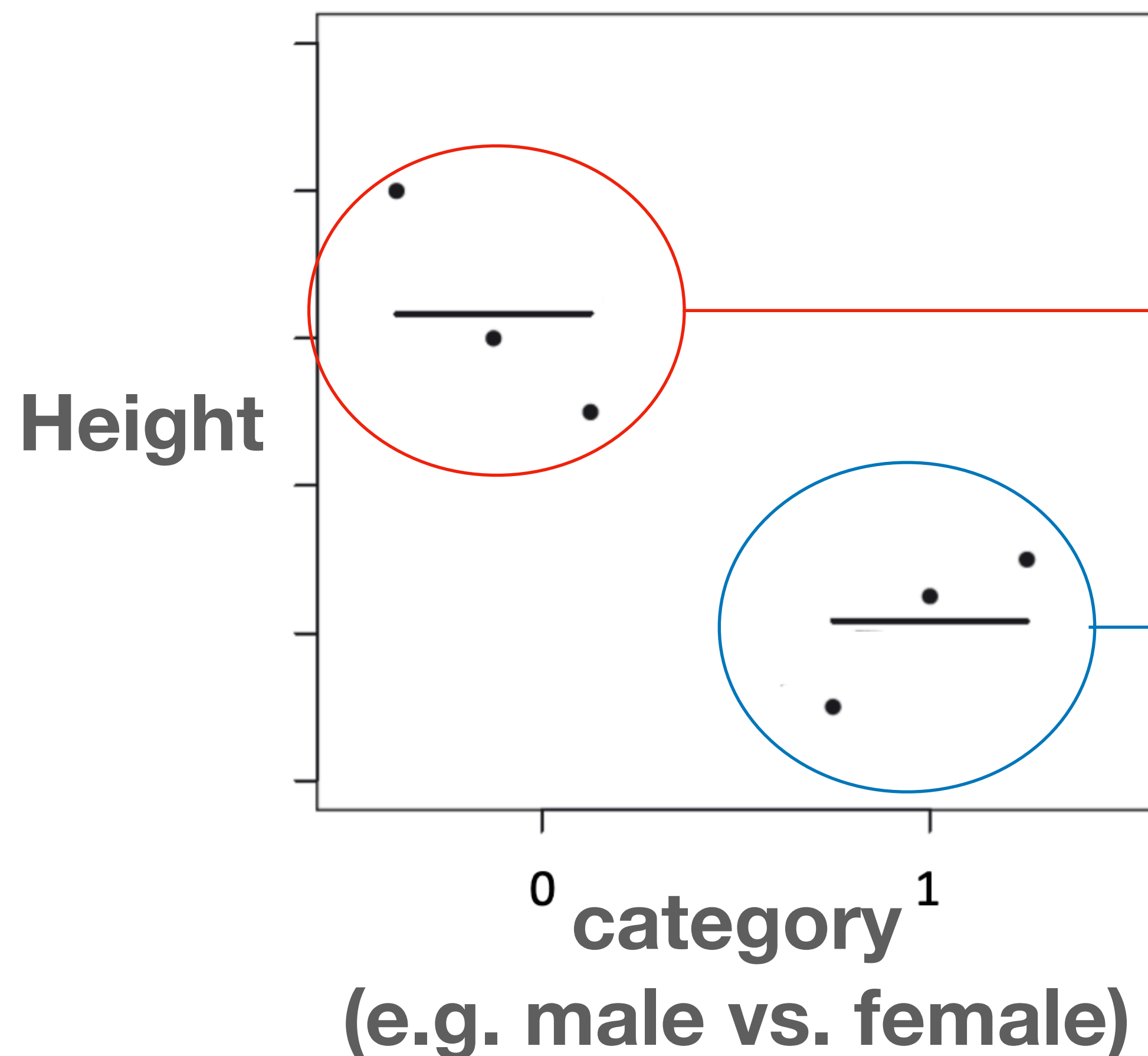We can modify our original linear model such that **x** is now a **dummy variable.**

These observations get assigned x = 0

Height

0          1
category
(e.g. male vs. female)

# Multiple linear regression with categorical predictors (ANOVA)

$$y = \beta_0 + \beta_1 x + \epsilon$$

We can modify our original linear model such that **x** is now a **dummy variable.**

These observations get assigned x = 0

$$y = \beta_0 + \epsilon$$

**Height**

0    category    1

**(e.g. male vs. female)**

# Multiple linear regression with categorical predictors (ANOVA)

$$y = \beta_0 + \beta_1 x + \epsilon$$

We can modify our original linear model such that **x** is now a **dummy variable.**

Height

category
(e.g. male vs. female)

These observations get assigned x = 0

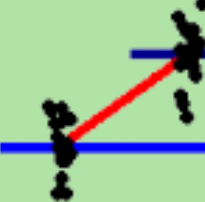$$y = \beta_0 + \epsilon$$

These observations get assigned x = 1

# Multiple linear regression with categorical predictors (ANOVA)

$$y = \beta_0 + \beta_1 x + \epsilon$$

We can modify our original linear model such that **x** is now a **dummy variable.**



**Height**

0    **category**    1
**(e.g. male vs. female)**

These observations get assigned x = 0

$$y = \beta_0 + \epsilon$$

These observations get assigned x = 1
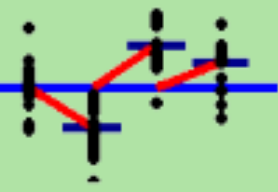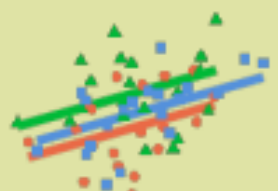
$$y = \beta_0 + \beta_1 + \epsilon$$

# A little secret…

## Common statistical tests are linear models
*Last updated: 02 April, 2019*

See worked examples and more details at the accompanying notebook: https://lindeloev.github.io/tests-as-linear

| | Common name | Built-in function in R | Equivalent linear model in R | Exact? | The linear model in words | Icon |
|---|---|---|---|---|---|---|
| **Simple regression: lm(y ~ 1 + x)** | **y is independent of x**<br>P: One-sample t-test<br>N: Wilcoxon signed-rank | t.test(y)<br>wilcox.test(y) | lm(y ~ 1)<br>lm(signed_rank(y) ~ 1) | ✓<br>for N >14 | One number (intercept, i.e., the mean) predicts **y**.<br>- (Same, but it predicts the *signed rank* of **y**.) | |
| | **P: Paired-sample t-test**<br>N: Wilcoxon matched pairs | t.test($y_1$, $y_2$, paired=TRUE)<br>wilcox.test($y_1$, $y_2$, paired=TRUE) | lm($y_2$ - $y_1$ ~ 1)<br>lm(signed_rank($y_2$ - $y_1$) ~ 1) | ✓<br>for N >14 | One intercept predicts the pairwise $y_2$-$y_1$ differences.<br>- (Same, but it predicts the *signed rank* of $y_2$-$y_1$.) | |
| | **y ~ continuous x**<br>P: Pearson correlation<br>N: Spearman correlation | cor.test(x, y, method='Pearson')<br>cor.test(x, y, method='Spearman') | lm(y ~ 1 + x)<br>lm(rank(y) ~ 1 + rank(x)) | ✓<br>for N >10 | One intercept plus **x** multiplied by a number (slope) predicts **y**.<br>- (Same, but with *ranked* **x** and **y**) | |
| | **y ~ discrete x**<br>P: Two-sample t-test<br>P: Welch's t-test<br>N: Mann-Whitney U | t.test($y_1$, $y_2$, var.equal=TRUE)<br>t.test($y_1$, $y_2$, var.equal=FALSE)<br>wilcox.test($y_1$, $y_2$) | lm(y ~ 1 + $G_2$)$^A$<br>gls(y ~ 1 + $G_2$, weights=…$^B$)$^A$<br>lm(signed_rank(y) ~ 1 + $G_2$)$^A$ | ✓<br>✓<br>for N >11 | An intercept for **group 1** (plus a difference if **group 2**) predicts **y**.<br>- (Same, but with one variance *per group* instead of one common.)<br>- (Same, but it predicts the *signed rank* of **y**.) | |

# A little secret…

# General multivariate regression

- **Several** predictor variables

- **Several** response variables

# General multivariate regression

- **Several** predictor variables

- **Several** response variables

$$\mathbf{y_1} = \beta_0 + \beta_{1,1}\mathbf{x_1} + \beta_{2,1}\mathbf{x_2} + \beta_{3,1}\mathbf{x_3} + \ldots + \epsilon_1$$

Variable $\mathbf{y_1}$ is a linear function of $\mathbf{x_1}$, $\mathbf{x_2}$, $\mathbf{x_3}$, **etc.** plus some **noise**

# General multivariate regression

- **Several** predictor variables

- **Several** response variables

$$y_1 = \beta_0 + \beta_{1,1}x_1 + \beta_{2,1}x_2 + \beta_{3,1}x_3 + \ldots + \epsilon_1$$

Variable **y₁** is a linear function of **x₁, x₂, x₃, etc.** plus some **noise**

$$y_2 = \beta_0 + \beta_{1,2}x_1 + \beta_{2,2}x_2 + \beta_{3,2}x_3 + \ldots + \epsilon_2$$

Variable **y₂** is a linear function of **x₁, x₂, x₃, etc.** plus some **noise**

# General multivariate regression

- **Several** predictor variables

- **Several** response variables

$$y_1 = \beta_0 + \beta_{1,1}x_1 + \beta_{2,1}x_2 + \beta_{3,1}x_3 + \ldots + \epsilon_1$$

Variable **y₁** is a linear function of **x₁, x₂, x₃, etc.** plus some **noise**

$$y_2 = \beta_0 + \beta_{1,2}x_1 + \beta_{2,2}x_2 + \beta_{3,2}x_3 + \ldots + \epsilon_2$$

Variable **y₂** is a linear function of **x₁, x₂, x₃, etc.** plus some **noise**

$$y_3 = \beta_0 + \beta_{1,3}x_1 + \beta_{2,3}x_2 + \beta_{3,3}x_3 + \ldots + \epsilon_3$$

Variable **y₃** is a linear function of **x₁, x₂, x₃, etc.** plus some **noise**

# Generalized linear models

# Generalized linear models

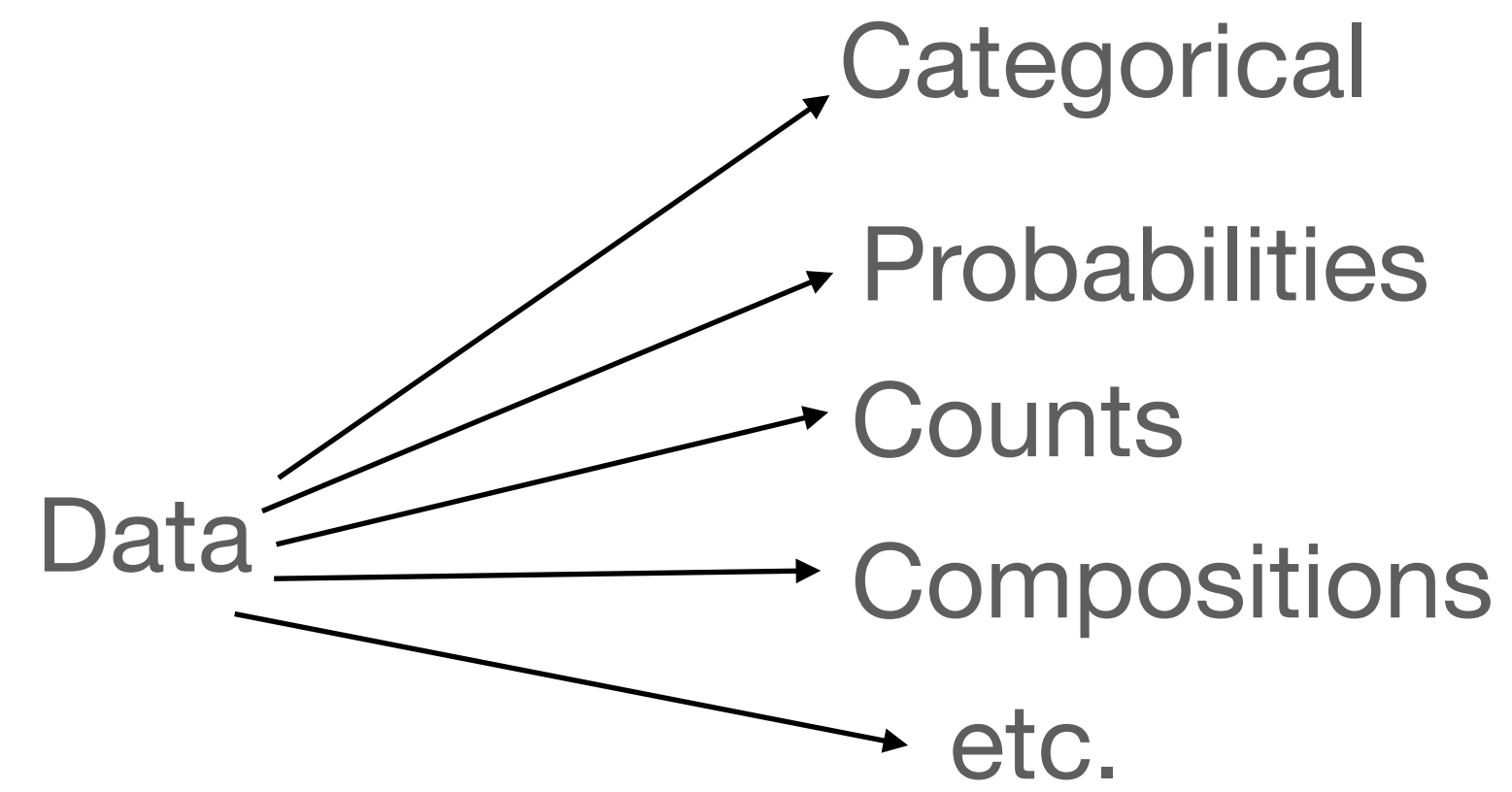$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \epsilon$$

# Generalized linear models

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \epsilon$$

y is continuous between -inf and inf

# Generalized linear models

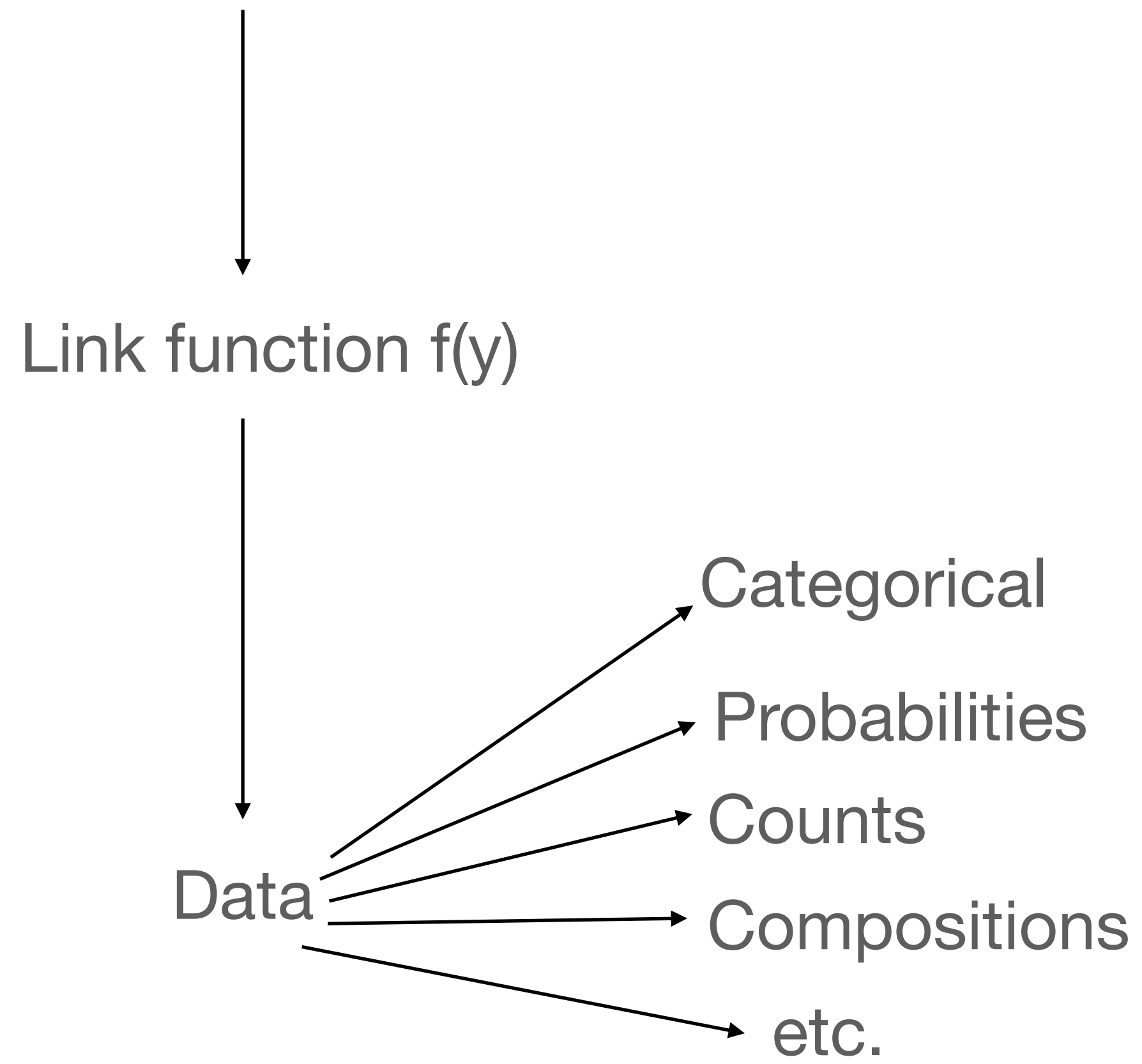$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x_1} + \beta_2 \mathbf{x_2} + \beta_3 \mathbf{x_3} + \ldots + \epsilon$$
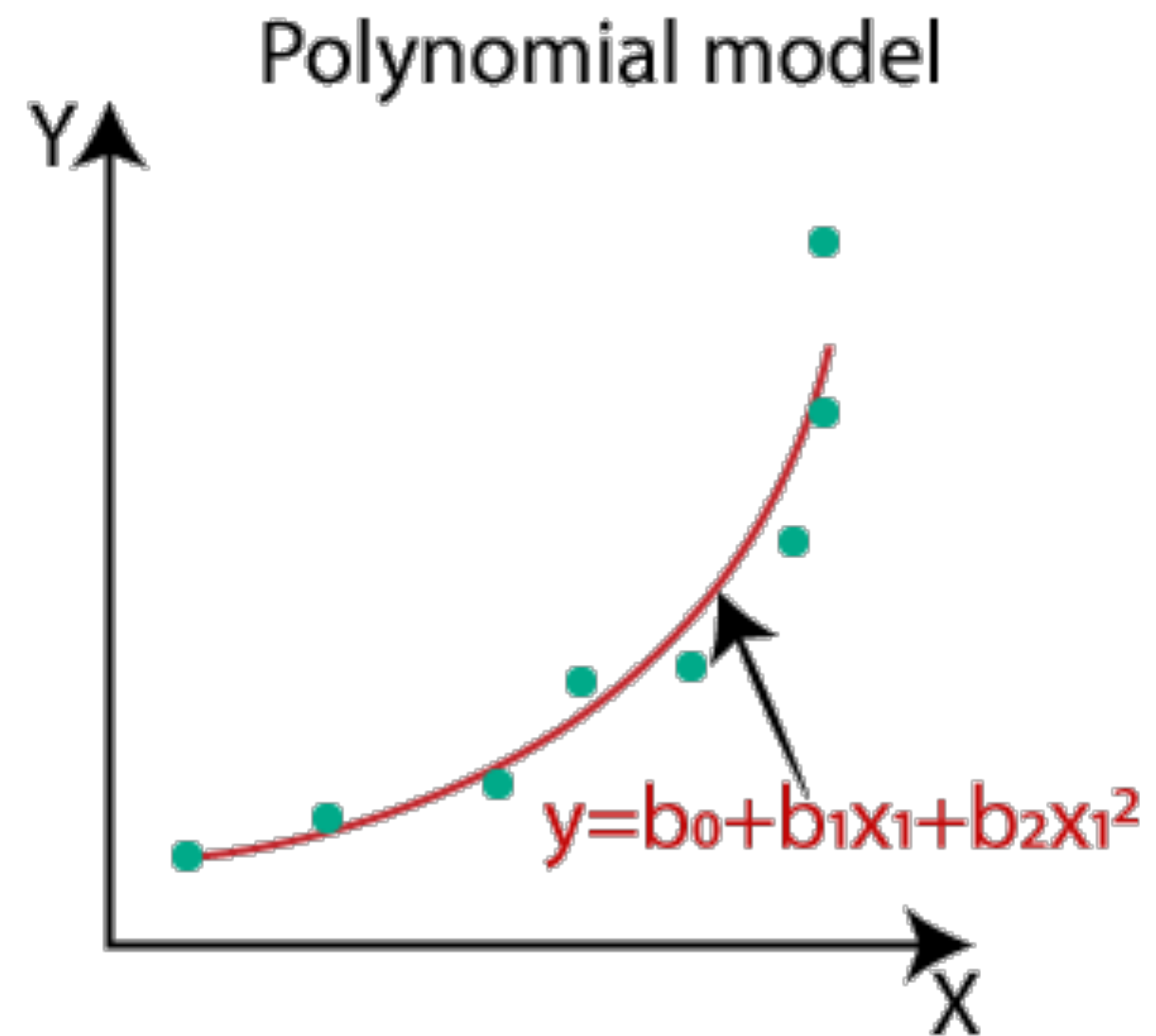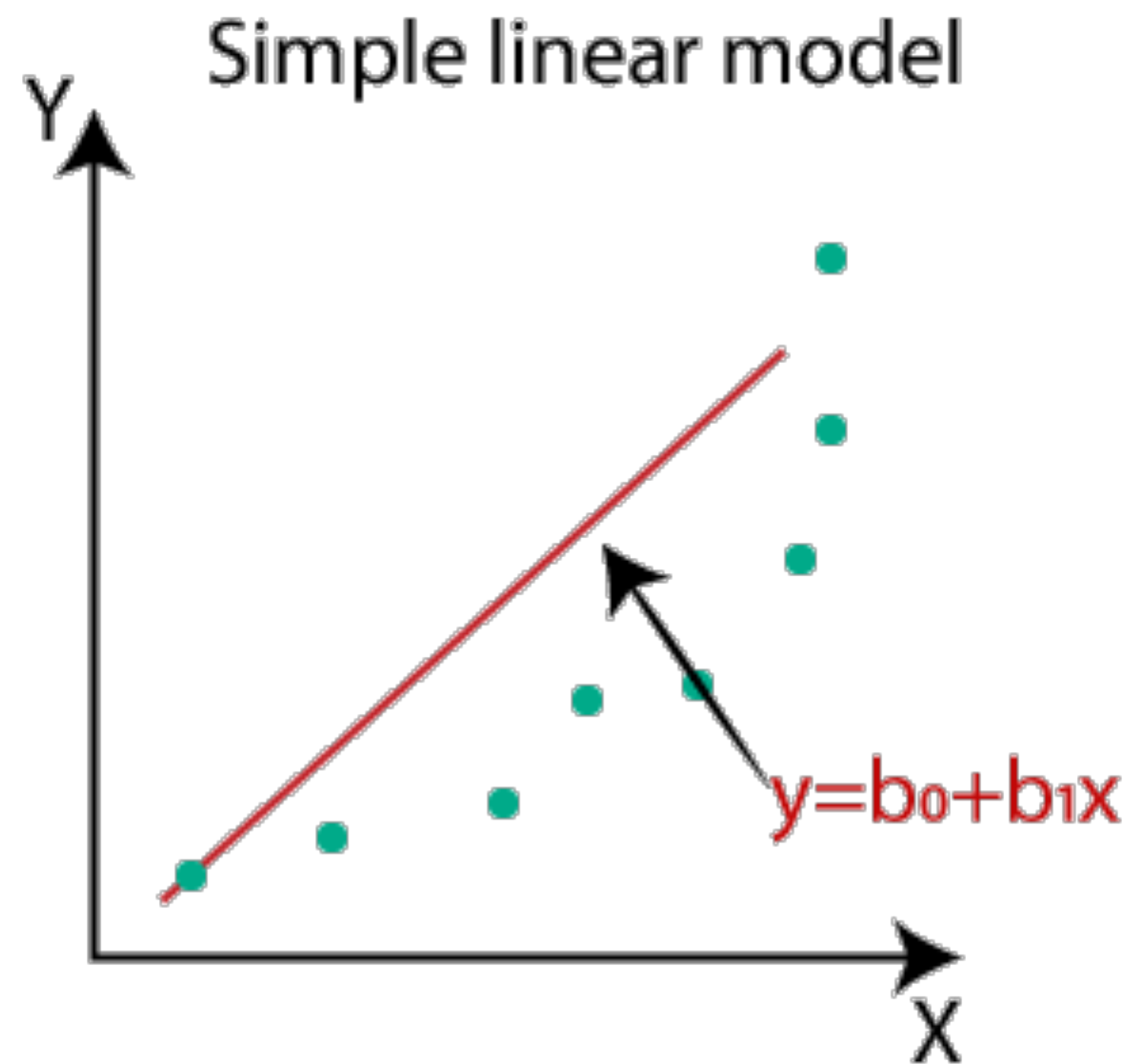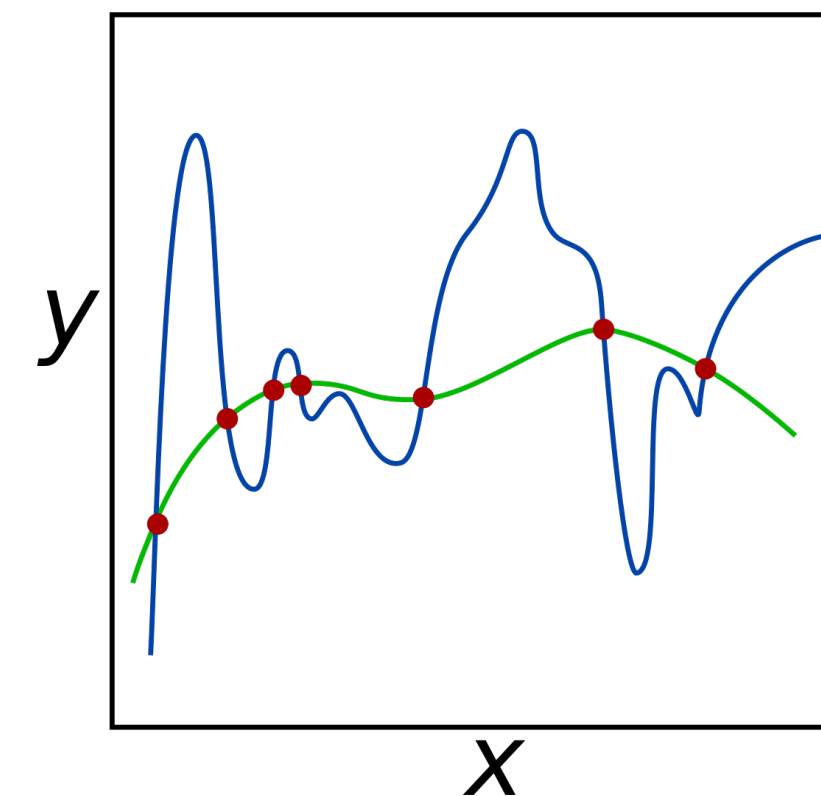
y is continuous between -inf and inf

Data

Categorical

Probabilities

Counts

Compositions

etc.

# Generalized linear models

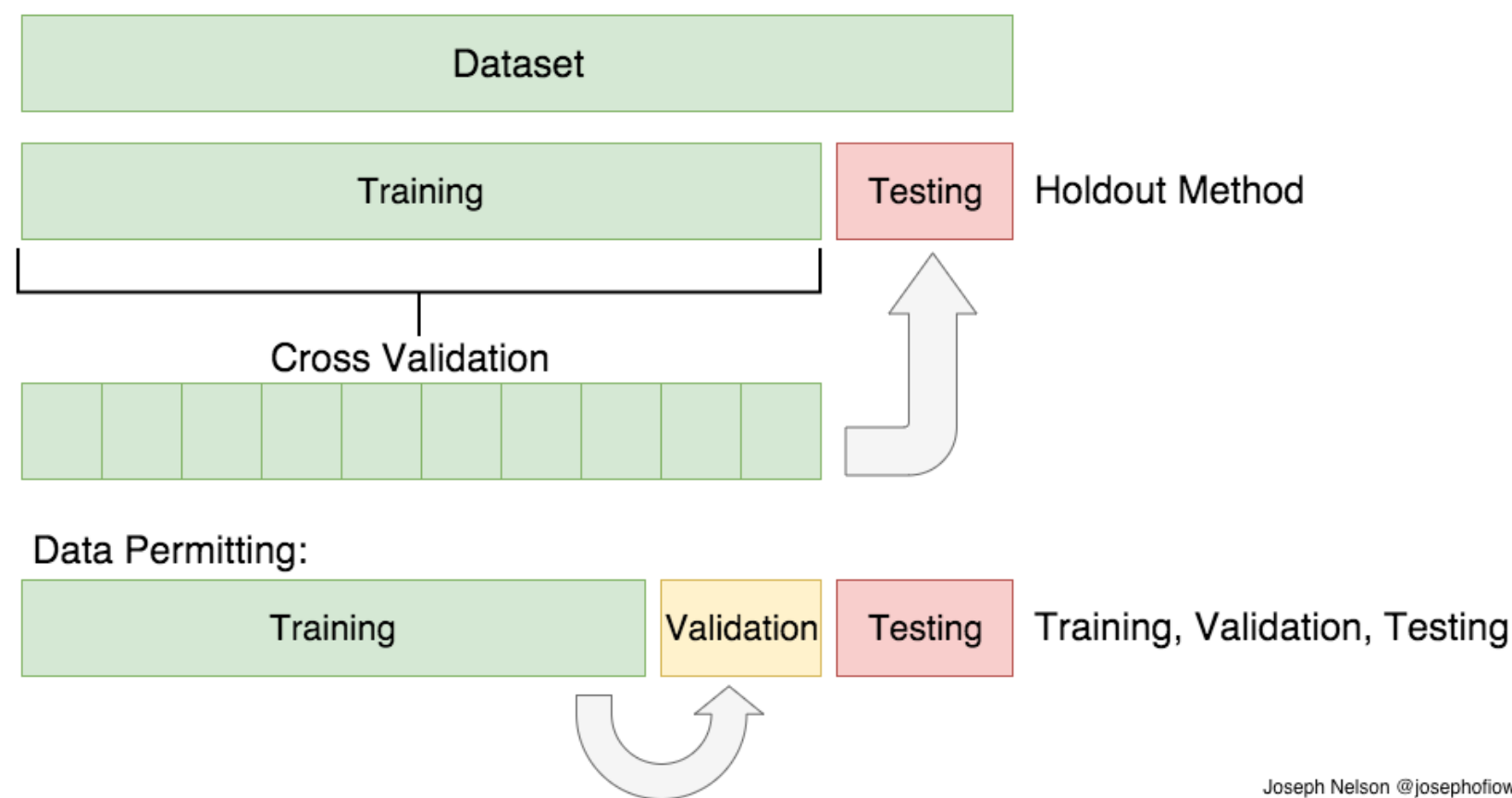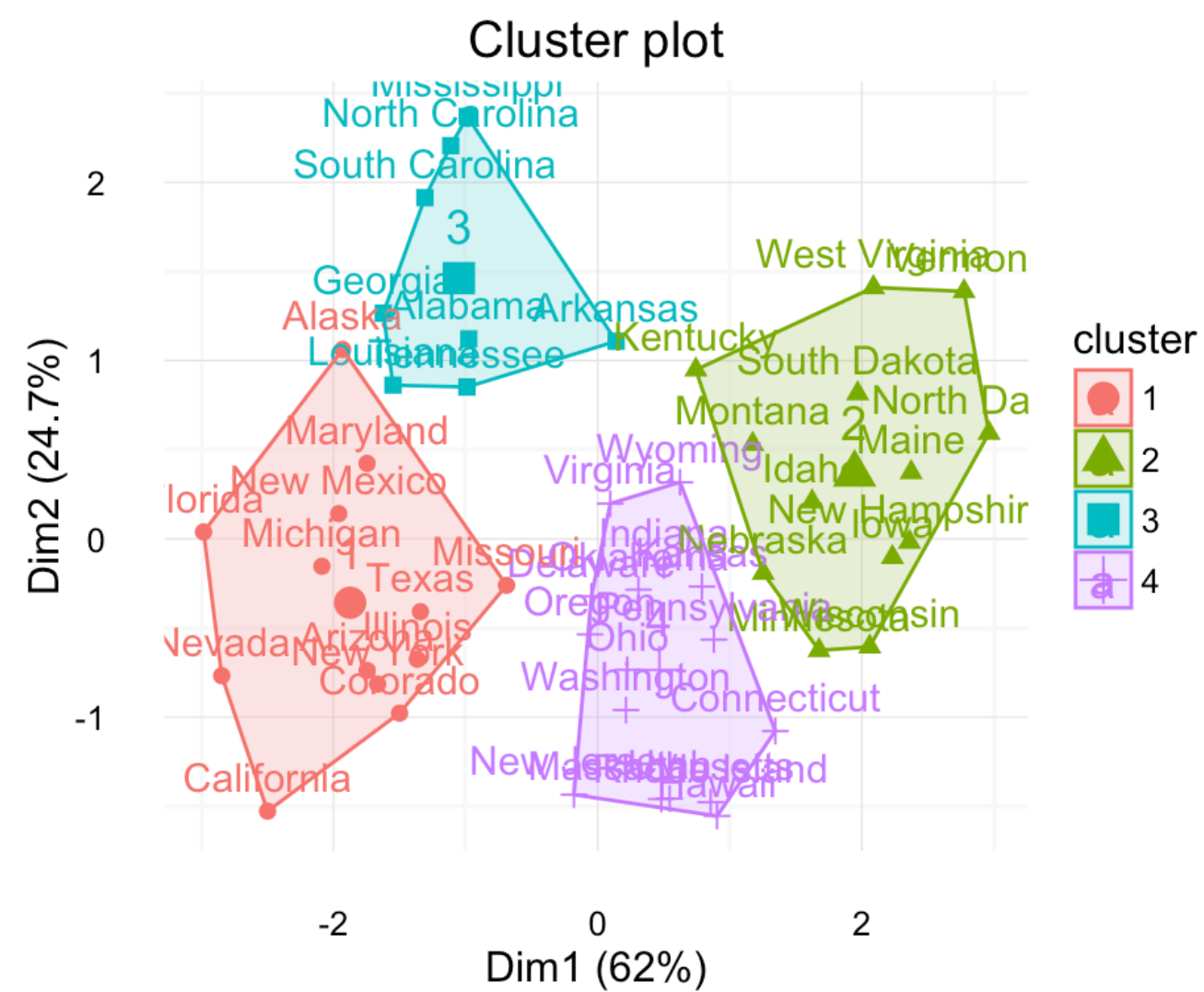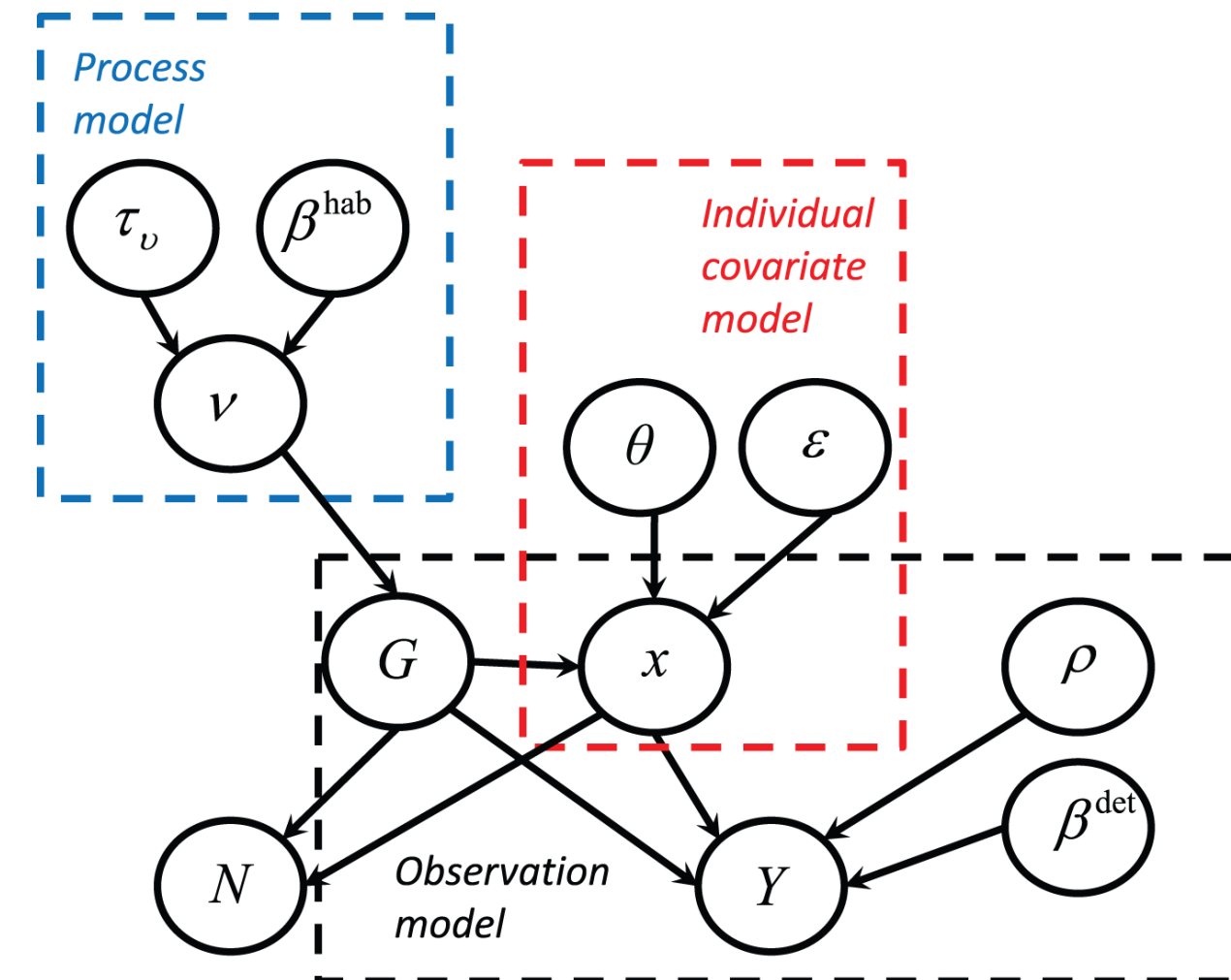$$\mathbf{y} = \beta_0 + \beta_1\mathbf{x_1} + \beta_2\mathbf{x_2} + \beta_3\mathbf{x_3} + \ldots + \epsilon$$

y is continuous between -inf and inf

Link function f(y)

Data

Categorical

Probabilities

Counts

Compositions

etc.

# Non-linear models



Simple linear model

$y = b_0 + b_1 x$

Polynomial model

$y = b_0 + b_1 x_1 + b_2 x_1^2$

# Advanced Topics in Data Analysis



Joseph Nelson @josephofiowa

- **Simulating a linear model in R**

- **Multivariate linear regression**