

# Lesson: Introduction to least squares estimation

Module: Linear regression parameter  
estimation

Photo by [Ryan Hoffman](#) on [Unsplash](#)



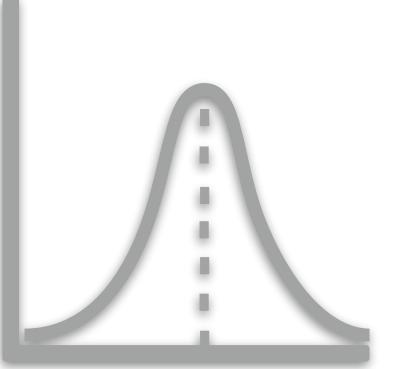
$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

measured response

Unknown, to be estimated

measured predictors

realization from a (normal) random variable

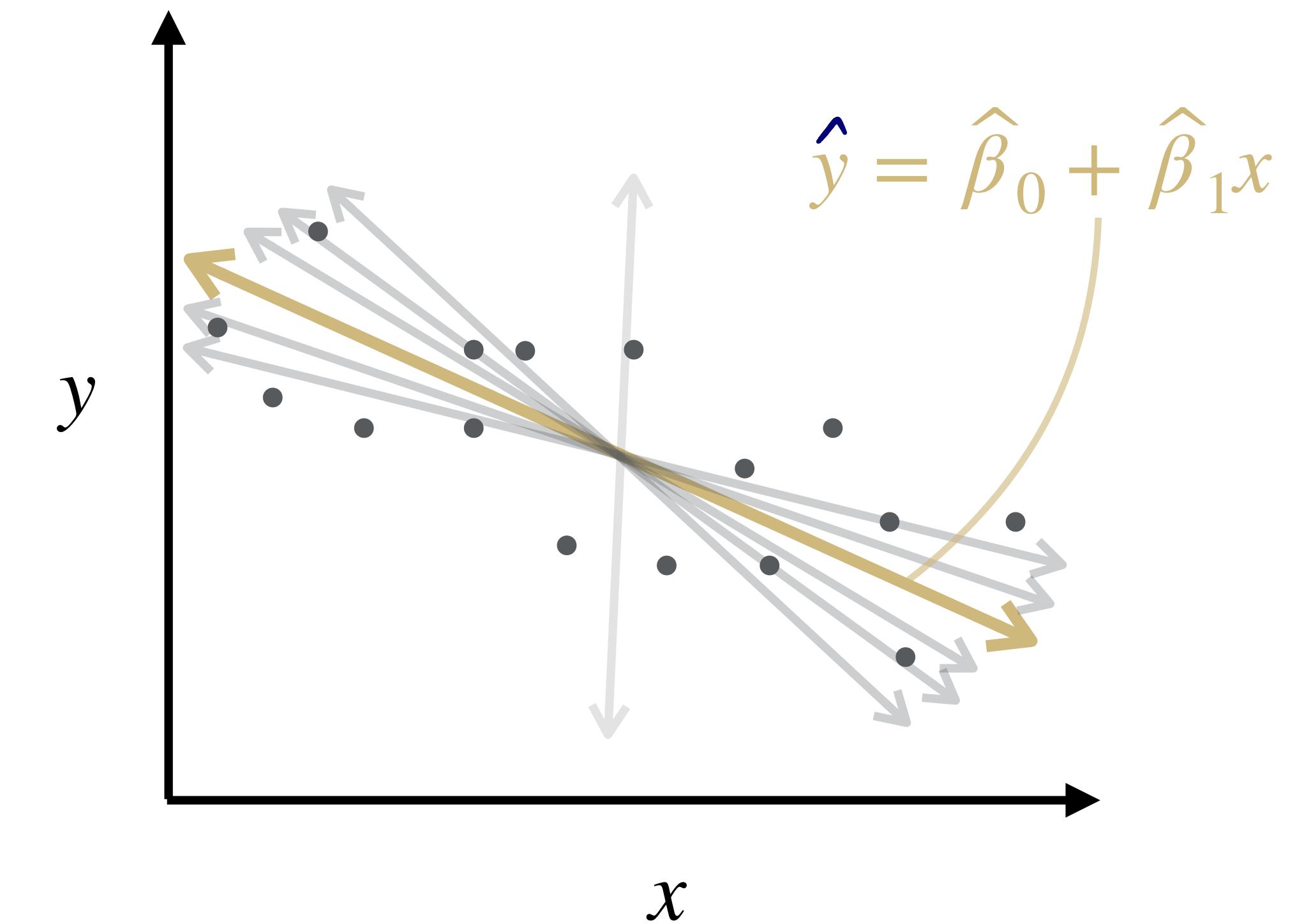

$$\varepsilon_i \sim N(0, \sigma^2)$$

# Estimating Model Parameters

The line of best fit to the data is the line that minimizes the sum of the squared vertical distances between the line and the observed points:

$$\arg \min_{\beta_0, \beta_1} \sum_{i=1}^n \left( y_i - [\beta_0 + \beta_1 x_i] \right)^2$$

measured response      Unknown, to be estimated      measured response



# Estimating Model Parameters

The regression model partitions the response into a systematic component and a random component.

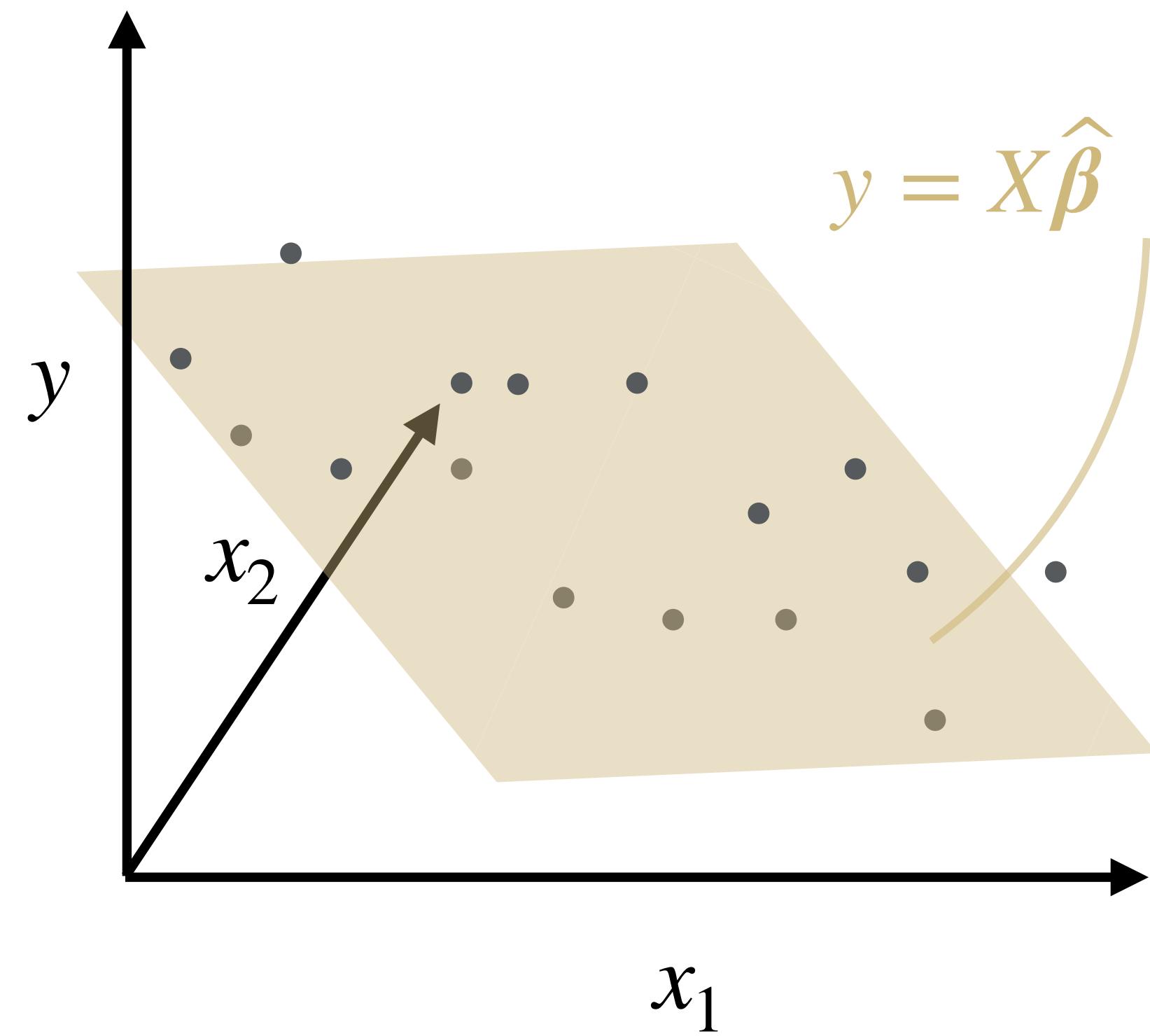
$$\begin{pmatrix} \mathbf{y} \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{pmatrix} \begin{pmatrix} \beta \\ \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon \\ \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$n \times 1$        $n \times (p+1)$        $(p+1) \times 1$

**measured response**      **measured predictors**      **Unknown, to be estimated**

The problem is the find a  $\beta$  so that  $X\beta$  is as close as possible to  $y$ .

# Estimating Model Parameters

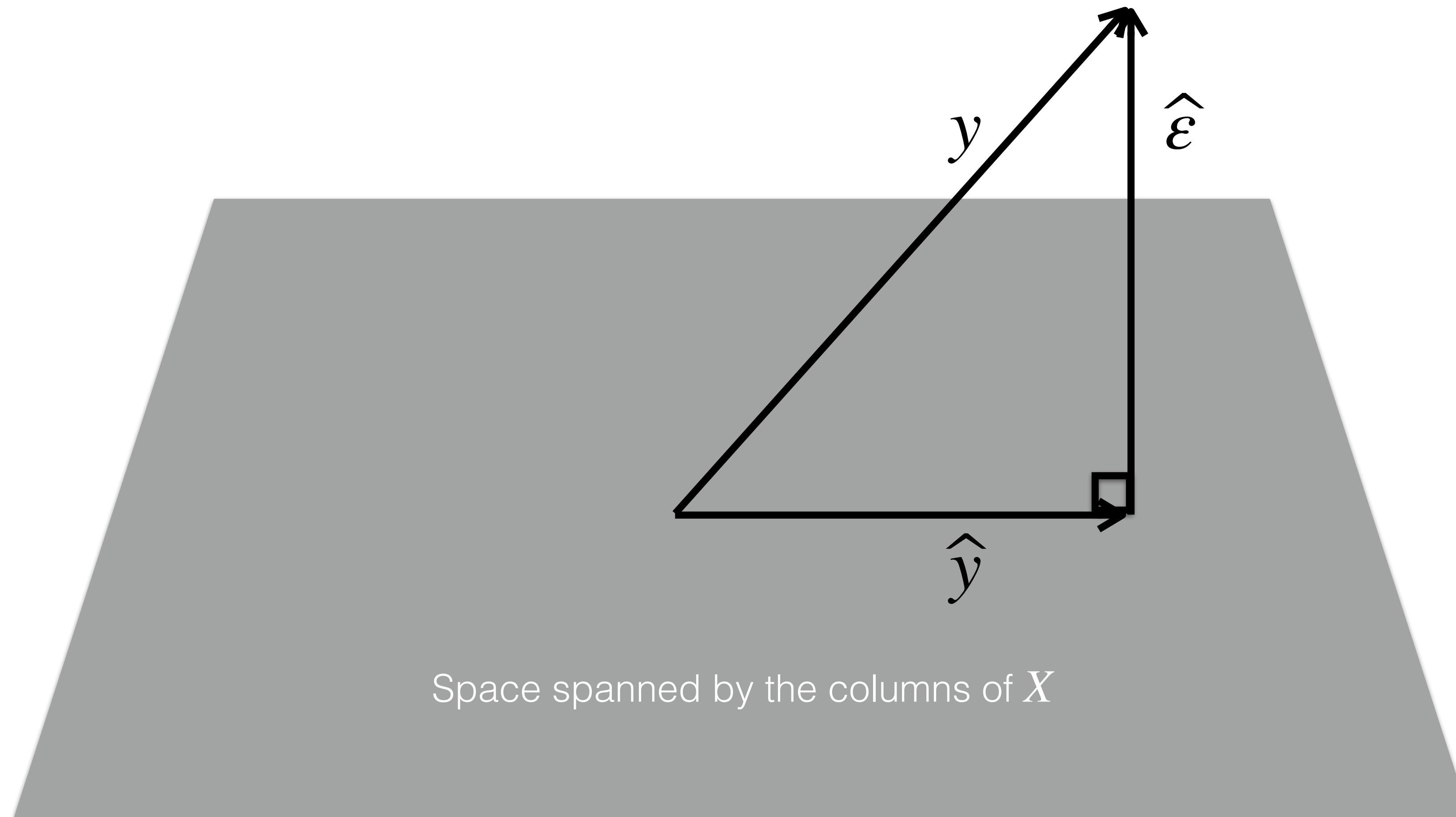


The surface of best fit to the data is the surface that minimizes the sum of the squared vertical distances between the surface and the observed points:

$$\arg \min_{\beta} \|Y - X\beta\|^2 = (Y - X\beta)^T (Y - X\beta)$$
$$\dots = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2$$

The problem is to find a  $\beta$  so that  $X\beta$  is as close as possible to  $y$ .

# Estimating Model Parameters



# Lesson: Linear algebra for least squares

Module: Linear regression parameter  
estimation

Photo by [Ryan Hoffman](#) on [Unsplash](#)



Let  $X$  be an  $m \times n$  matrix,  $\mathbf{v}$  be  $n \times 1$ , and  $\mathbf{y}$  be  $m \times 1$ . Then:

$$\begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ x_{m1} & \ddots & \ddots & x_{mn} \end{pmatrix}$$

1. Lemma 1: Then  $X^T X$  is symmetric, i.e.,  $(X^T X)^T = X^T X$ .

$$(X^T X)^T = X^T (X^T)^T = X^T X$$

2. Lemma 2: Let  $\mathbf{y} = X\mathbf{v}$ . Then  $\frac{\partial \mathbf{y}}{\partial \mathbf{v}} = X$  and  $\frac{\partial \mathbf{y}^T}{\partial \mathbf{v}} = X^T$

$$y_i = x_{i1}v_1 + x_{i2}v_2 + \dots + x_{in}v_n$$

3. Lemma 3: Let  $c = \mathbf{v}^T (X^T X) \mathbf{v}$ . Then  $\frac{\partial c}{\partial \mathbf{v}} = 2X^T X \mathbf{v}$

$$(f(x) = 5x^2, f'(x) = 2 \cdot 5x)$$

# Lesson: Deriving the least squares solution

Module: Linear regression parameter estimation

Photo by [Ryan Hoffman](#) on [Unsplash](#)





**Definition:** The *residuals* are defined as:

$$\hat{\varepsilon}_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})$$

$$\hat{\varepsilon} = (y - \hat{y})$$

**Definition:** The *fitted values* are defined as:

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

*as a function  
of  $\beta_j$*

**Definition:** The *hat matrix*,  $H$ , is defined as:

$$H = X(X^T X)^{-1} X^T$$

The hat matrix is useful in theoretical calculations related to least squares



Photo by Raymond T. on [Unsplash](#)

*Least Squares Estimation:* We define the best estimate of  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  as the one that minimized the sum of the squared residuals:

$$\begin{aligned} RSS &= (\underline{y} - \underline{x}\underline{\beta})^T(\underline{y} - \underline{x}\underline{\beta}) \\ &= (\underline{y}^T - \underline{\beta}^T \underline{x}^T)(\underline{y} - \underline{x}\underline{\beta}) = \underline{y}^T \underline{y} - \underline{y}^T \underline{x} \underline{\beta} - \underline{\beta}^T \underline{x}^T \underline{y} + \underline{\beta}^T \underline{x}^T \underline{x} \underline{\beta} \\ &= \underline{y}^T \underline{y} - 2\underline{\beta}^T \underline{x}^T \underline{y} + \underline{\beta}^T \underline{x}^T \underline{x} \underline{\beta} \end{aligned}$$

Differentiating with respect to  $\underline{\beta}$ , we get:

$$\frac{\partial RSS}{\partial \underline{\beta}} = \underline{0} - 2\underline{x}^T \underline{y} + 2\underline{x}^T \underline{x} \underline{\beta} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \underline{x}^T \underline{x} \underline{\beta} = \underline{x}^T \underline{y}$$

$$\Rightarrow \hat{\underline{\beta}} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y}$$

when  $(\underline{x}^T \underline{x})^{-1}$   
exists

---

$$(\underline{y}^T \underline{x} \underline{\beta})^T = \underline{\beta}^T \underline{x}^T \underline{y}$$



In order to use least squares, we assume that:

1.  $E(\varepsilon_i) = 0$  for all  $i = 1, \dots, n$ .  
$$\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}$$
2.  $E(Y_i) = \underbrace{\mathbf{x}_i^T \boldsymbol{\beta}}$  for all  $i = 1, \dots, n$ .
3.  $Cov(\varepsilon_i, \varepsilon_j) = \begin{cases} 0 & i \neq j \\ \sigma^2 & i = j \end{cases}$
4.  $(X^T X)^{-1}$  exists.

Lesson: Justifying least squares: the Gauss-Markov theorem and maximum likelihood estimation

Module: Linear regression parameter estimation





**The Gauss-Markov Theorem:** Suppose that:

1.  $E(\varepsilon_i) = 0$  for all  $i = 1, \dots, n$ .
2.  $E(Y_i) = \mathbf{x}_i^T \boldsymbol{\beta}$  for all  $i = 1, \dots, n$ .
3.  $Cov(\varepsilon_i, \varepsilon_j) = \begin{cases} 0 & i \neq j \\ \sigma^2 & i = j \end{cases}$
4.  $(X^T X)^{-1}$  exists.

Then  $\hat{\boldsymbol{\beta}}$  is the “best” *unbiased estimator* of  $\boldsymbol{\beta}$ .



# The maximum likelihood estimator.

Suppose that  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ . Then:

1. marginal pdf:

$$f(y_i; \beta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu_y)^2\right\}$$

2. joint pdf:

$$f(\mathbf{y}; \beta) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_y)^2\right\}$$

3. log-likelihood:

$$\ell(\beta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_y)^2$$

RSS

# Lesson: Sums of squares and estimating the error variance

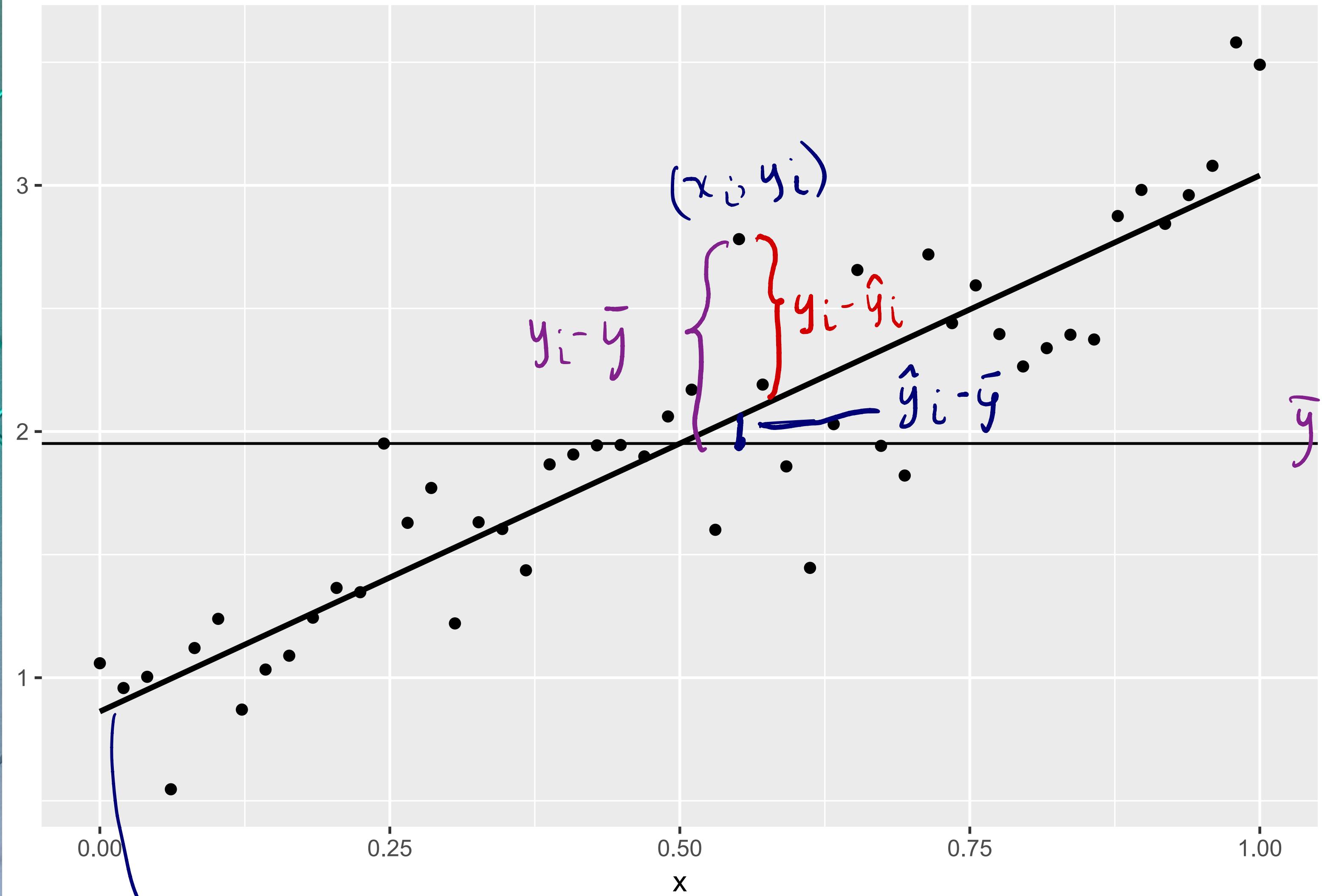
Module: Linear regression parameter estimation

Photo by [Ryan Hoffman](#) on [Unsplash](#)





Photo by Nils Rasmusson on Unsplash



Least squares  
regression line :  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$



Sums of squares:

1. RSS: Residual sum of squares:

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

2. ESS: *Explained* (or regression) sum of squares:

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

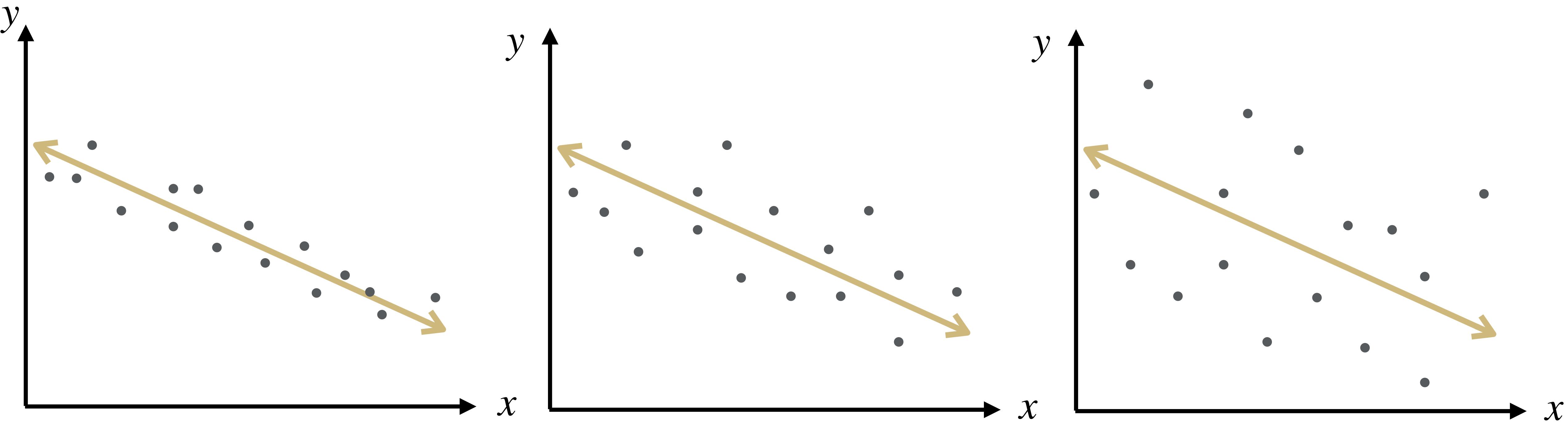
3. TSS: *Total* sum of squares:

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

$$\boxed{TSS = ESS + RSS}$$

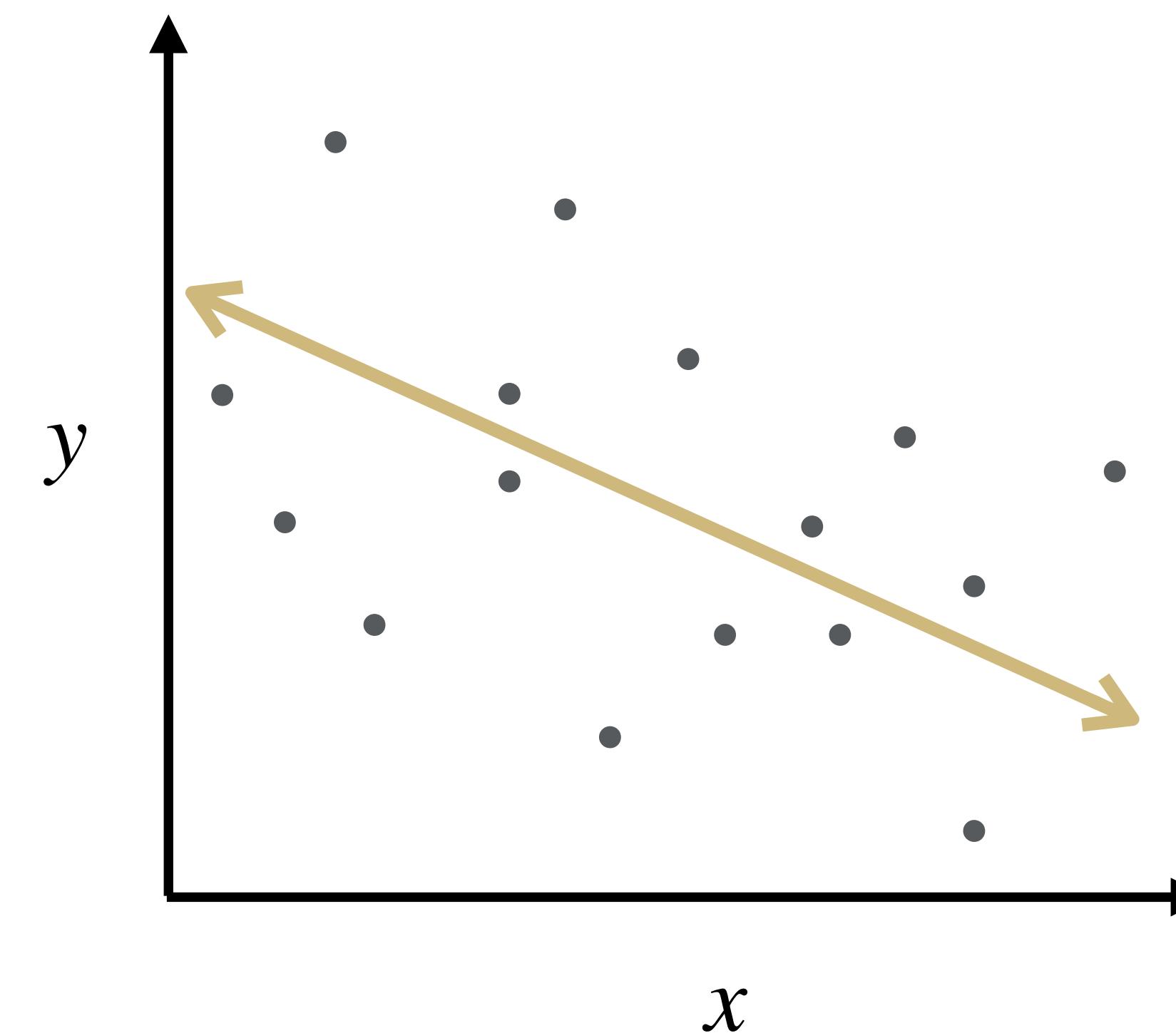
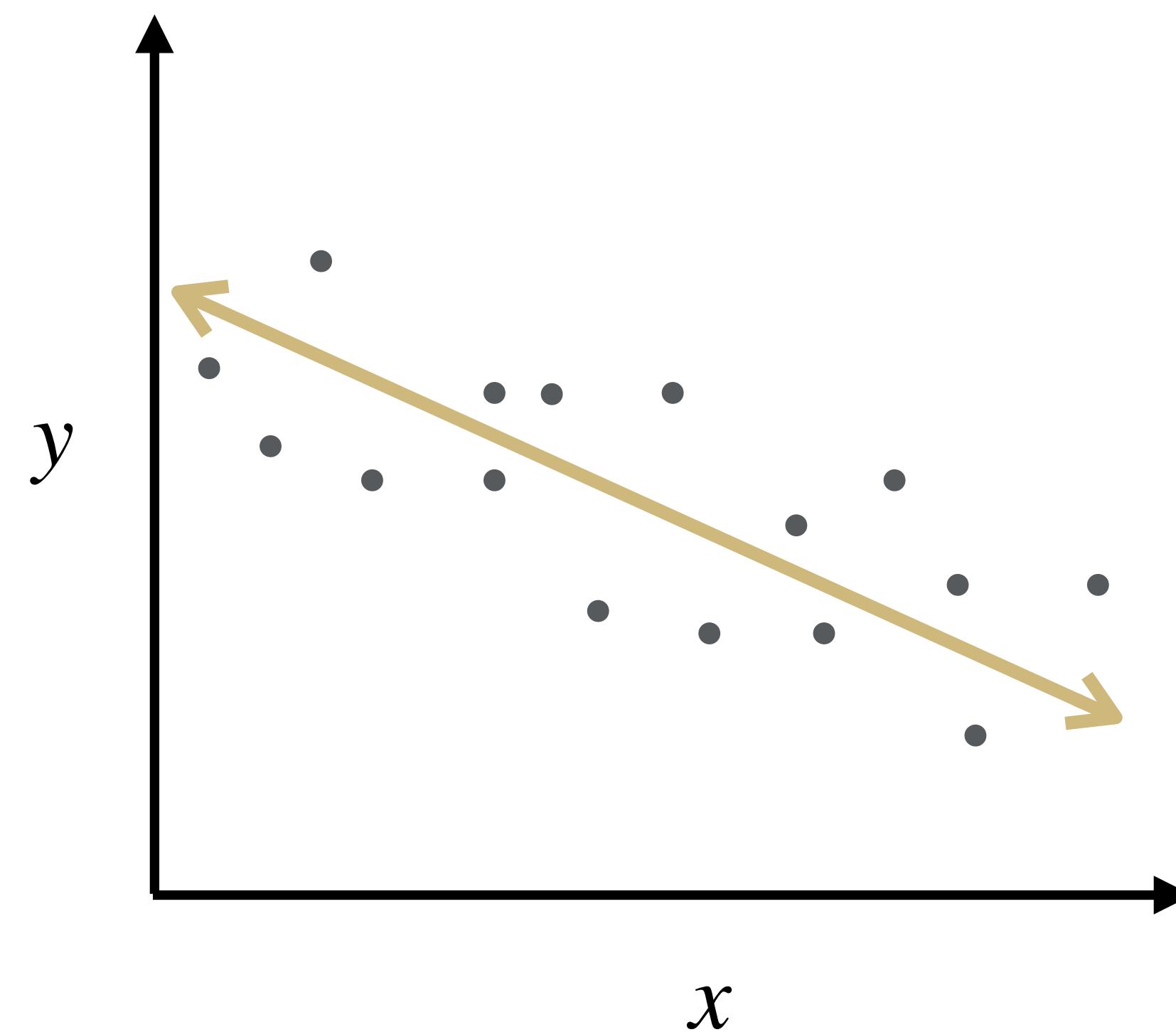
# Sums of Squares

The residual sum of squares RSS can be interpreted as a measure of how much variation in  $y$  is left *unexplained by the model*—that is, how much cannot be attributed to a linear relationship.



# Estimating $\sigma^2$

The parameter  $\sigma^2$  determines the amount of spread about the **true regression line**. Two separate examples:





An estimate of  $\sigma^2$  will be used in statistical inference (e.g., confidence interval formulas and hypothesis testing), presented in the next two sections.

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

Given our interpretation of the RSS, it seems plausible that contribute to an estimator of  $\sigma^2$ :

$$\hat{\sigma}^2 = \frac{RSS}{n - (p+1)}$$



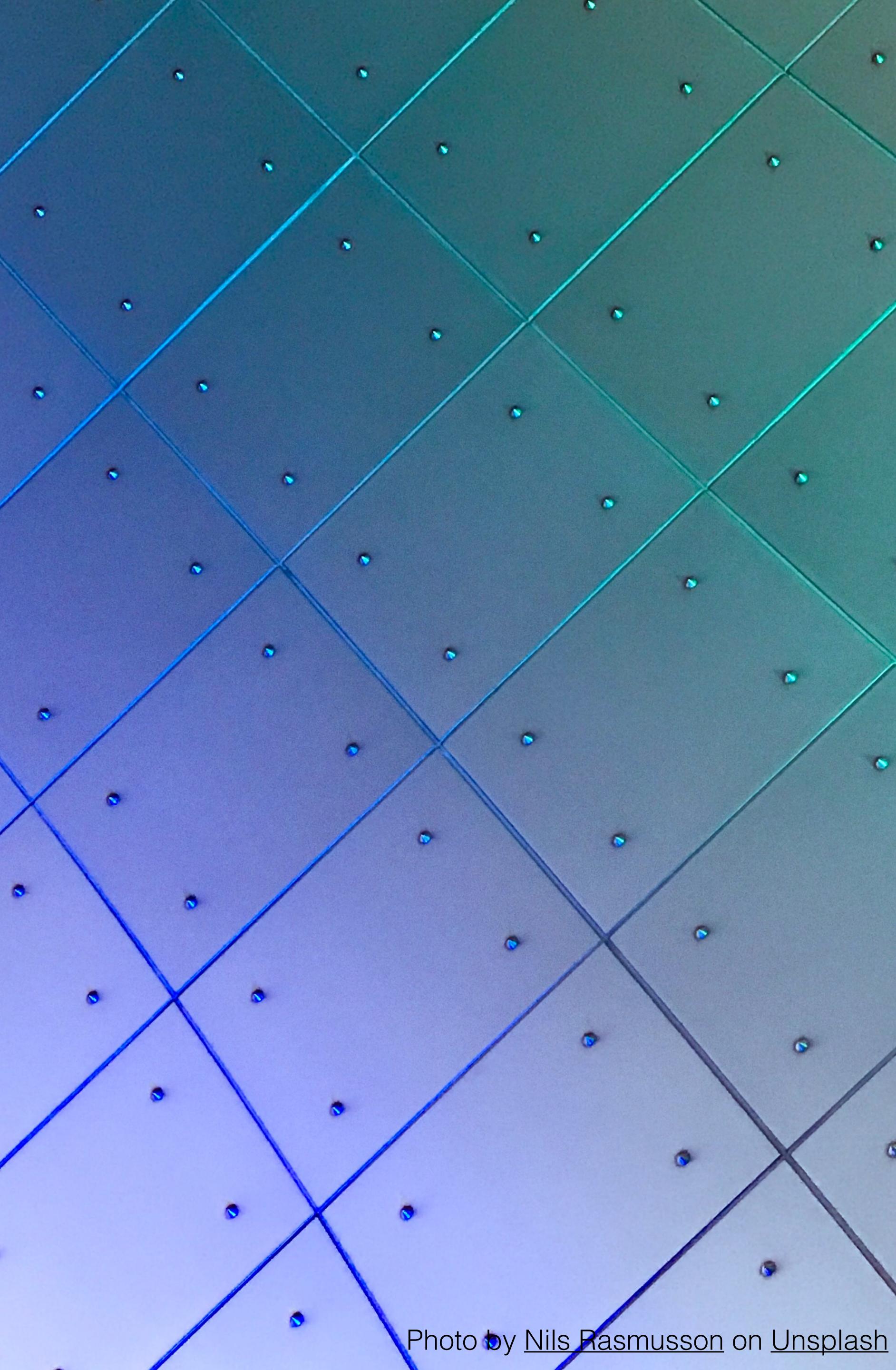
### Note:

1. The divisor  $n - (p + 1)$  in is the number of degrees of freedom (df) associated with RSS and  $\hat{\sigma}^2$ .
2. The RSS has  $n - (p + 1)$  df because  $p + 1$  parameters must first be estimated to compute it, which results in a loss of  $p + 1$  df.
3. Replacing each  $y_i$  in the formula for  $\hat{\sigma}^2$  by the r.v.  $Y_i$  gives a random variable.
4. It can be shown that the r.v.  $\hat{\sigma}^2$  is an *unbiased estimator* for  $\hat{\sigma}^2$ .

Lesson: The coefficient  
of determination, "R-  
squared"

Module: Linear regression parameter  
estimation





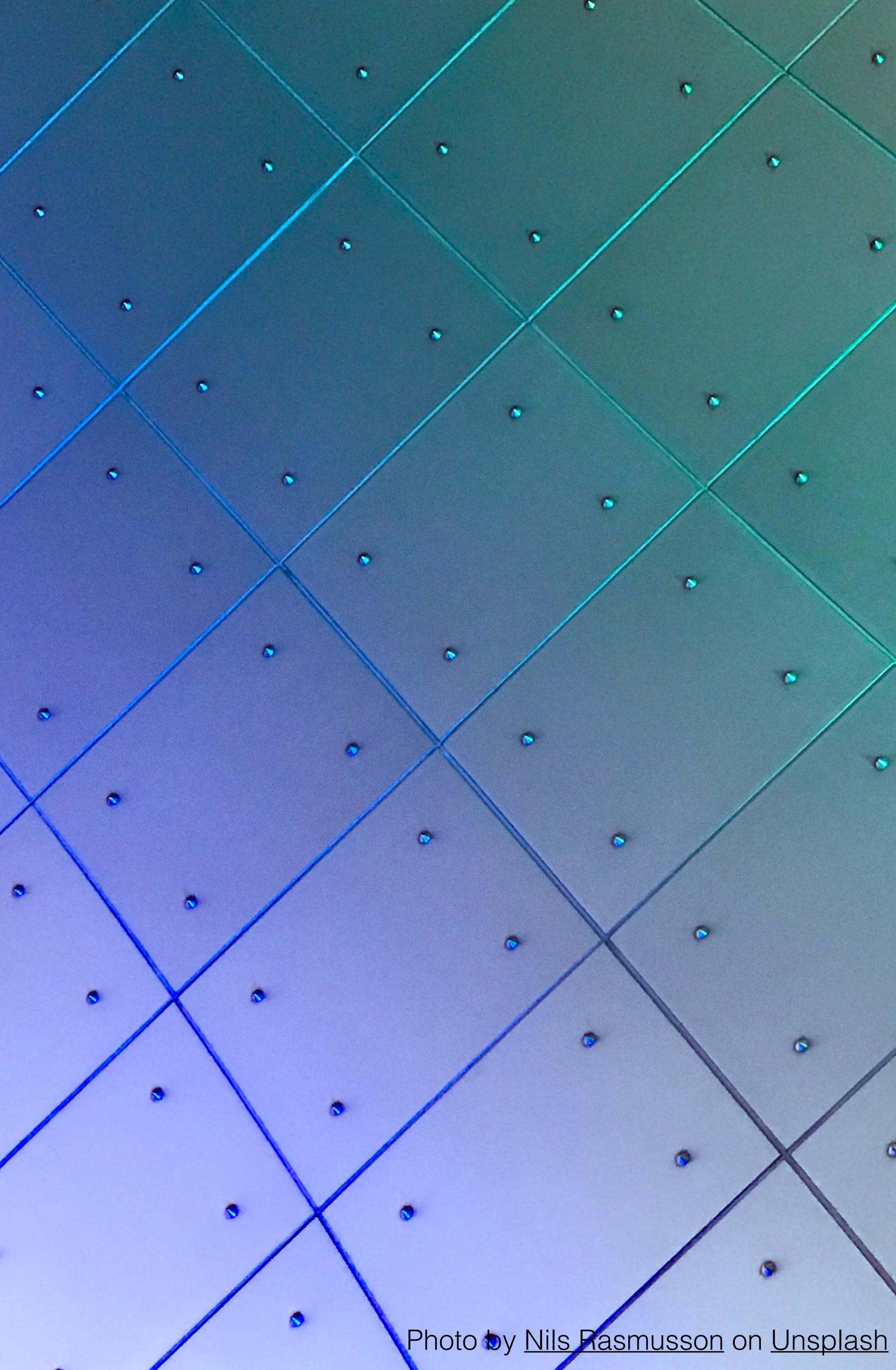
$$TSS = RSS + ESS$$

The *coefficient of determination*,  $R^2$ , is defined as:

$$R^2 = 1 - \frac{RSS}{TSS} = \frac{ESS}{TSS}$$

Note:

- $0 \leq R^2 \leq 1$
- Assuming that the model is correct,  $R^2$  is interpreted as the proportion of observed variation in  $y$  explained by the model.



## Warnings about $R^2$ :

1.  $R^2$  can be close to 1 but the model is the wrong fit for the data.
2.  $R^2$  can be close to 0 even when the model is the correct fit for the data.
3.  $R^2$  should **not** be used to compare models with a different number of predictors.
4.  $R^2$  says nothing about the causal relationship between the predictors and the response.

# Lesson: The problem of non-identifiability

Module: Linear regression parameter  
estimation

Photo by [Ryan Hoffman](#) on [Unsplash](#)





The least squares estimate is the solution to the *normal* equations:

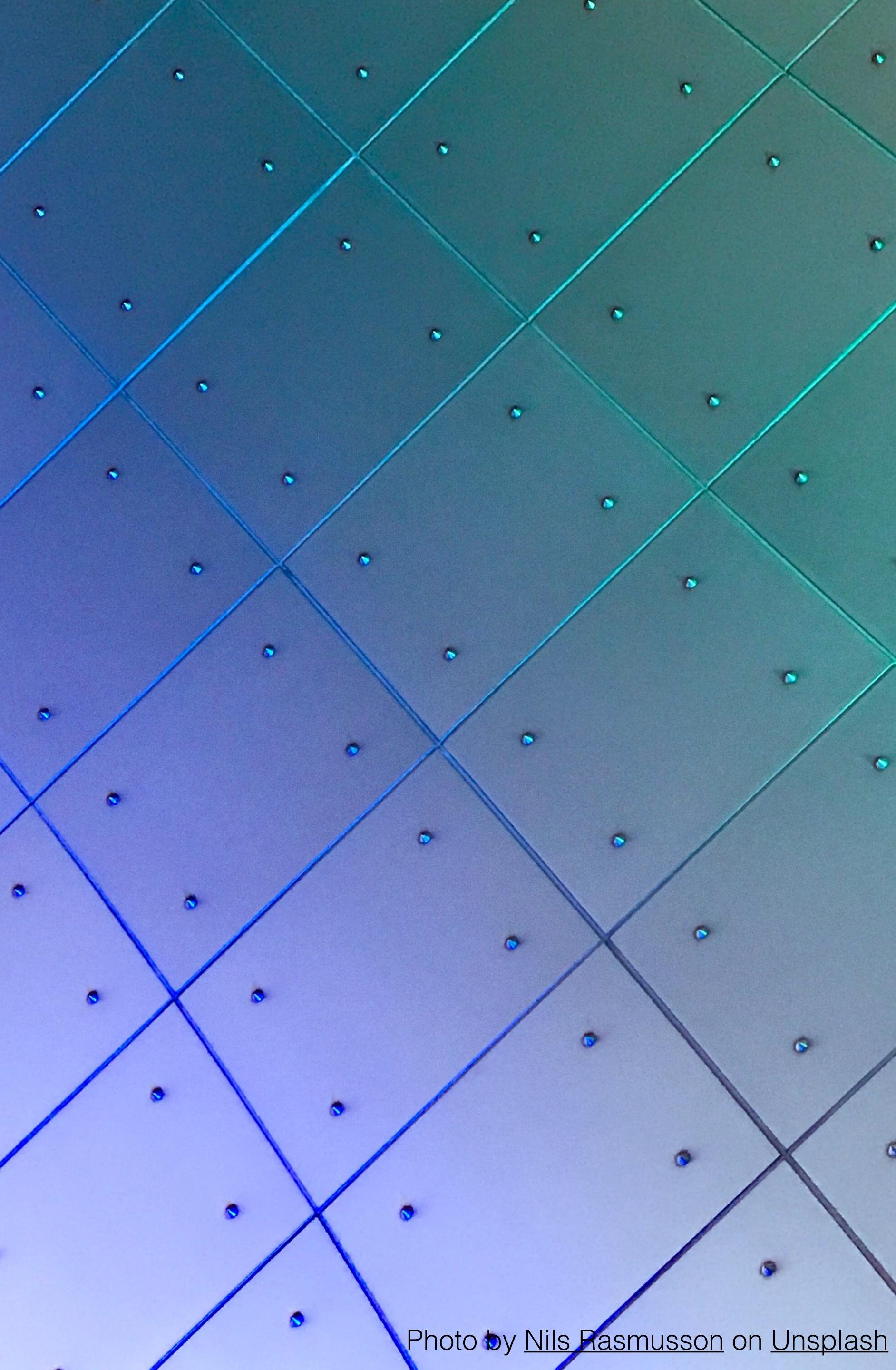
$$X^T X \beta = X^T Y.$$

1. When  $(X^T X)^{-1}$  exists, there is a unique solution,  $\hat{\beta}$ .
2. When  $(X^T X)^{-1}$  does not exist, there will be infinitely many solutions.



**Definition:** When  $(X^T X)^{-1}$  does not exist, the regression model is said to be *non-identifiable* (or, *unidentifiable*).

When does this happen?



## Why might we have non-identifiability?

1. One variable is just a multiple of another.
2. One variable is a *linear combination* of several others.  $\underline{x}_1 = a \underline{x}_2 + b \underline{x}_3 + \dots + p \underline{x}_P$
3. There are more variables than members in the sample.

Note: Near non-identifiability is trickier than exact non-identifiability.