# Linear Regression

- *Learning Types*
- *Linear Regression*

James Balamuta

INMAS Statistical Methods Workshop Fall 2021
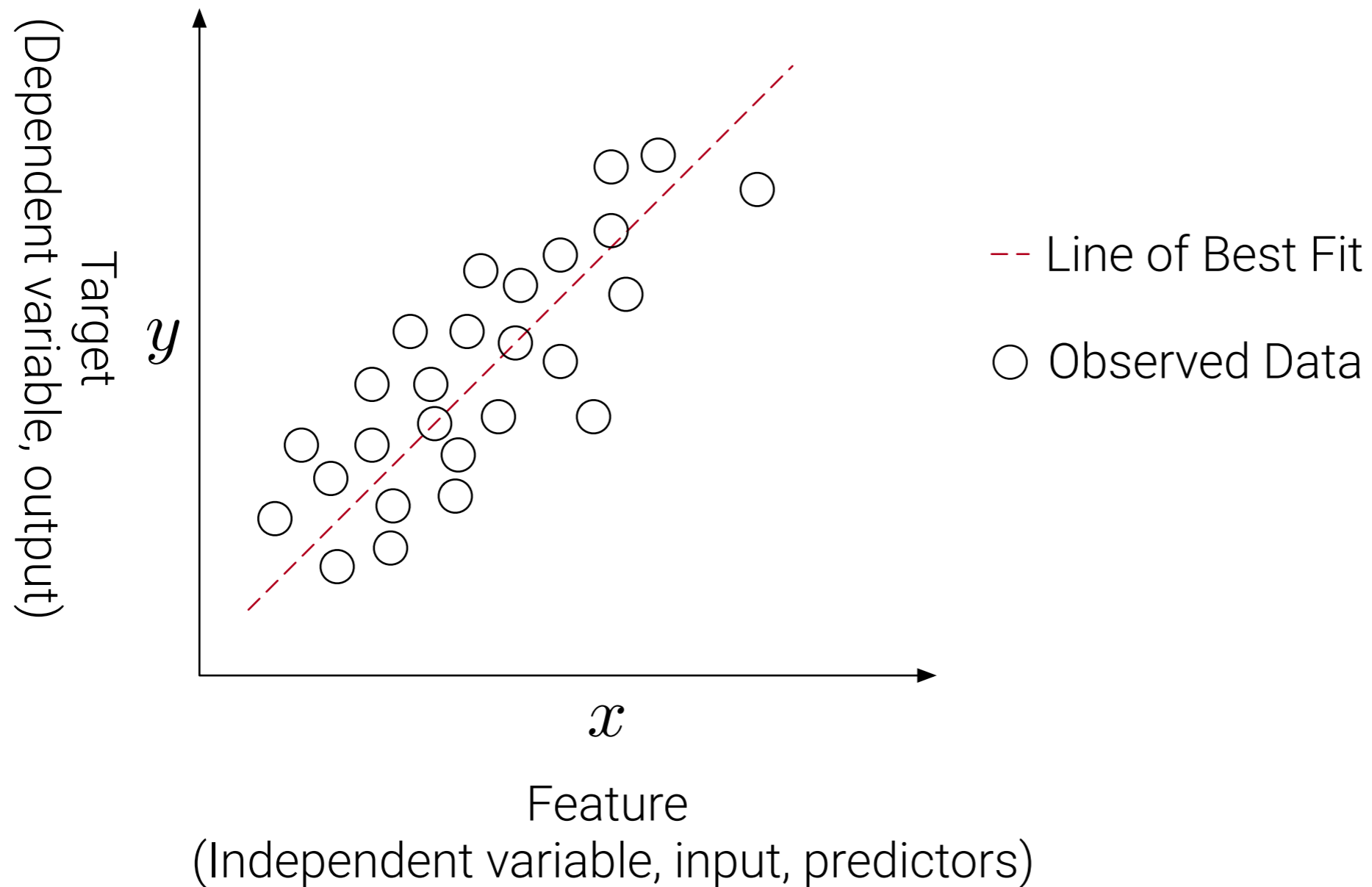
# Lecture Objectives

- *Emphasize* differences between Supervised Learning and Unsupervised Learning

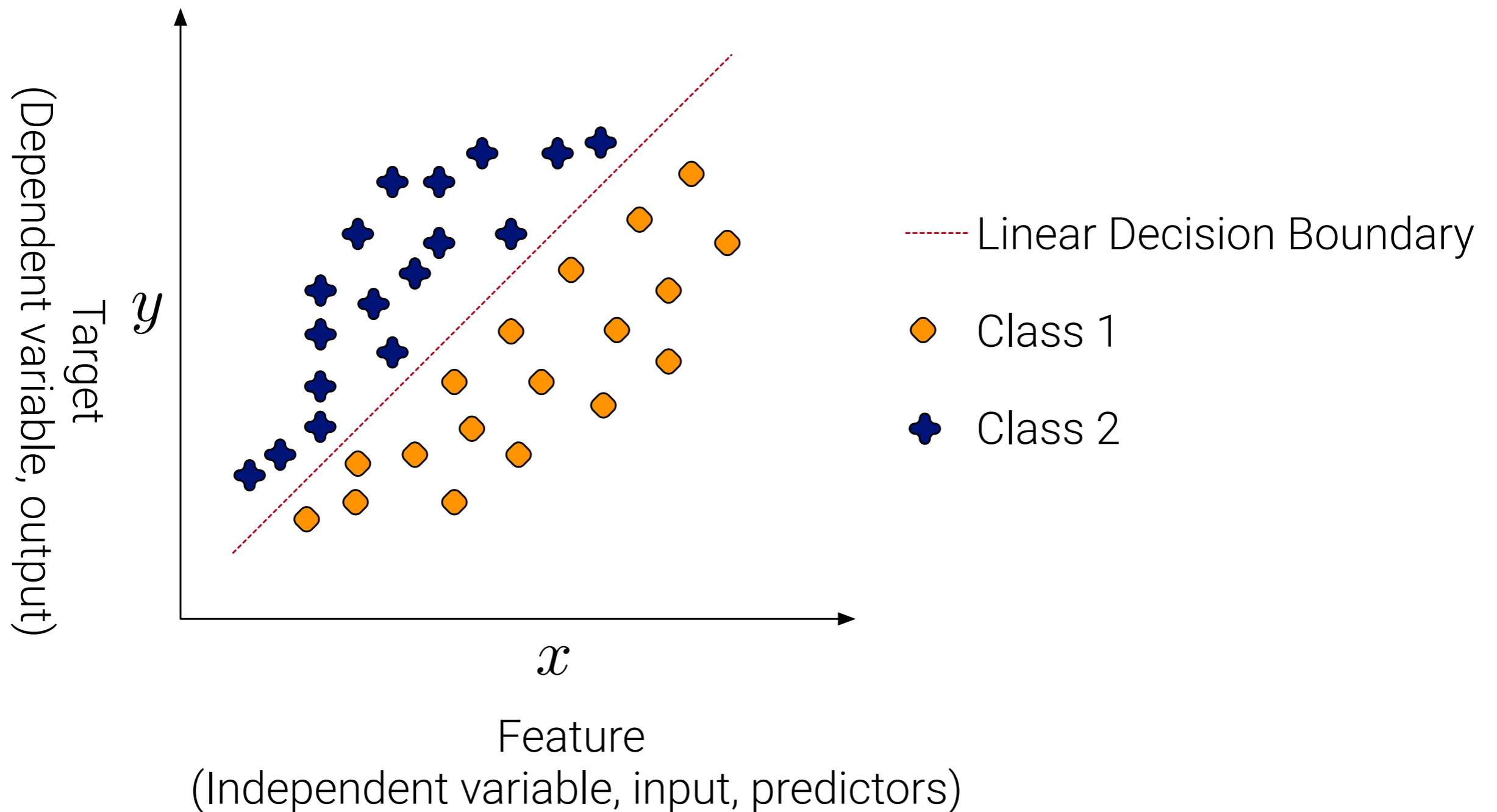- *Discuss* linear regression.

# Learning Types

# Learning Types

- **Supervised Learning** with <u>Labeled</u> Data. (Today)

  - Methods: Regression or classification

  - Objective: To predict a response or outcome.

- **Unsupervised Learning** with <u>Unlabeled</u> Data.

  - Methods: Clustering, Principle Component Analysis (PCA), autoencoders, generative adversarial networks (GANs)

  - Objective: Identify patterns in the data or understand how data was created.

- Best distinction between the two:

  - **Is there a response variable Y?**

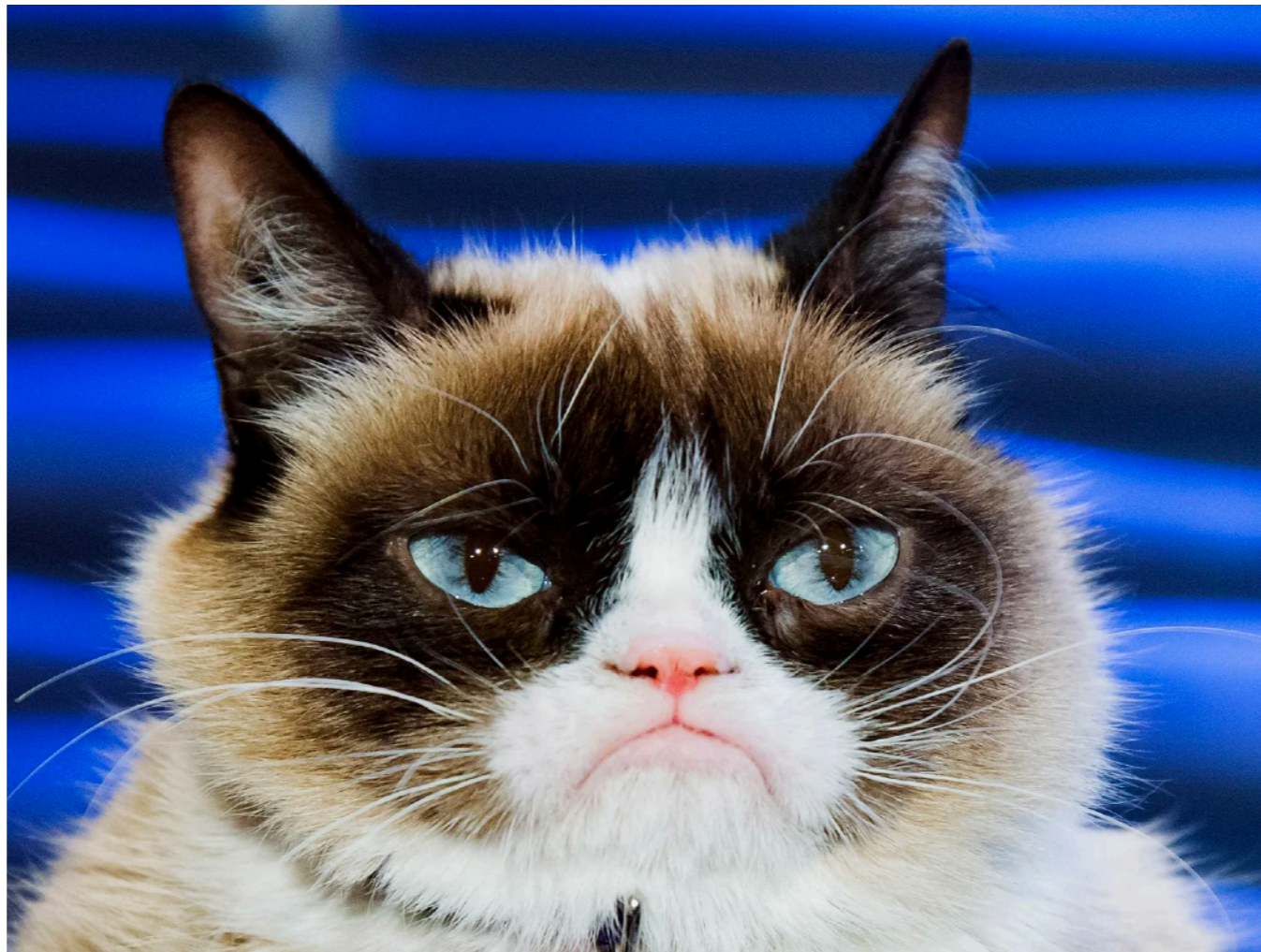# **Supervised Learning**:
# Regression

# **Supervised Learning**:
# Classification



Target
(Dependent variable, output)

$y$

$x$

Feature
(Independent variable, input, predictors)

----- Linear Decision Boundary

○ Class 1

✚ Class 2

# Is Dog?



380

506
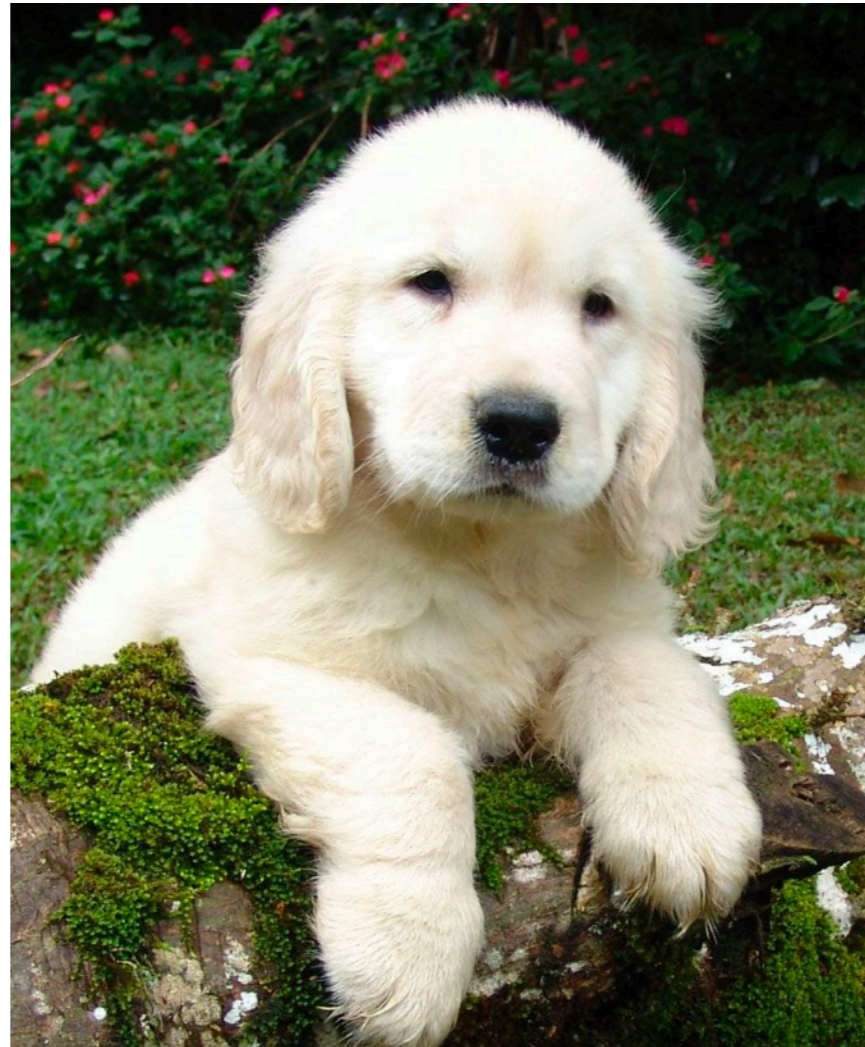
Image dimensions: (380, 506, **3**)

# Is Dog?



965

800

Image dimensions: (965, 800, **3**)

# Unsupervised Learning:
## Clustering

**Data**



$x_2$

$x_1$

**Clustered, K = 4**



$x_2$

$x_1$

# Clustering **is not** Classification

$x$



$y$

**Dog**

**Data**

$P(Y = \mathrm{Dog})$

**Model**

**Output**

# Linear Regression

# Viewpoints

... inference or predictions ...

**Statistician**

**Machine Learning**

**Inference**

**Prediction**
$$\hat{y}_i$$

**Prediction**
$$\hat{y}_i$$

**Inference**

How good are ...
... the predictors selected?
... the diagnostic plots?

How good are ...
... future predictions?

$$Y = f(X) + \varepsilon$$

True
Response

Unknown
Relationship

Unlearnable
Noise

# Simulating Data



```python
import numpy as np
import matplotlib.pyplot as plt

# Set parameters
theta_0 = 3
theta_1 = 2
n = 100

# Generate design matrix
X = np.arange(0, n, 1)

# Create relationship
Y = theta_1*X + theta_0

# Add error
error = 45*np.random.rand(n)
Y = Y + error

# Show graph
plt.scatter(X, Y)
plt.show()
```
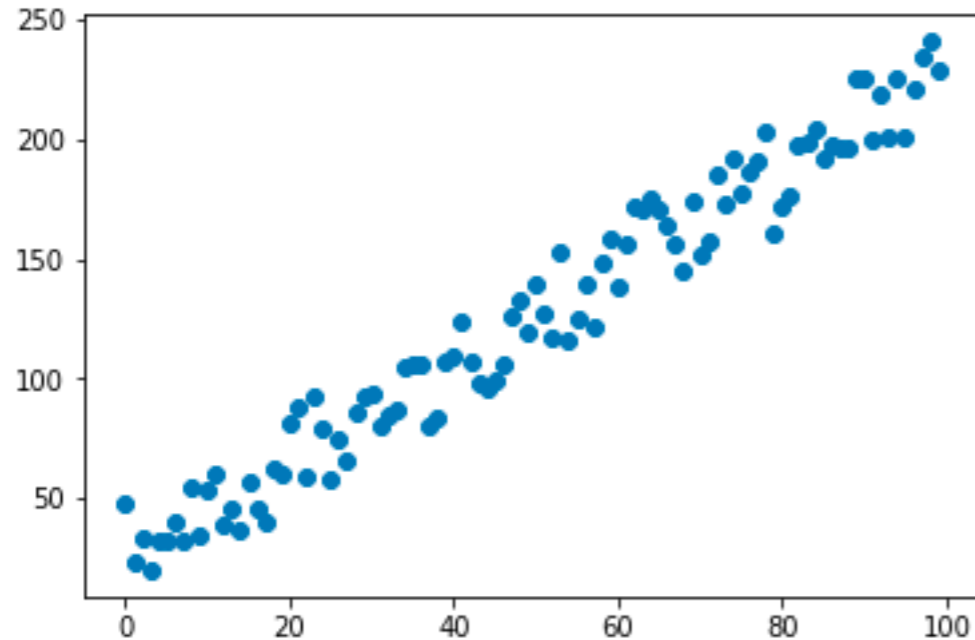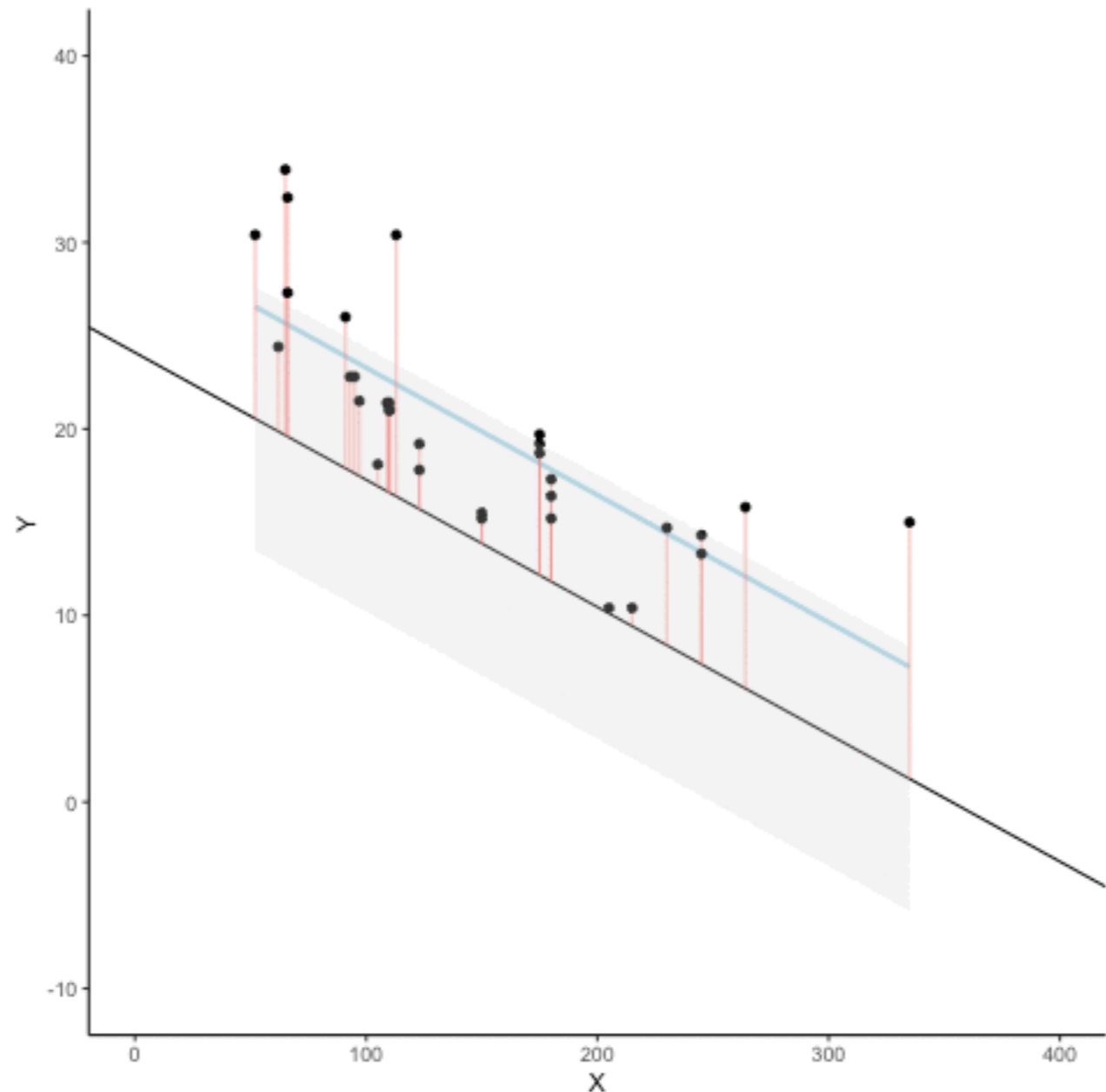
# Simple Linear Regression

## Line of best fit between two variables

\*    The **blue** line represents the optimal line of best fit.

\*\*   The **black** line represents the current line of fit.

\*\*\* The **red** lines represent distance from points. The goal is to *minimize* these values.

# Simple Linear Regression
## Mathematical formulation

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

$$\underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{\substack{n\times 1 \\ \text{Responses}}} = \underbrace{\begin{pmatrix} 1 & x_{1,1} \\ \vdots & \vdots \\ 1 & x_{n,1} \end{pmatrix}}_{\substack{n\times 2 \\ \text{Design Matrix}}} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}}_{\substack{2\times 1 \\ \text{Parameters}}} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\substack{n\times 1 \\ \text{Error}}}$$

$$Y_{n\times 1} = X_{n\times 2}\beta_{2\times 1} + \varepsilon_{n\times 1}$$

---

$n$ is the number of observations, there are 2 variables, **X** provides the design matrix for the variables, **y** is response vector, $\beta$ is the parameter or coefficient vector and $\varepsilon$ is the random error vector

# SLR
## Components



**Real** $y_1$

**Residual** $e_1 = (y_1 - \hat{y}_1)$

**Predicted** $\hat{y}_1$

Best Line of Fit: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

$y_n$

$e_n = (y_n - \hat{y}_n)$

$\hat{y}_n$

A line of fit: $\hat{y}_i^* = \hat{\beta}_0^* + \hat{\beta}_1^* x_i$

**Definition:**
A *closed-form* or *analytical expression* is a formula that provides an answer to a mathematical statement involving infinity that is finite.

The **Quadratic Formula** given by
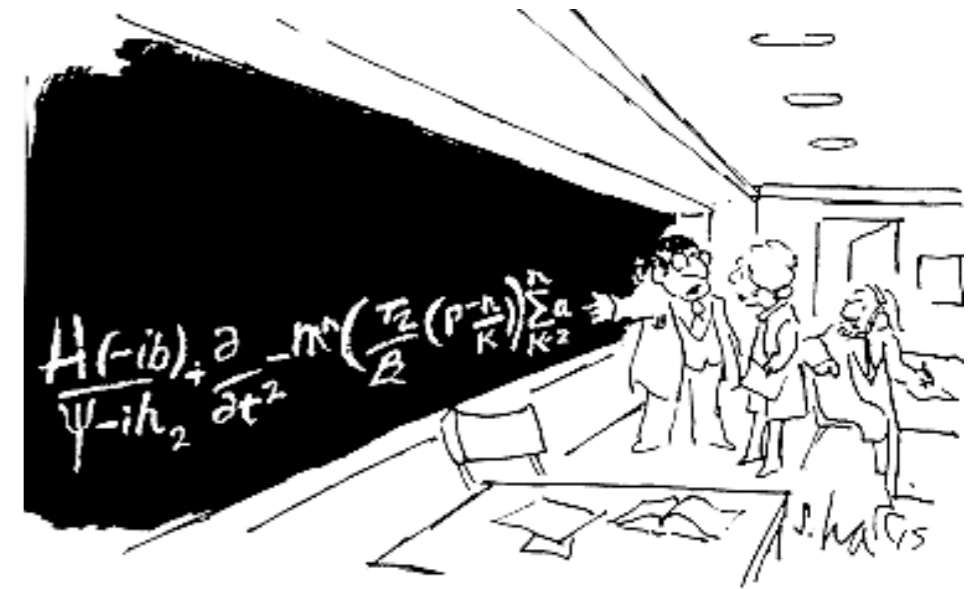$$0 = ax^2 + bx + c$$

has a solution of
$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$



"But this *is* the simplified version for the general public."

**Sidney Harris**

# Estimating Parameters

**Minimize the <u>R</u>esidual <u>S</u>um <u>S</u>quared** (red lines)

$$\hat{\beta} = \underset{\beta_0, \beta_1 \in \mathbb{R}}{\arg\min} \sum_{i=1}^{n} \left( y_i - \underbrace{(\beta_0 + \beta_1 x_i)}_{\hat{y}_i} \right)^2$$

**Analytical Solutions**

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

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

$$\overline{x} = \frac{1}{n}\sum_{i=1}^{n} x_i \qquad \overline{y} = \frac{1}{n}\sum_{i=1}^{n} y_i$$

# Through the Mean

The closed-form solution for the intercept is given as:

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

As a result, we note that **every linear regression line** must go through means of x and y

$$\overline{x} = \frac{1}{n}\sum_{i=1}^{n}x_i \qquad \overline{y} = \frac{1}{n}\sum_{i=1}^{n}y_i$$

# Normal **Equation**

$$\hat{\boldsymbol{\beta}} = \left(\boldsymbol{X}^\top \boldsymbol{X}\right)^{-1} \boldsymbol{X}^\top \boldsymbol{y}$$

- For more than two parameters, the solution to linear regression is given by the normal equation.

- This is a closed-form solution that works well as it is a well-studied problem.

# Extra

# Notations

**Training Set** $\qquad\qquad \mathcal{D} = \left\{ \langle \boldsymbol{x}^{(i)}, y^{(i)} \rangle, i = 1, \dots, n \right\}$

**Unknown Function** $\qquad f(\boldsymbol{x}) = y$

**Hypothesis** $\qquad\qquad h(\boldsymbol{x}) = \hat{y}$

$h : \mathbb{R}^m \rightarrow \mathcal{Y}, \mathcal{Y} = \{1, \dots, k\}$ $\qquad\qquad h : \mathbb{R}^m \rightarrow \mathbb{R}$

**Classification** $\qquad\qquad\qquad\qquad\qquad$ **Regression**

# Regression Hypothesis
## Simple Linear Regression

$$h(\boldsymbol{x}) = \beta_0 + \beta_1 x_1$$
$$= \beta_0 x_0 + \beta_1 x_1, \text{ where } x_0 = 1$$
$$= \sum_{i=0}^{p} \beta_i x_i$$
$$= \boldsymbol{\beta}^T \boldsymbol{x}$$

**Training**

$$\mathcal{D} = \left\{ \langle \boldsymbol{x}^{(i)}, y^{(i)} \rangle, i = 1, \ldots, n \right\}$$
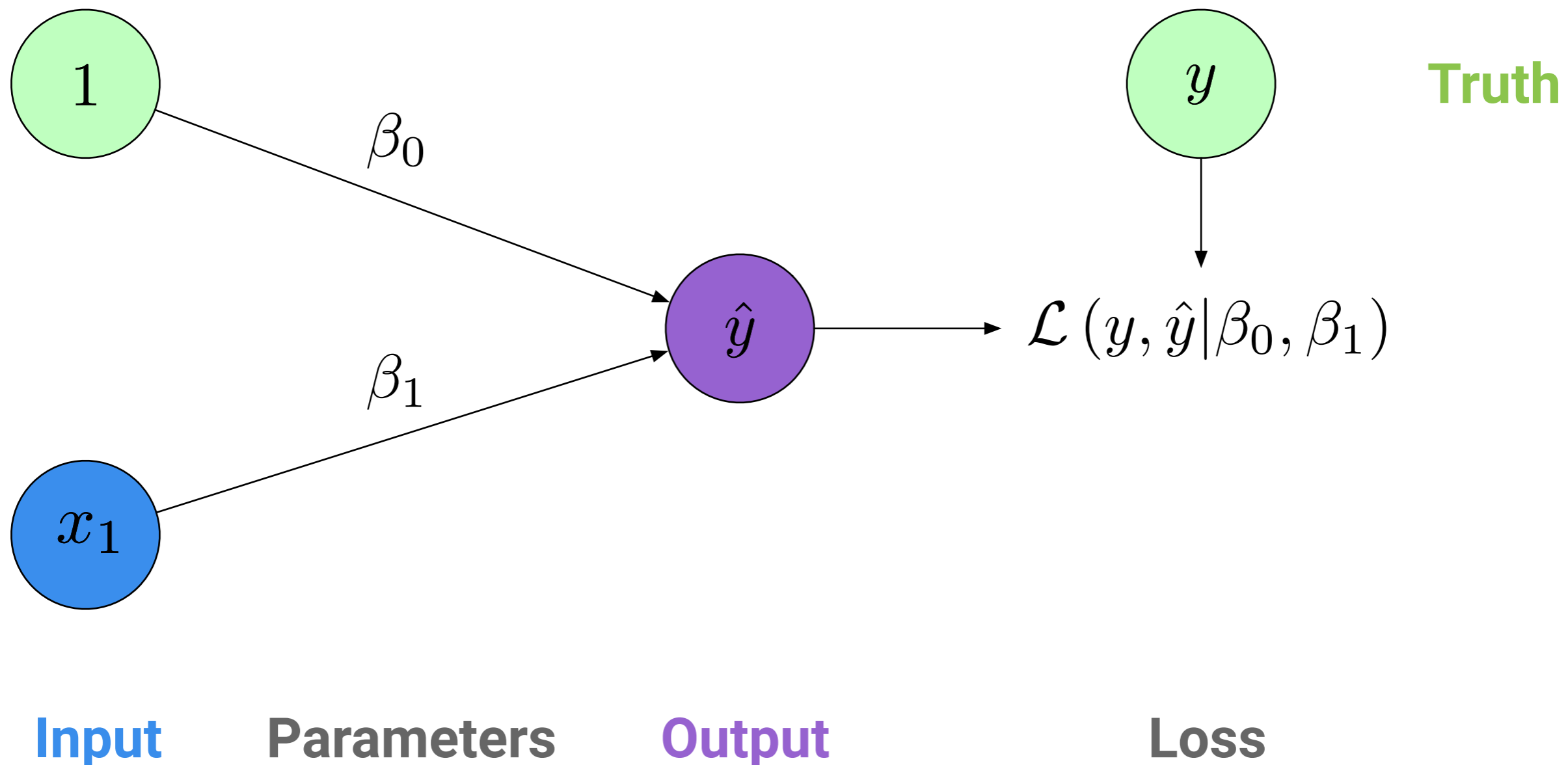
**Unknown Function**

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

# Computational Graph
## SLR shown in the context of a graph

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



**Input**     **Parameters**     **Output**     **Loss**

# Organizing Data

⭐ **Row-major**

$$\mathbf{y}_{1 \times n} = \boldsymbol{\beta}_{1 \times 2} \mathbf{X}_{2 \times n}$$

$$\boldsymbol{\beta}_{1 \times 2} = \begin{bmatrix} \beta_0 & \beta_1 \end{bmatrix}_{1 \times 2}$$

$$\mathbf{X}_{2 \times n} = \begin{bmatrix} 1 & \cdots & 1 \\ \mathbf{x}^{(1)} & \cdots & \mathbf{x}^{(n)} \end{bmatrix}_{2 \times n}$$

**Column-major**

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times 2} \boldsymbol{\beta}_{2 \times 1}$$

$$\boldsymbol{\beta}_{2 \times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}_{2 \times 1}$$

$$\mathbf{X}_{n \times 2} = \begin{bmatrix} 1 & \mathbf{x}^{(1)} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(n)} \end{bmatrix}_{n \times 2}$$

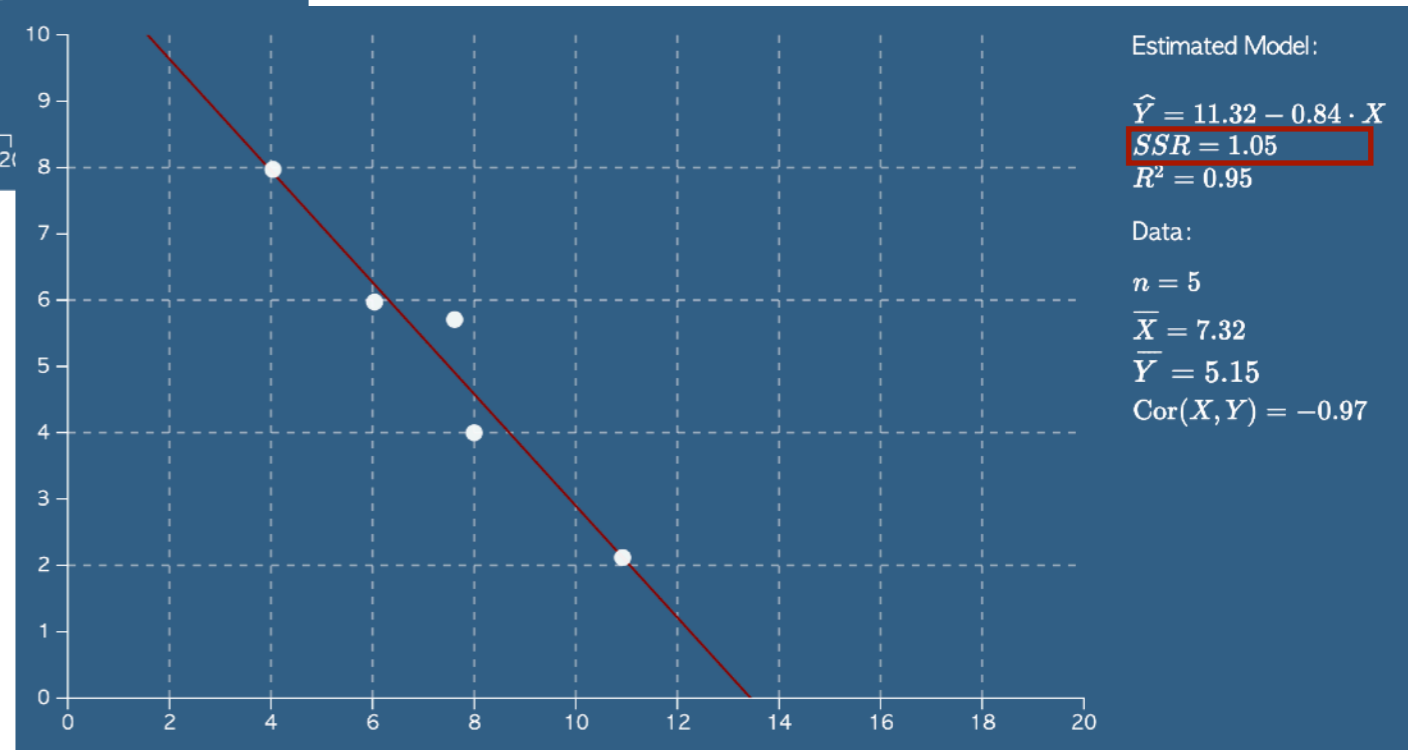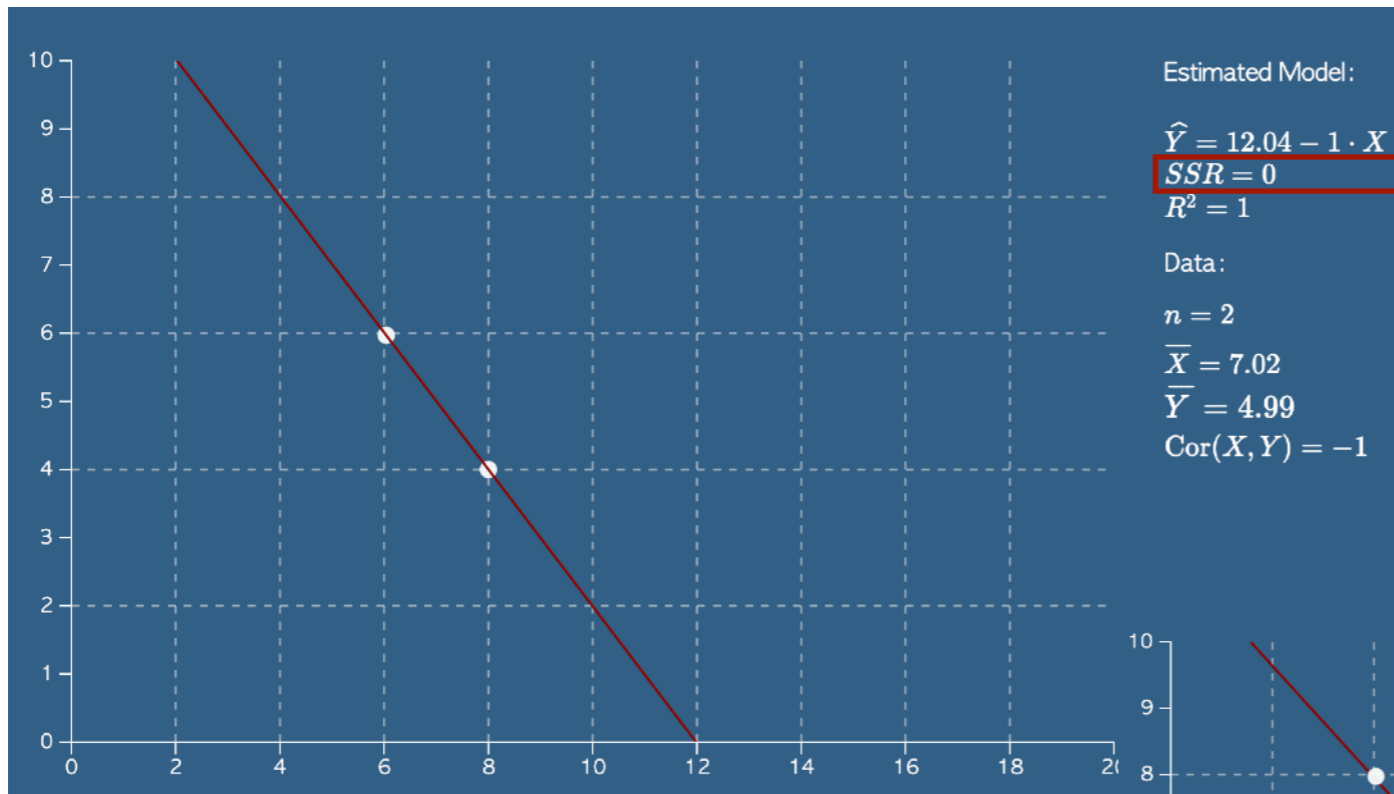**Rows** are **features**
**Columns** are **observations**

**Rows** are **observations**
**Columns** are **features**

NumPy, TensorFlow, PyTorch uses row-major form

# Dynamic SLR
## Observing how data changes affect estimates



Estimated Model:

$\widehat{Y} = 12.04 - 1 \cdot X$
$SSR = 0$
$R^2 = 1$

Data:

$n = 2$
$\overline{X} = 7.02$
$\overline{Y} = 4.99$
$\text{Cor}(X, Y) = -1$

Estimated Model:

$\widehat{Y} = 11.32 - 0.84 \cdot X$
$SSR = 1.05$
$R^2 = 0.95$

Data:

$n = 5$
$\overline{X} = 7.32$
$\overline{Y} = 5.15$
$\text{Cor}(X, Y) = -0.97$

Interactive Web Widget

# Jargon/Terms across Fields

| Machine + Deep Learning | Statistical Learning | Description |
|---|---|---|
| Network, Graphs | Model | Assumption on how things work |
| Label, Target, Y | Dependent, Response, Output | Value being predicted |
| Feature, X | Independent, Explanatory, Input | Used to make predictions |
| Feature Engineering | Data Wrangling/Transformation | Changing the data to get more predictions |
| Learning | Fitting | Creating the model |
| Generalization | Test Set Performance | How the model works with new data |
| Weights | Parameters | Learned values |
| 1D, 2D, 3D, …, nD | Dimensionality | Number of variables |
| Supervised Learning | Regression/Classification | Learning with a target |
| Unsupervised Learning | Density Estimation/Clustering | Learning without a target |
| Large grant: $1 million | Large grant: $50k | Statisticians are bad marketers |
| Conferences: French Alps | Conferences: Las Vegas in August | ML community goes on vacations |

An extension of Rob Tibshirani's Glossary

# Summary

- Differentiated between supervised and unsupervised learning.

- Discussed the closed-form solution to regression.

This work is licensed under the