



SINGAPORE UNIVERSITY OF
TECHNOLOGY AND DESIGN

10.020 Data Driven World

Linear Regression: Training

Peng Song, ISTD

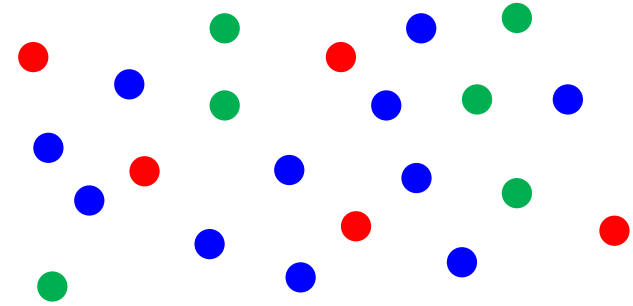
Week 9, Lesson 1, 2021

some slides are from Jia-Bin Huang

Revision: Supervised Learning

- **Classification**

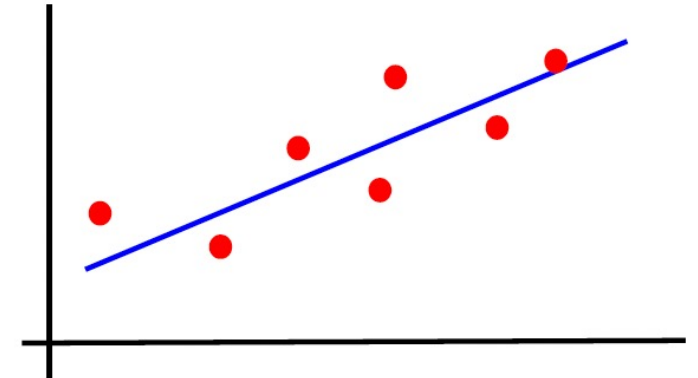
- Given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function $f(x)$ to predict y given x
- y is **categorical**



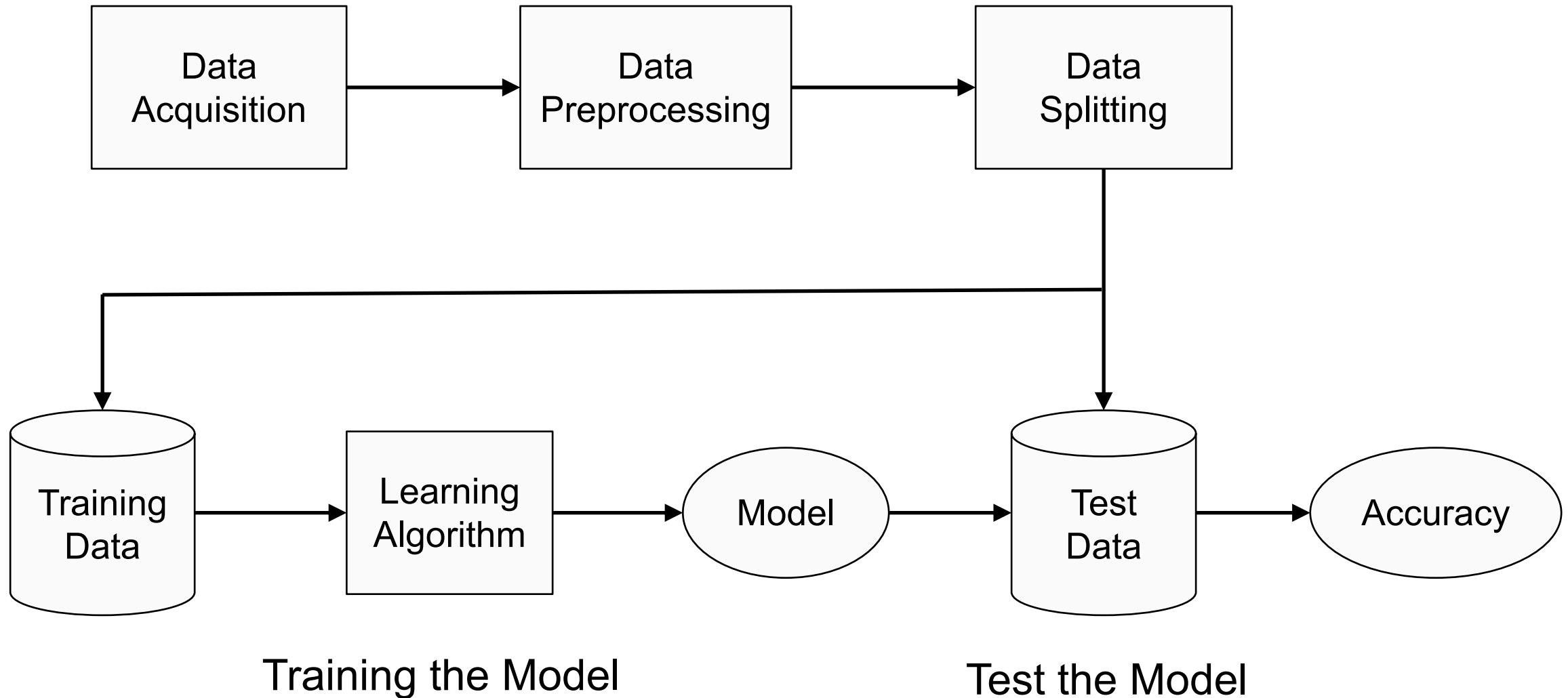
- **Regression**

- Given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function $f(x)$ to predict y given x
- y is **numeric**

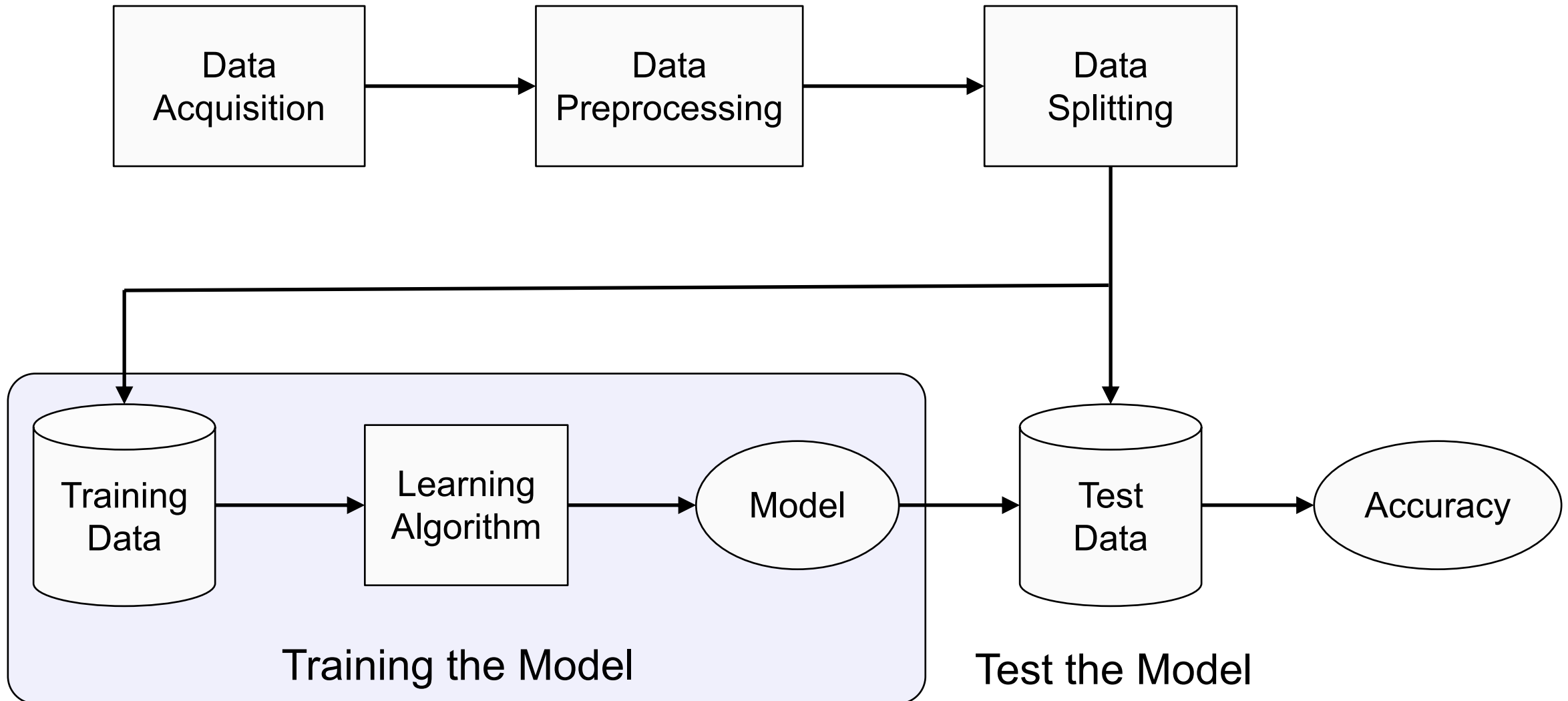
We will focus on Regression this week



Revision: Supervised Learning Process



Today: Linear Regression



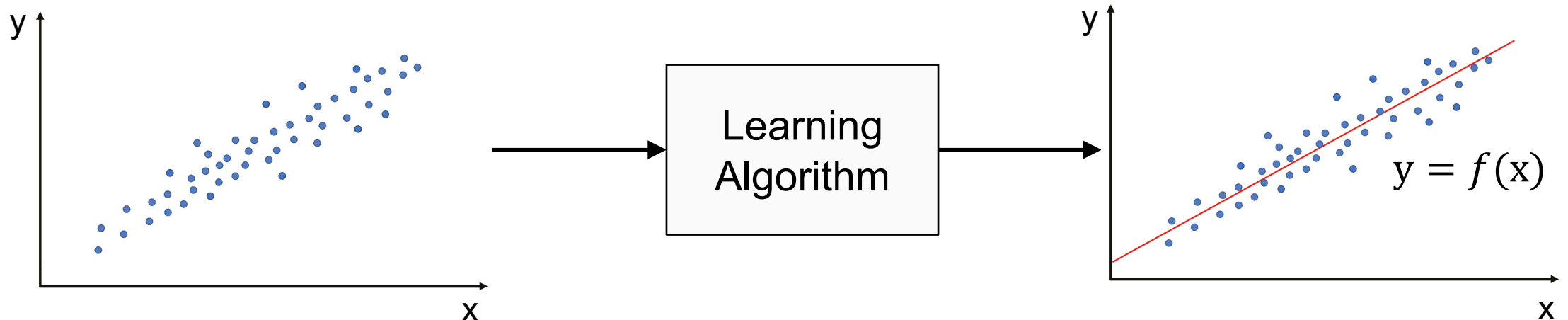
Regression Problem: House Pricing Prediction

- Given samples (x, y) , where
 - x is house size
 - y is house price
- Learn a function $y = f(x)$
 - Assume f is a linear function



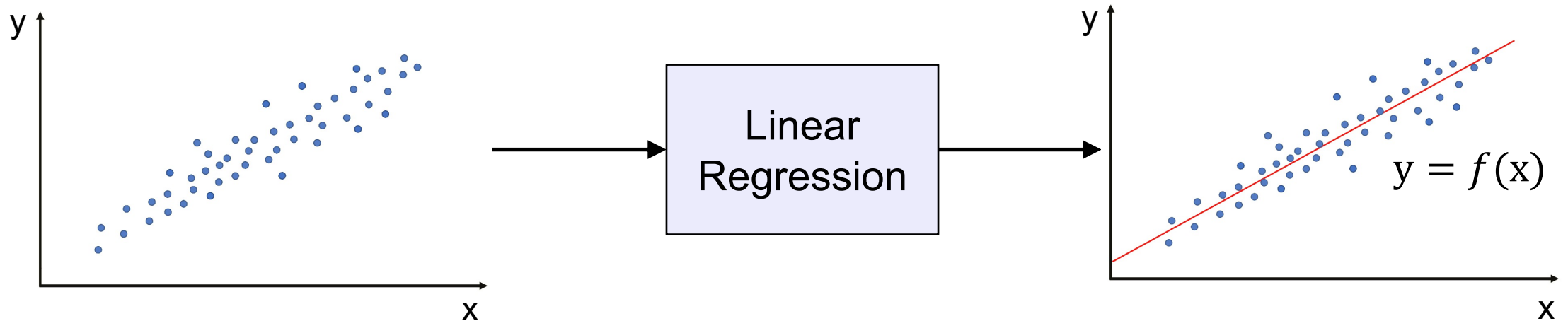
Regression Problem: House Pricing Prediction

- Given samples (x, y) , where
 - x is house size
 - y is house price
- Learn a function $y = f(x)$
 - Assume f is a linear function



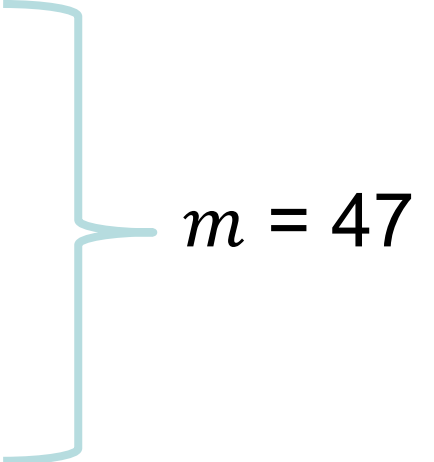
Linear Regression

Linear regression is a machine learning algorithm that assumes a **linear** relationship between the input variable x and the single output variable y .



Training Set

Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...



- m = Number of training examples
- x = Input variable / feature
- y = Output variable / target variable
- $(x^i, y^i) = i^{th}$ training example

Examples:

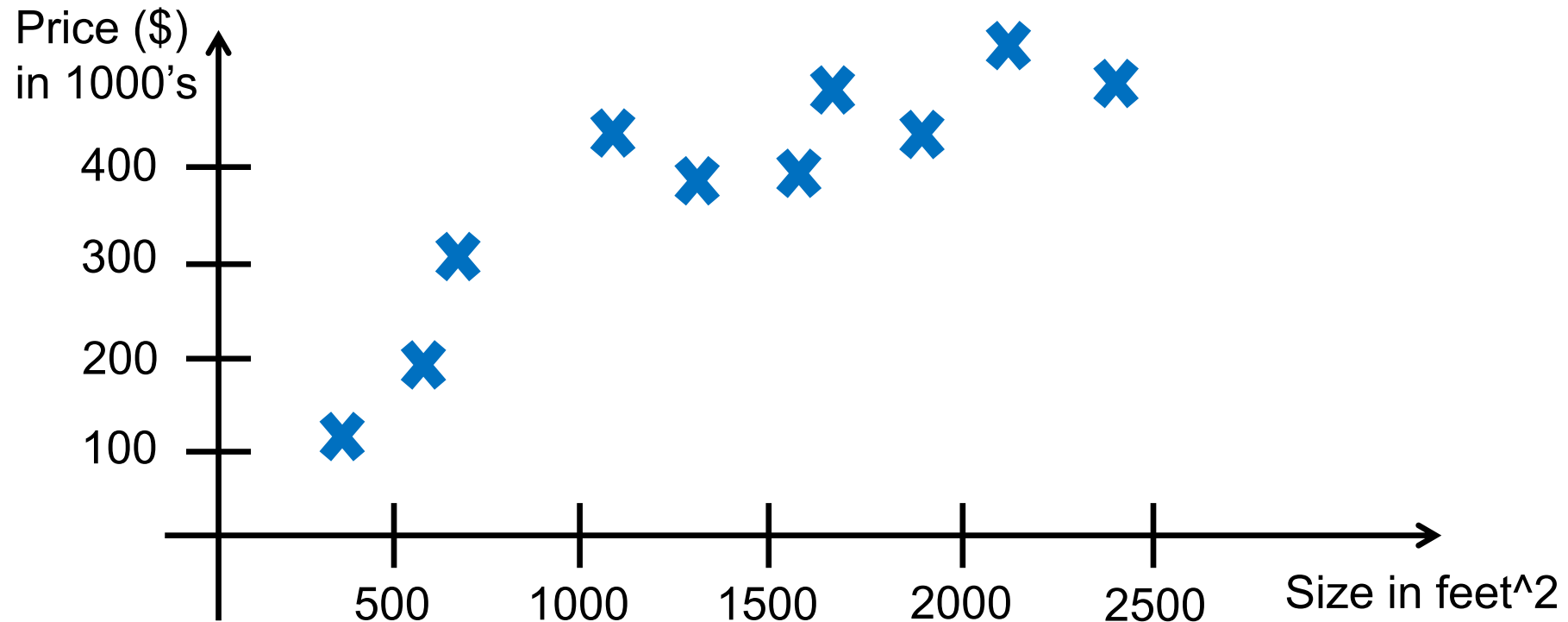
$$x^1 = 2104$$

$$x^2 = 1416$$

$$y^1 = 460$$

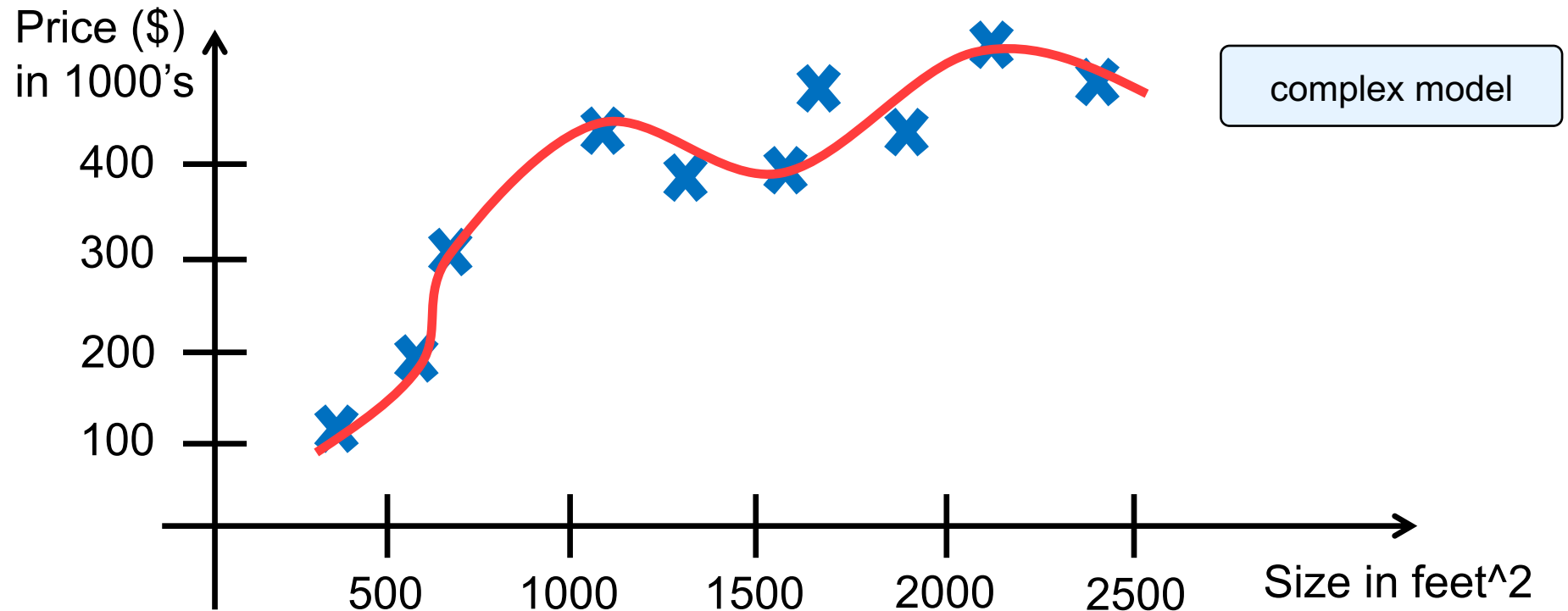
Training Set

Visualize the training set as a scatter plot.



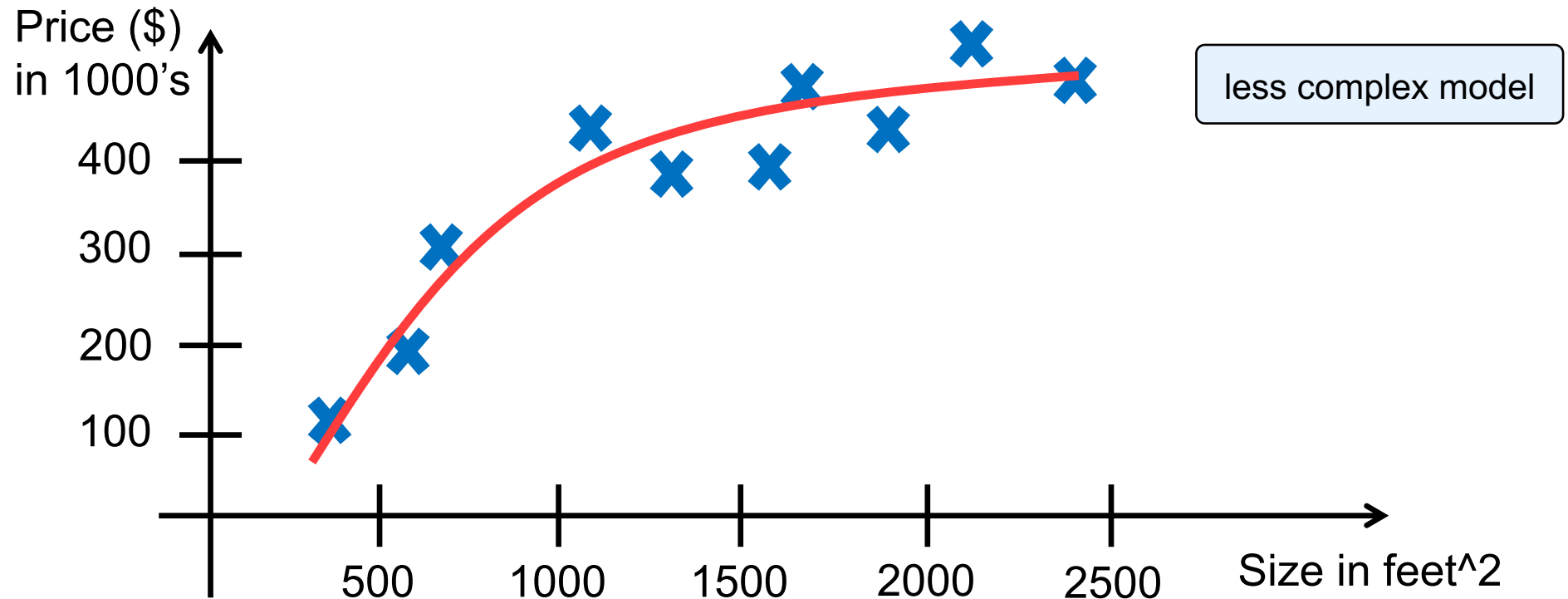
Hypothesis

Hypothesis is a **candidate model** that approximates a target function for mapping examples of inputs to outputs.



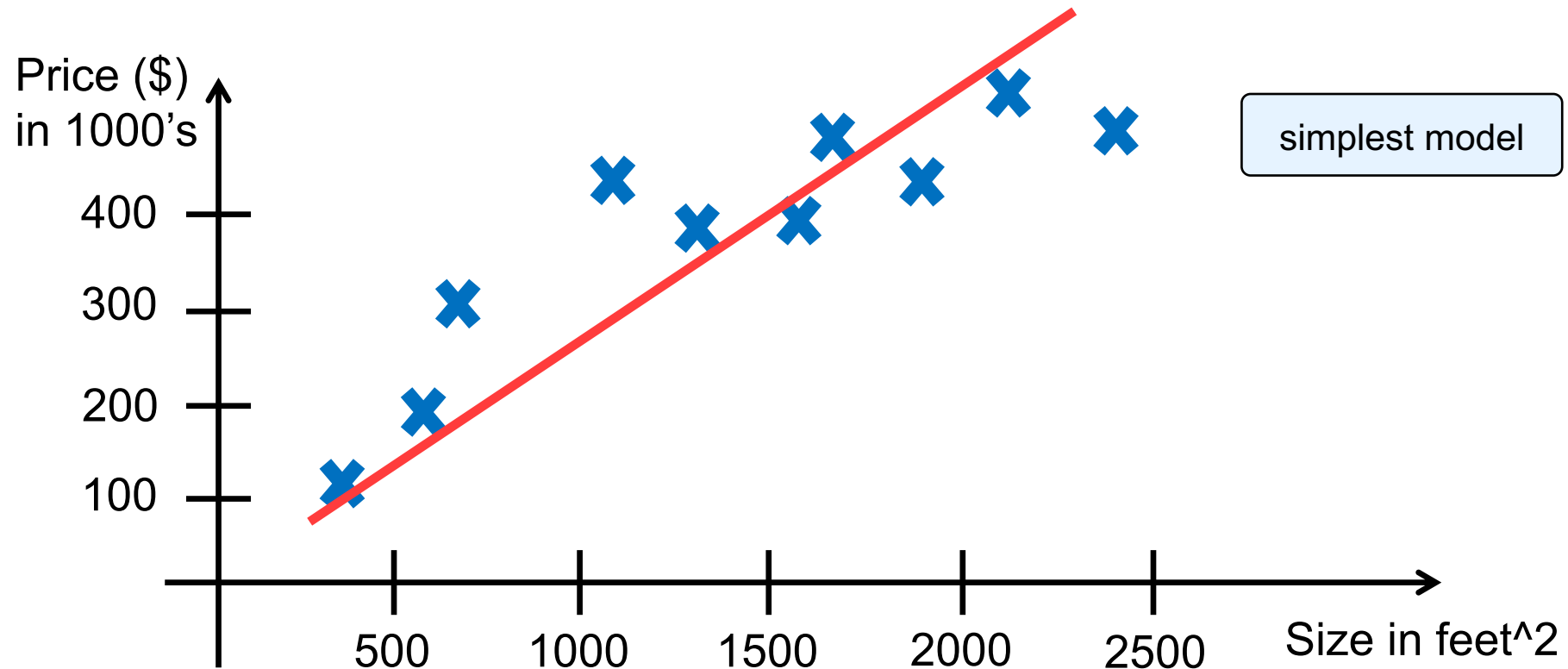
Hypothesis

Hypothesis is a **candidate model** that approximates a target function for mapping examples of inputs to outputs.



Hypothesis

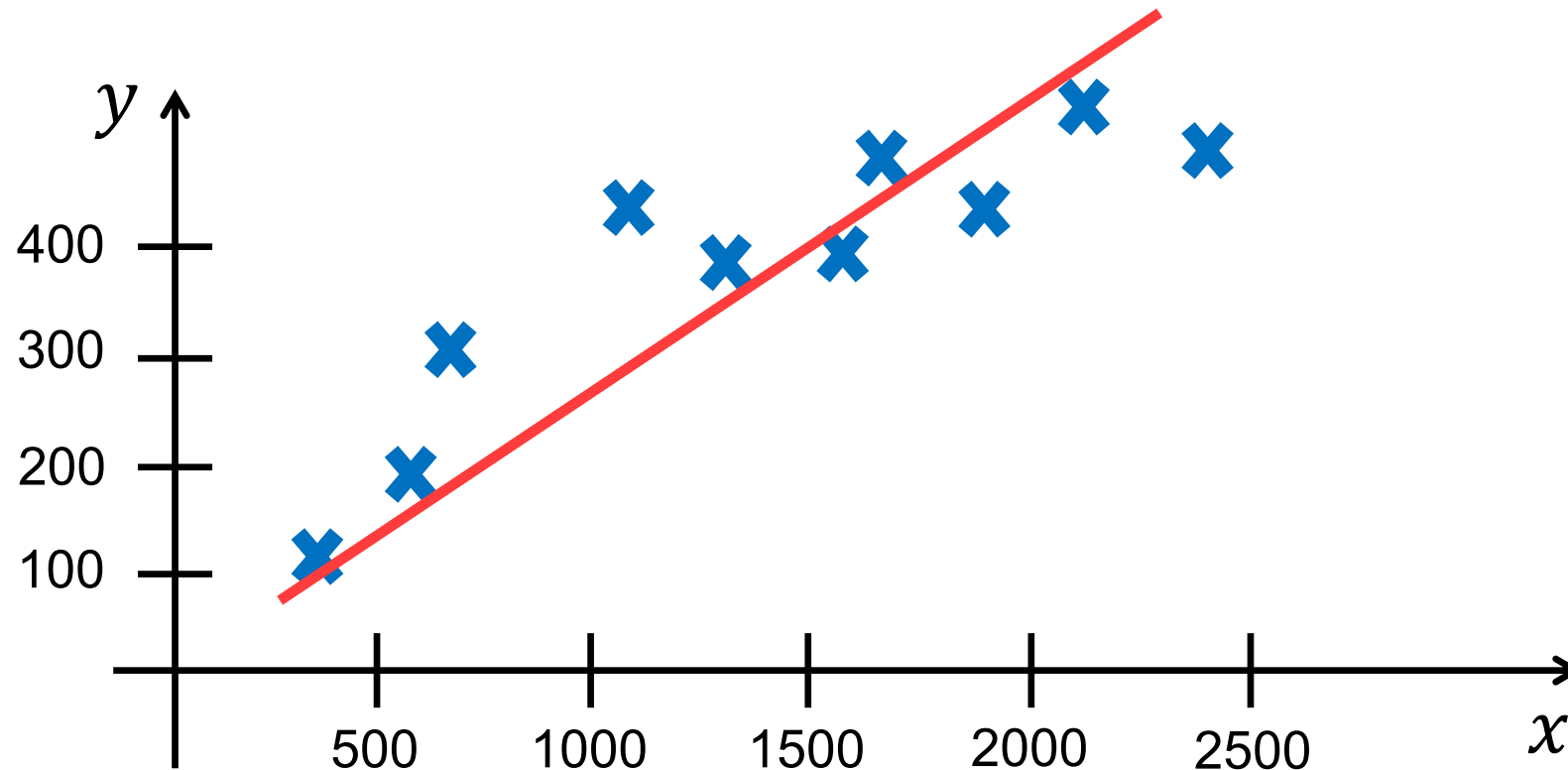
Hypothesis is a **candidate model** that approximates a target function for mapping examples of inputs to outputs.



Hypothesis in Linear Regression

$$y = h_{\beta}(x) = \beta_0 + \beta_1 x$$

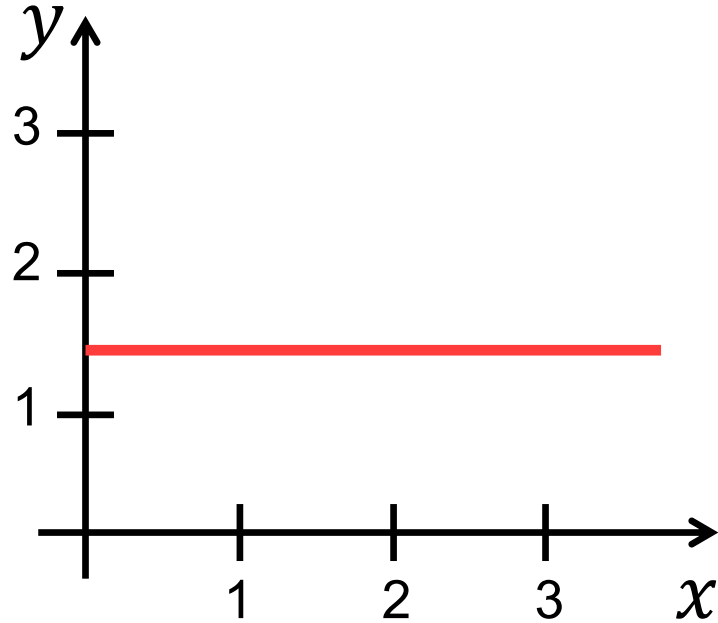
β_0 and β_1 are unknown parameters of the hypothesis



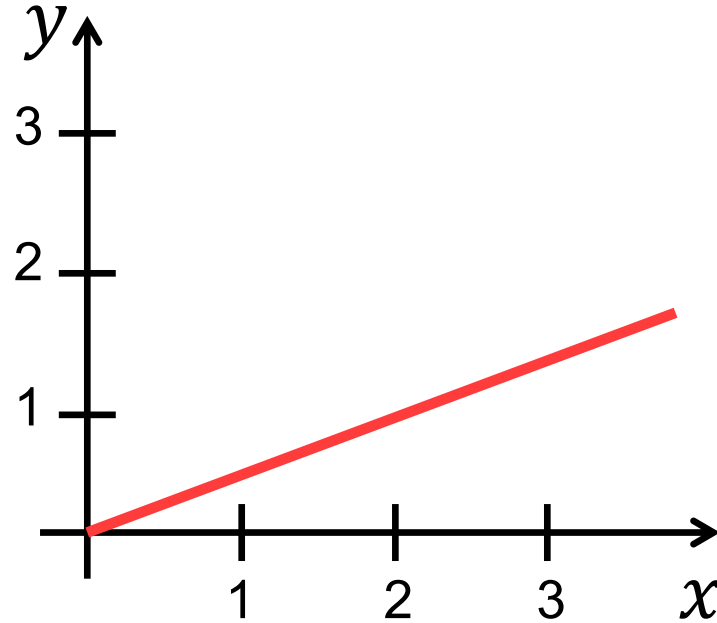
Hypothesis in Linear Regression

$$y = h_{\beta}(x) = \beta_0 + \beta_1 x$$

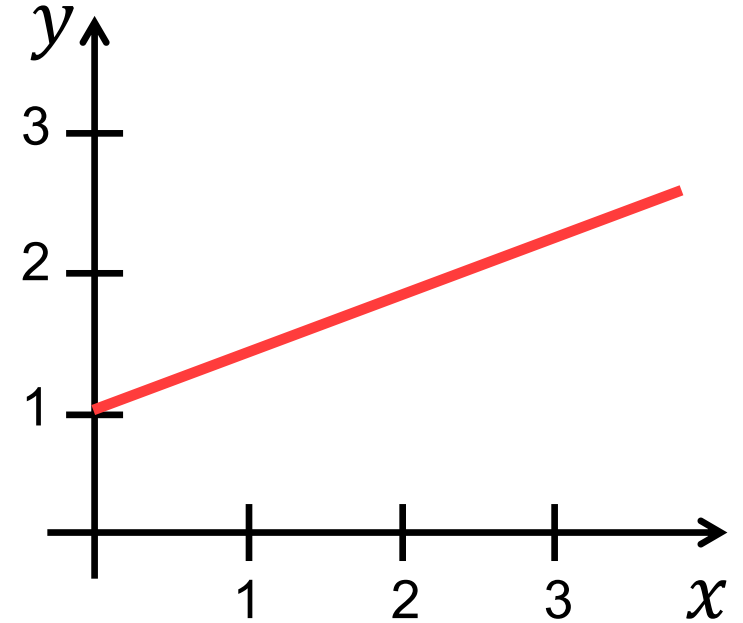
How to choose β_0 and β_1 ?



$$\begin{aligned}\beta_0 &= 1.5 \\ \beta_1 &= 0\end{aligned}$$



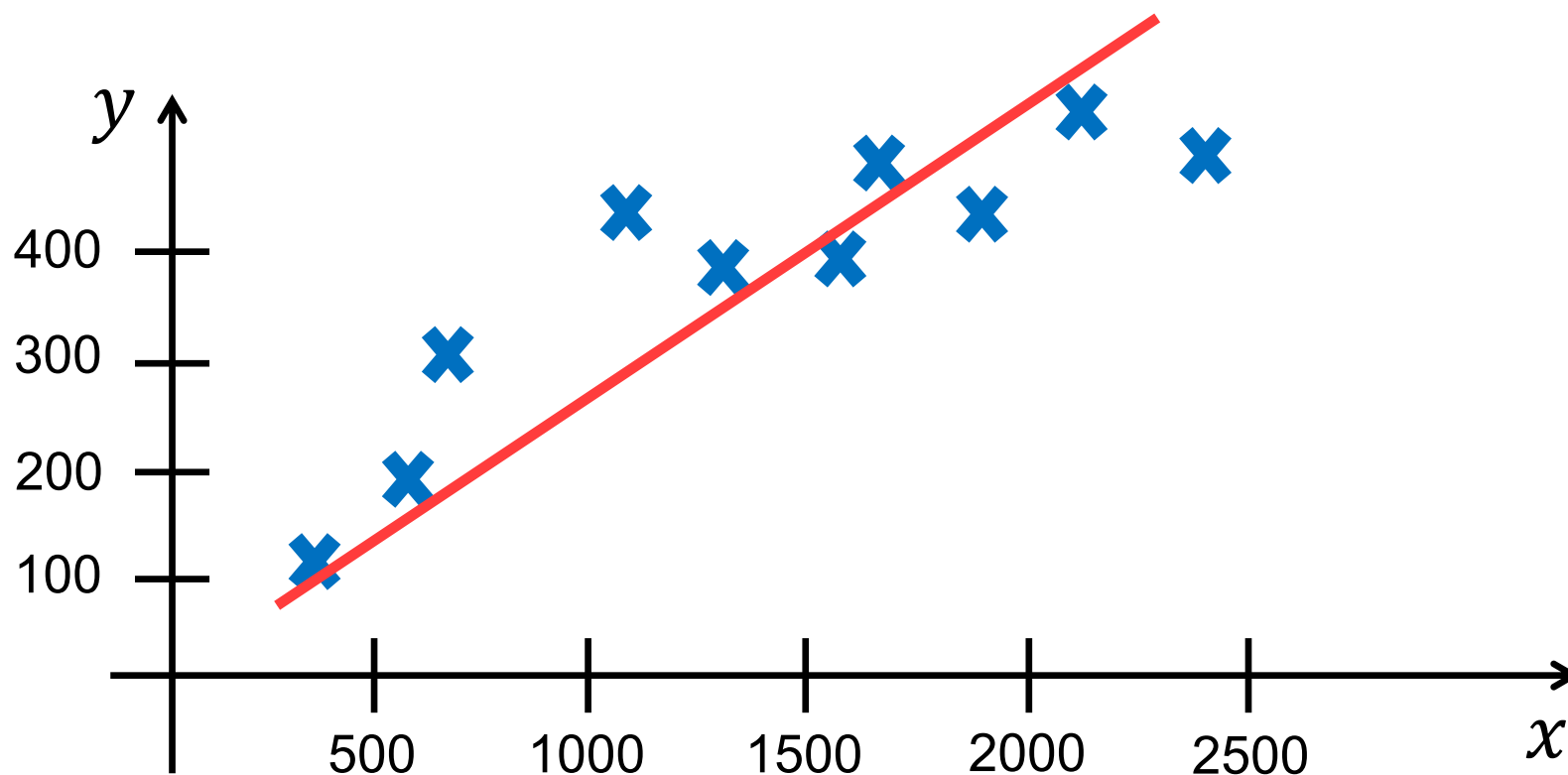
$$\begin{aligned}\beta_0 &= 0 \\ \beta_1 &= 0.5\end{aligned}$$



$$\begin{aligned}\beta_0 &= 1 \\ \beta_1 &= 0.5\end{aligned}$$

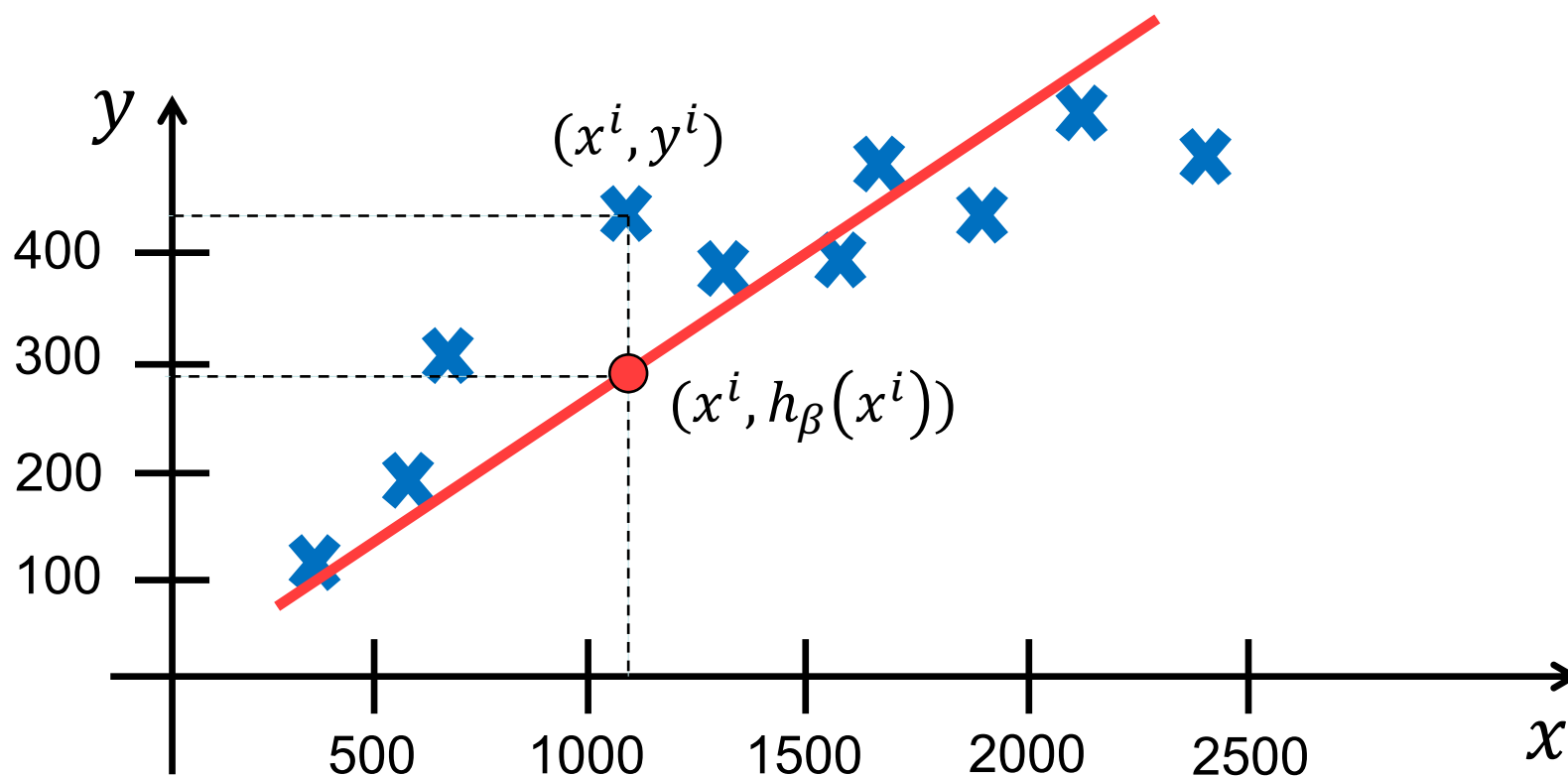
Cost function

Idea: Choose β_0, β_1 so that our predicted value $h_{\beta}(x^i)$ is close to the observed value y^i for our training examples $\{(x^i, y^i)\}$



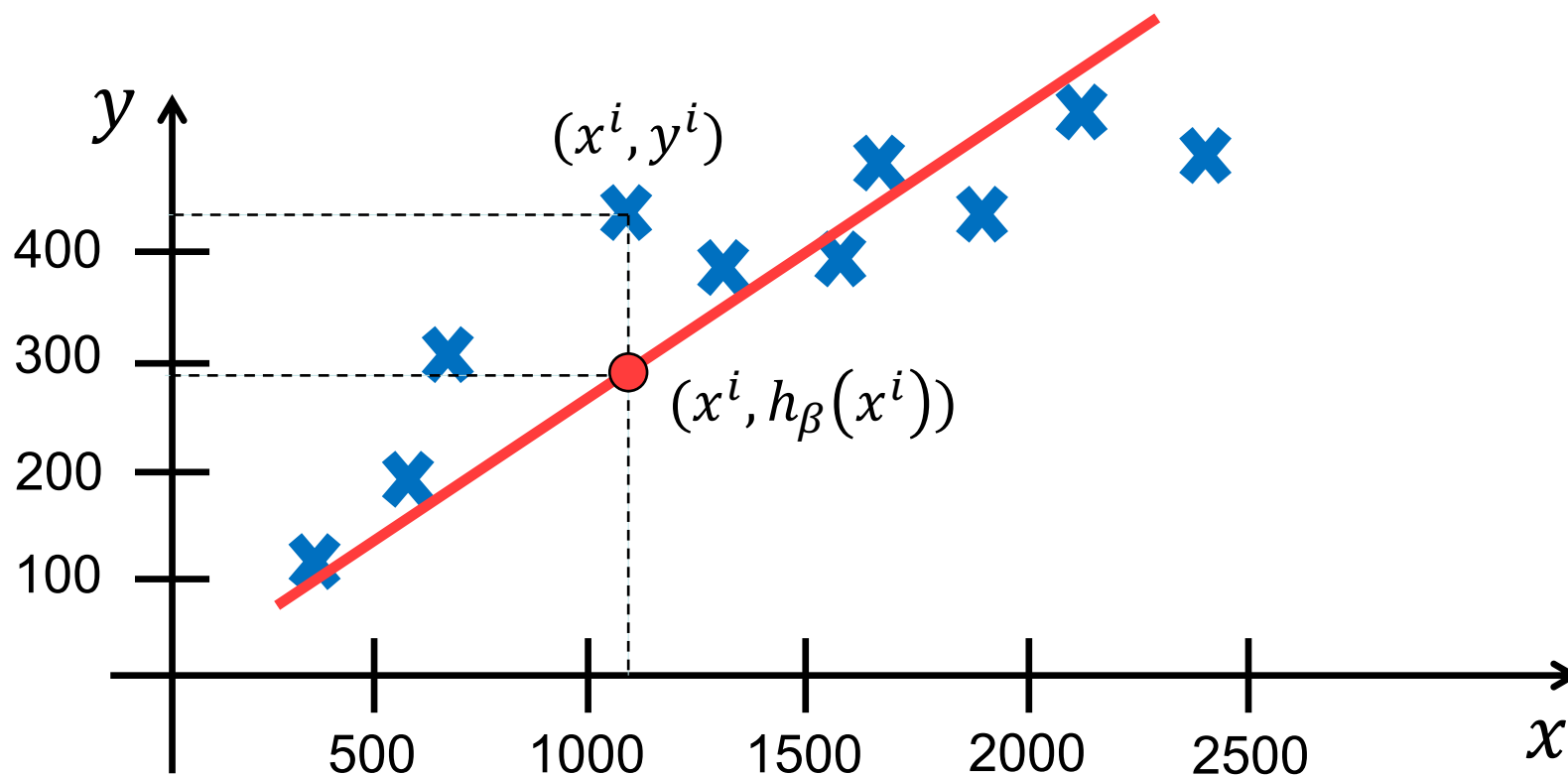
Cost function

Idea: Choose β_0, β_1 so that our predicted value $h_\beta(x^i)$ is close to the observed value y^i for our training examples $\{(x^i, y^i)\}$



Cost function

$$\underset{\beta_0, \beta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2 \quad \text{where } h_{\beta}(x^i) = \beta_0 + \beta_1 x^i$$



Cost function

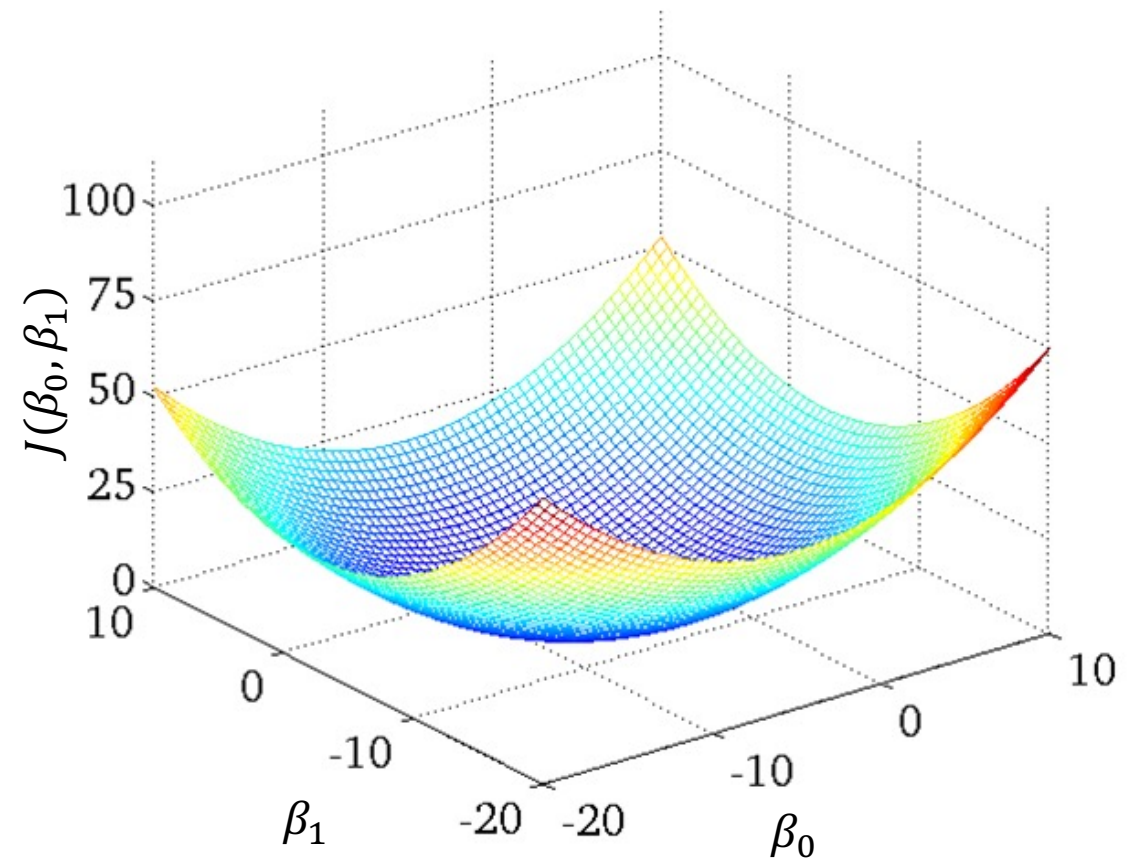
$$\underset{\beta_0, \beta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$$

$$\text{where } h_{\beta}(x^i) = \beta_0 + \beta_1 x^i$$

$$J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$$

Cost function

$$\underset{\beta_0, \beta_1}{\text{minimize}} \quad J(\beta_0, \beta_1)$$



Cost Function

- **Hypothesis:** $h_{\beta}(x) = \beta_0 + \beta_1 x$
- **Parameters:** β_0, β_1
- **Cost function:** $J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$
- **Goal:** minimize $J(\beta_0, \beta_1)$
 β_0, β_1

Simplified Hypothesis

- **Hypothesis:**

$$h_{\beta}(x) = \beta_0 + \beta_1 x \longrightarrow$$

- **Hypothesis:**

$$h_{\beta}(x) = \beta_1 x \quad \beta_0 = 0$$

- **Parameters:**

$$\beta_0, \beta_1 \longrightarrow$$

- **Parameters:**

$$\beta_1$$

- **Cost function:**

$$J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2 \longrightarrow$$

- **Cost function:**

$$J(\beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$$

- **Goal:**

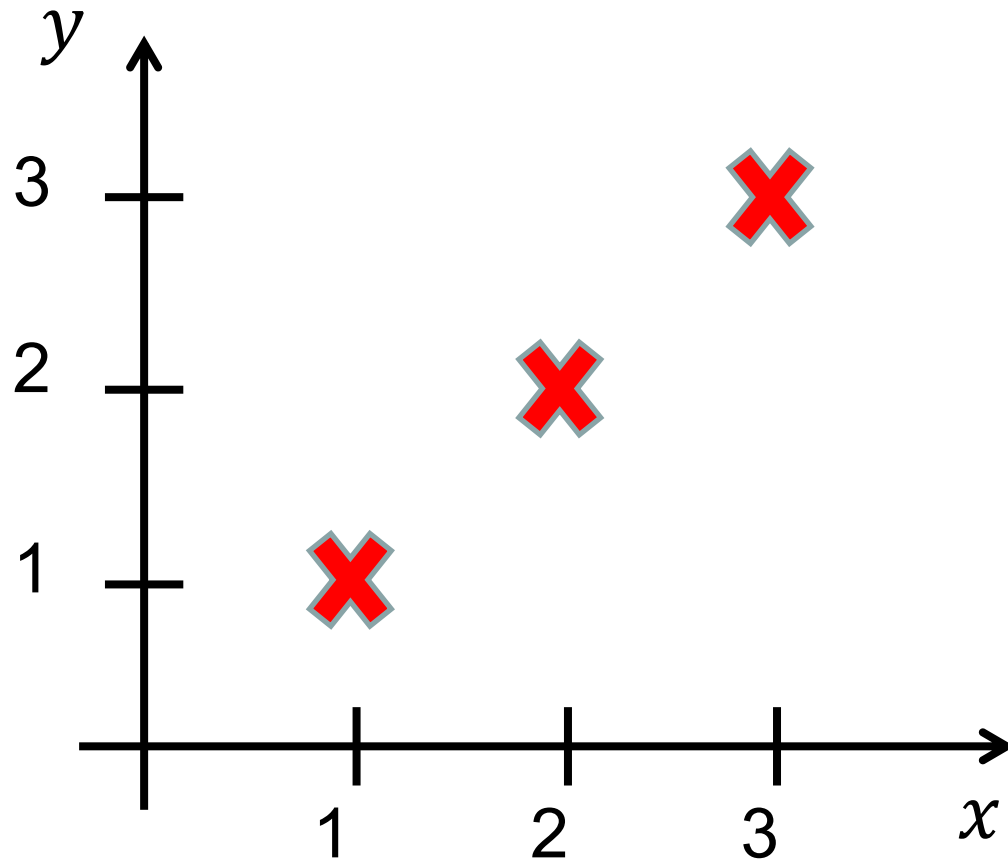
$$\underset{\beta_0, \beta_1}{\text{minimize}} \quad J(\beta_0, \beta_1) \longrightarrow$$

- **Goal:**

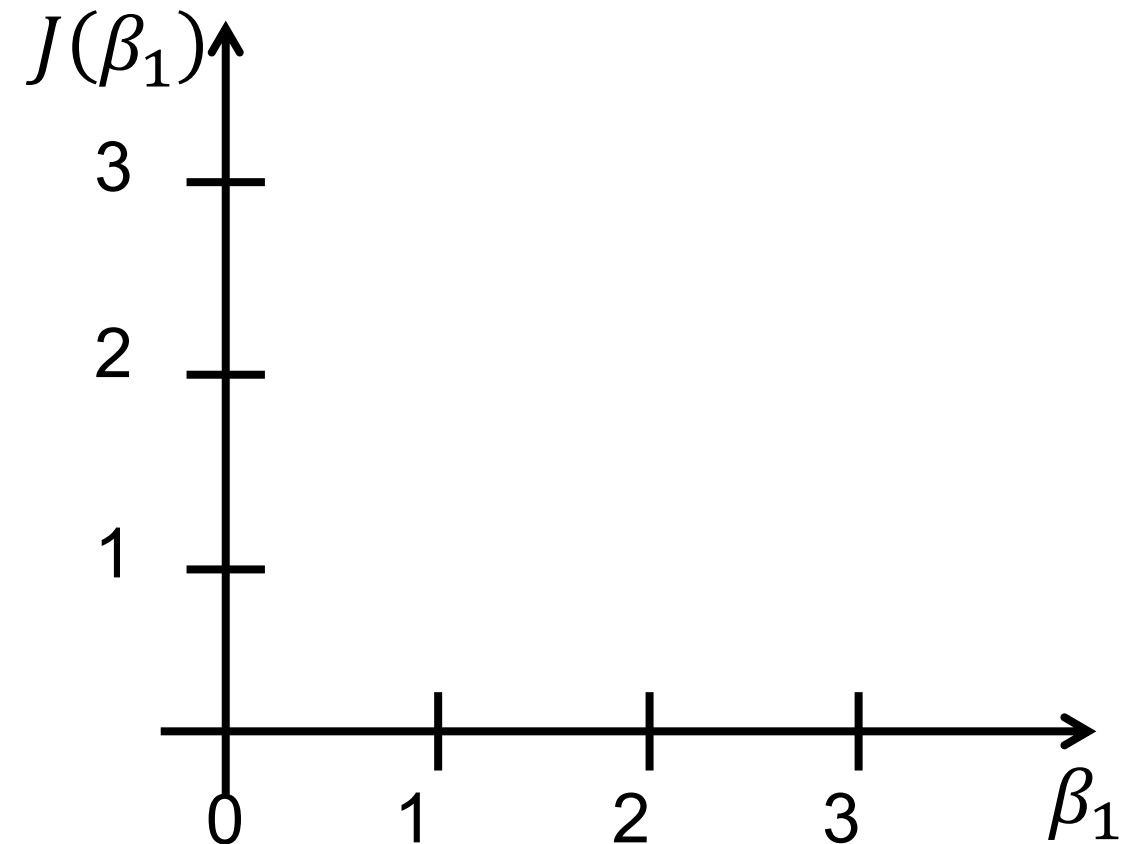
$$\underset{\beta_1}{\text{minimize}} \quad J(\beta_1)$$

Simplified Hypothesis

$h_{\beta}(x)$, function of x

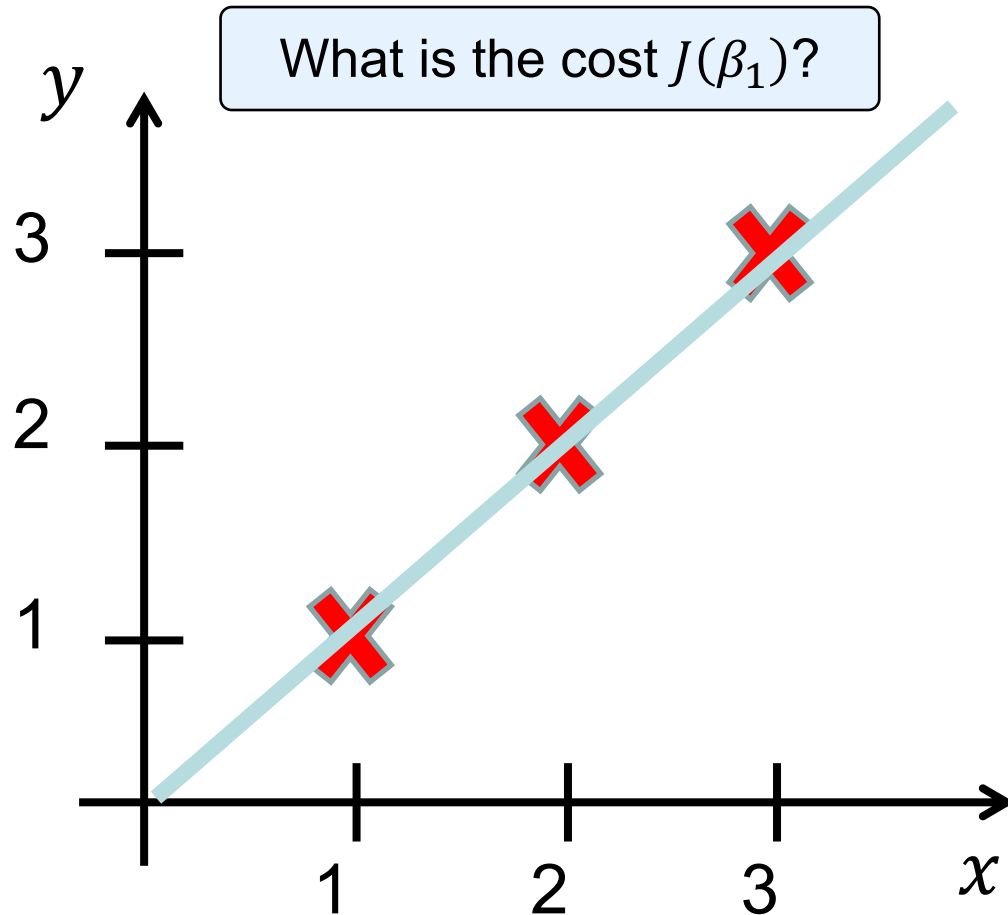


$J(\beta_1)$, function of β_1

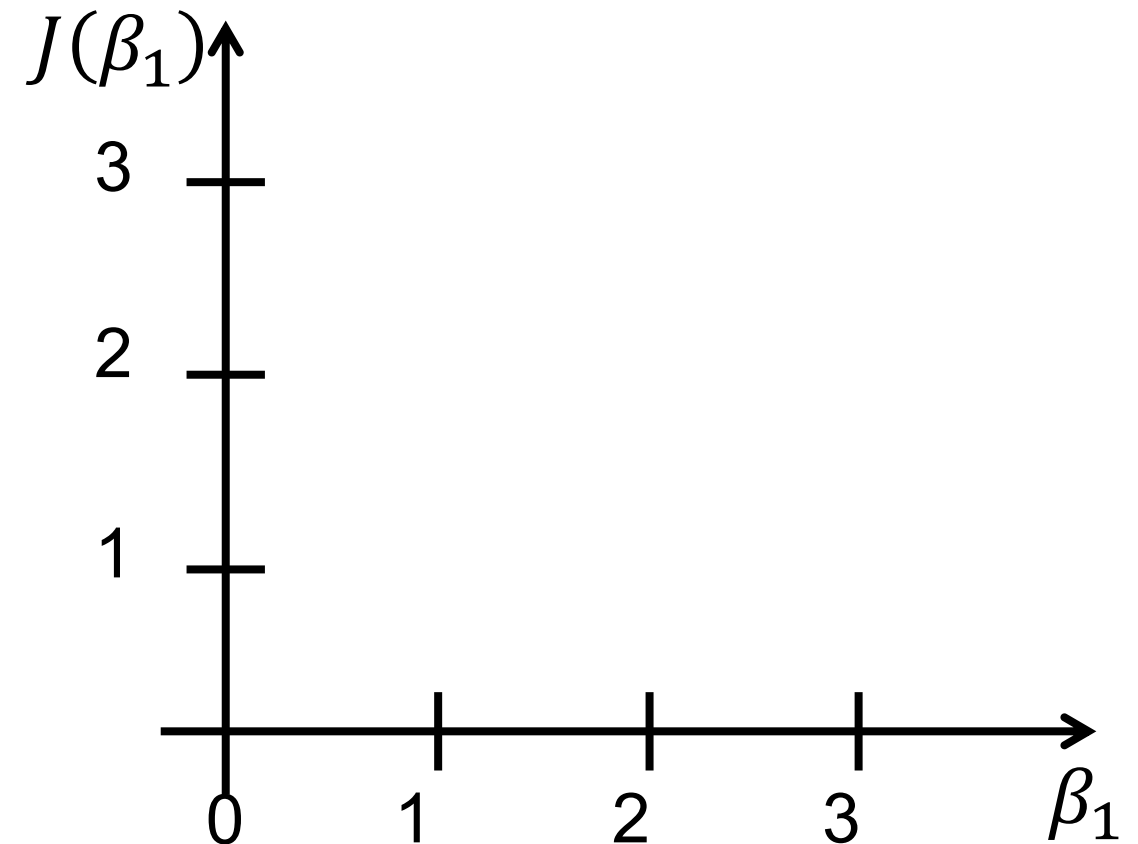


Simplified Hypothesis

$h_{\beta}(x)$, function of x

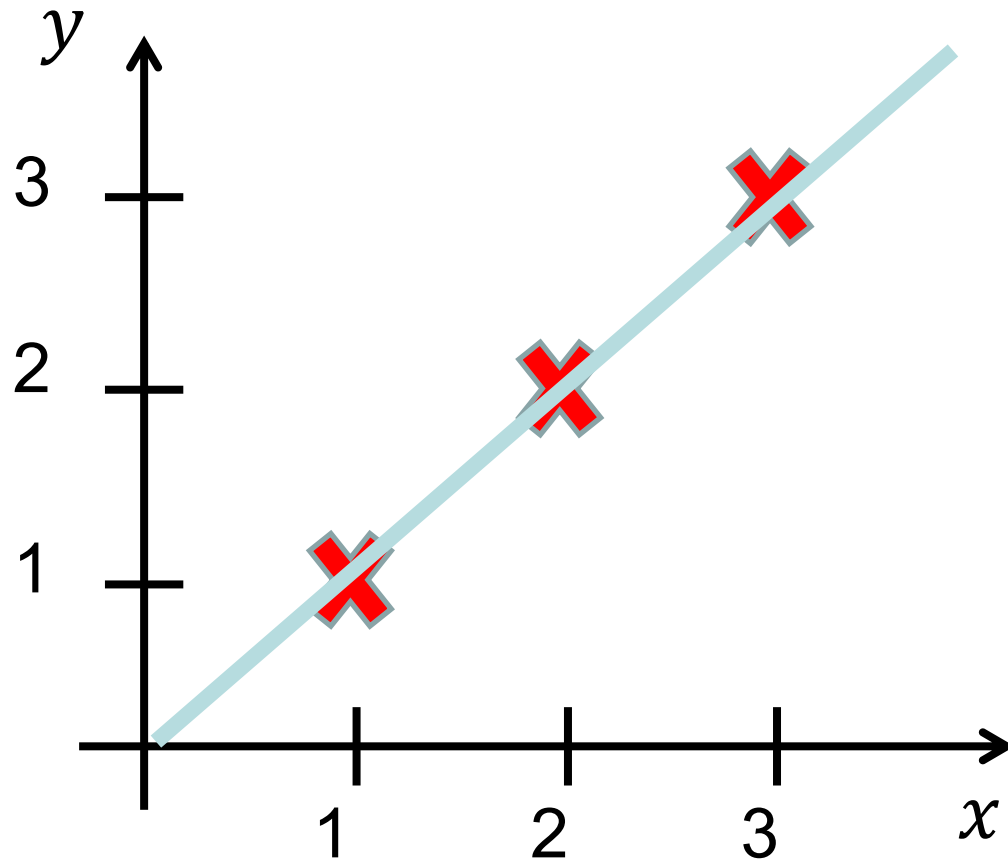


$J(\beta_1)$, function of β_1

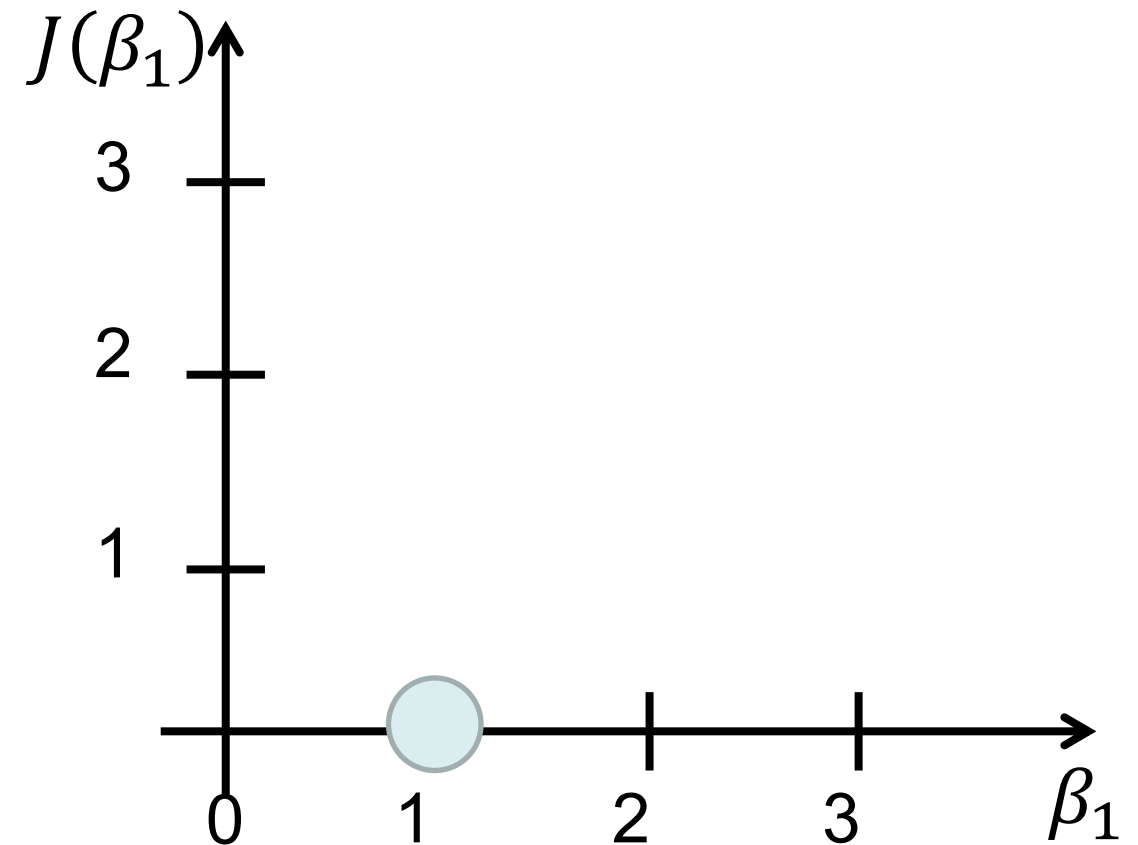


Simplified Hypothesis

$h_{\beta}(x)$, function of x

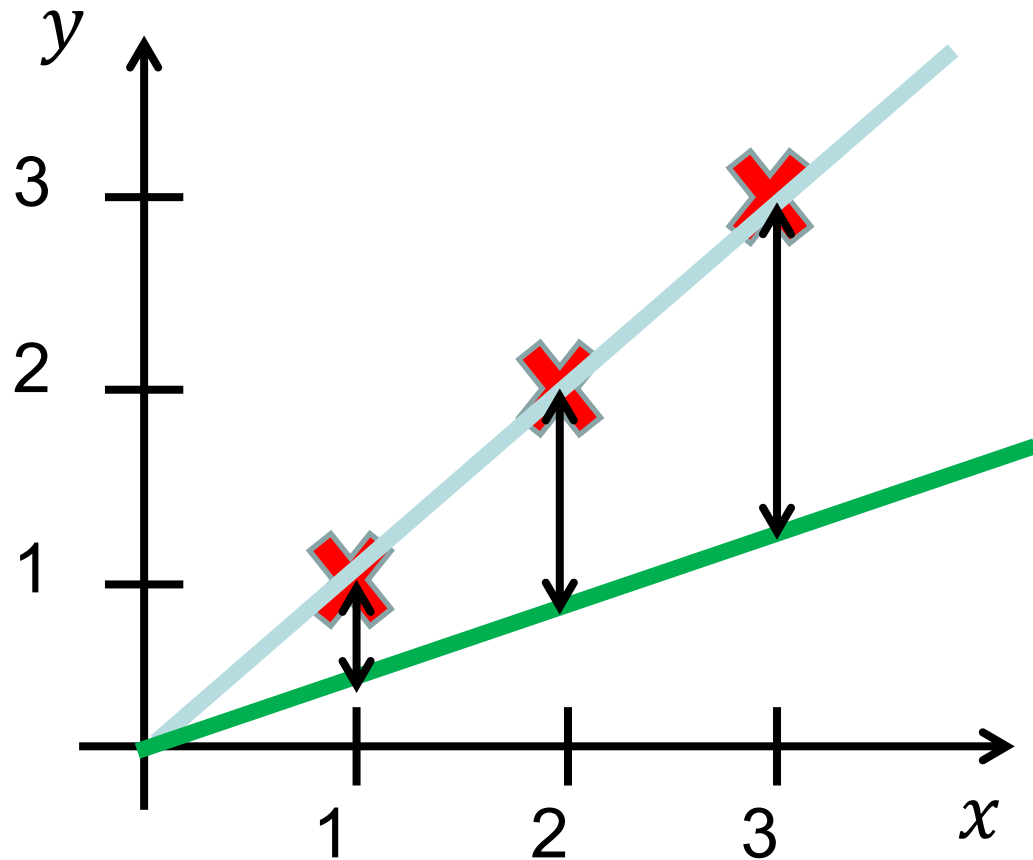


$J(\beta_1)$, function of β_1

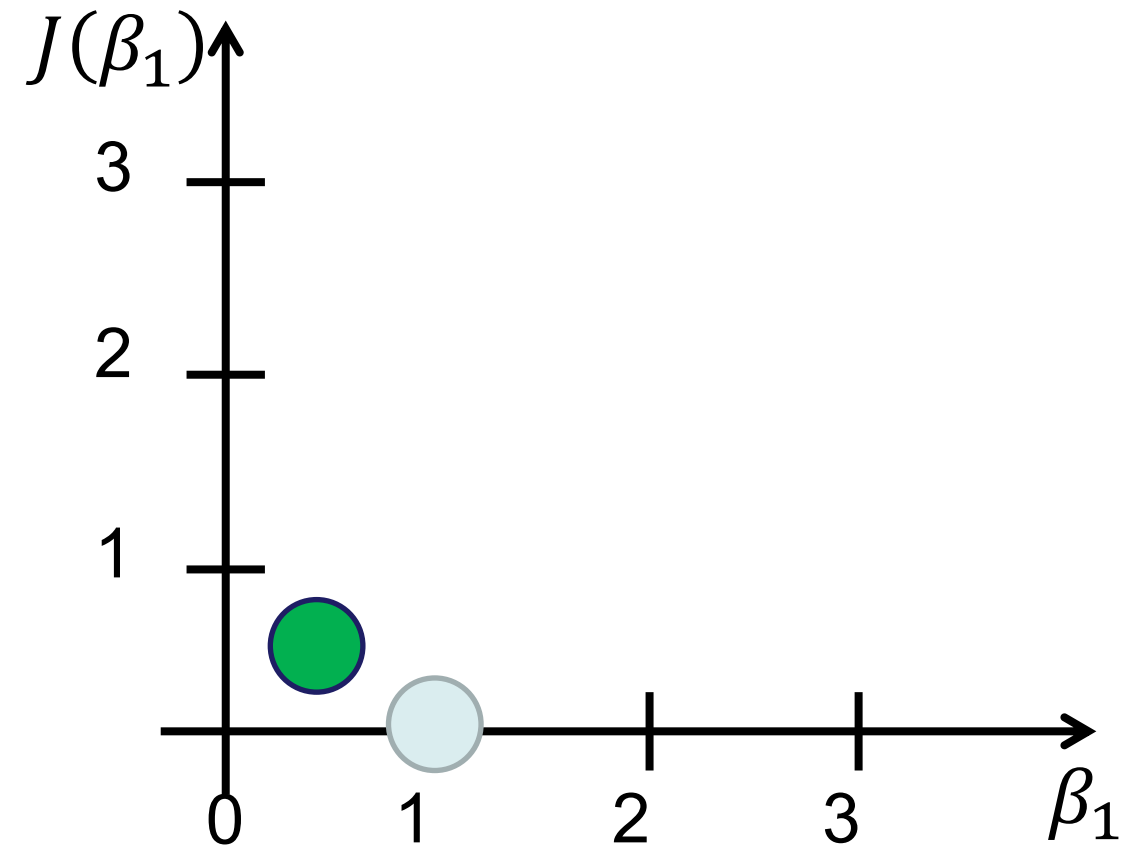


Simplified Hypothesis

$h_{\beta}(x)$, function of x

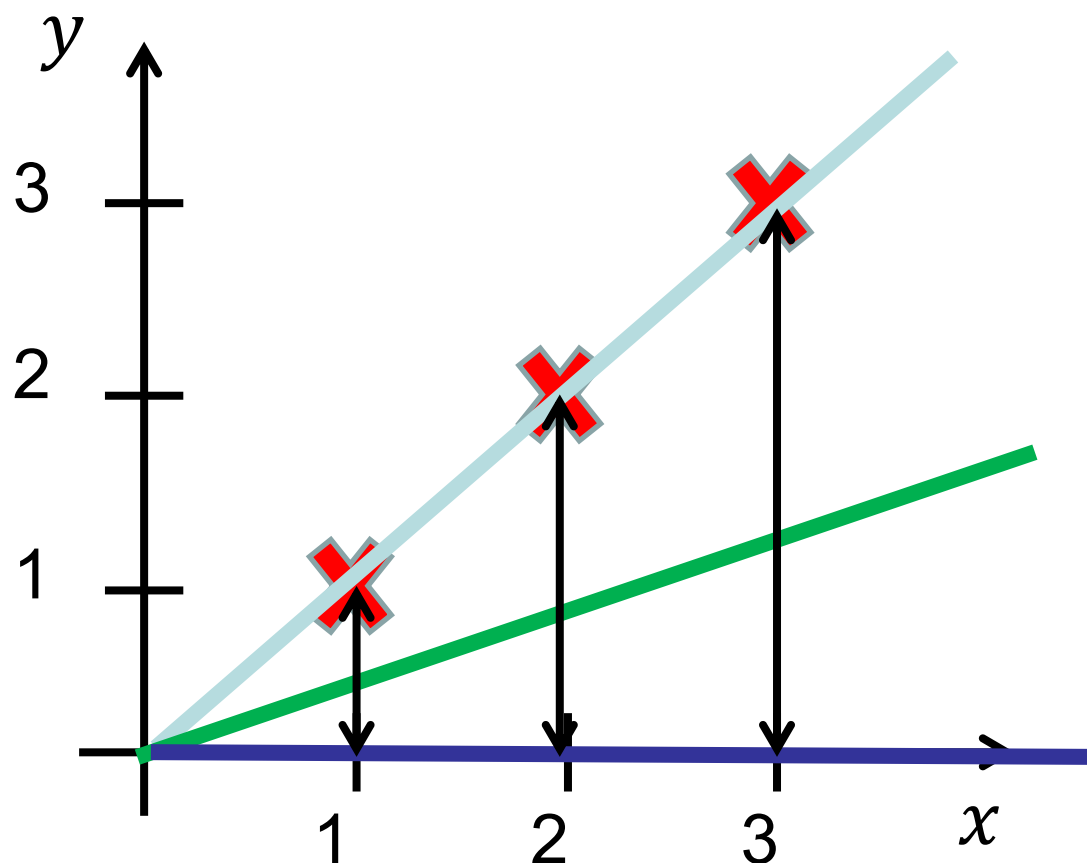


$J(\beta_1)$, function of β_1

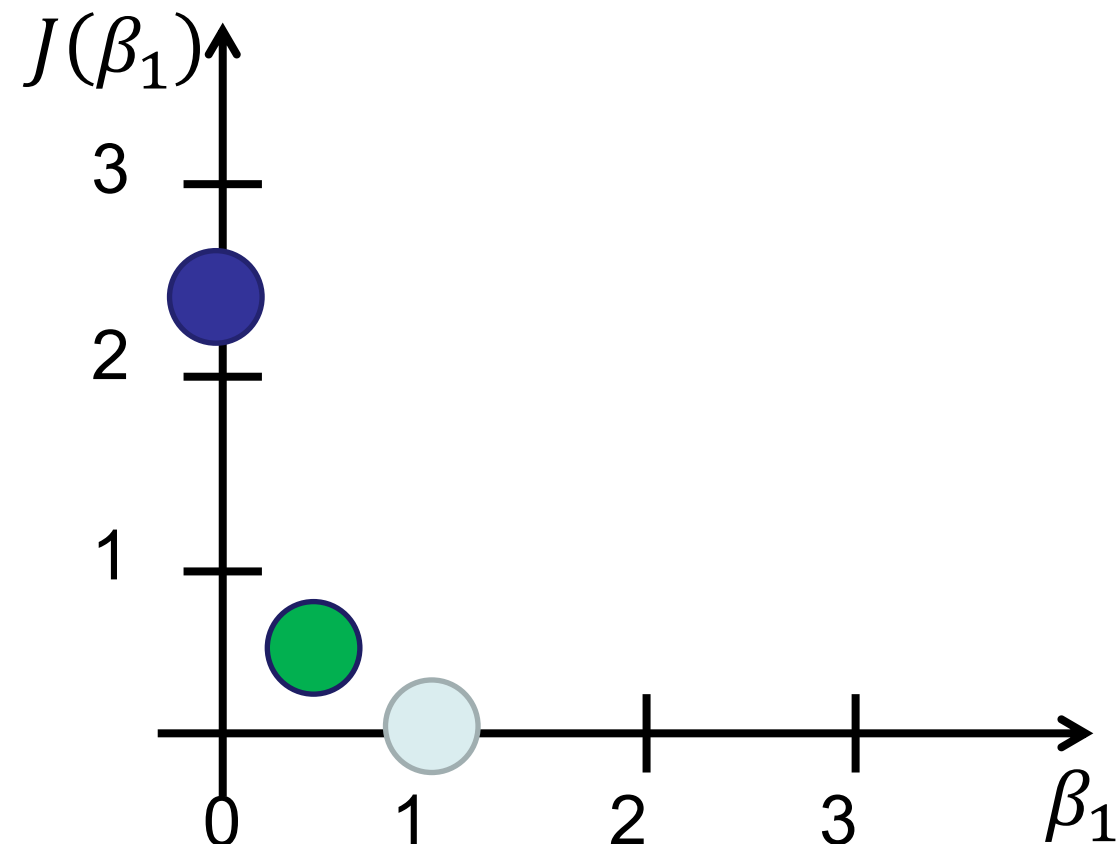


Simplified Hypothesis

$h_{\beta}(x)$, function of x

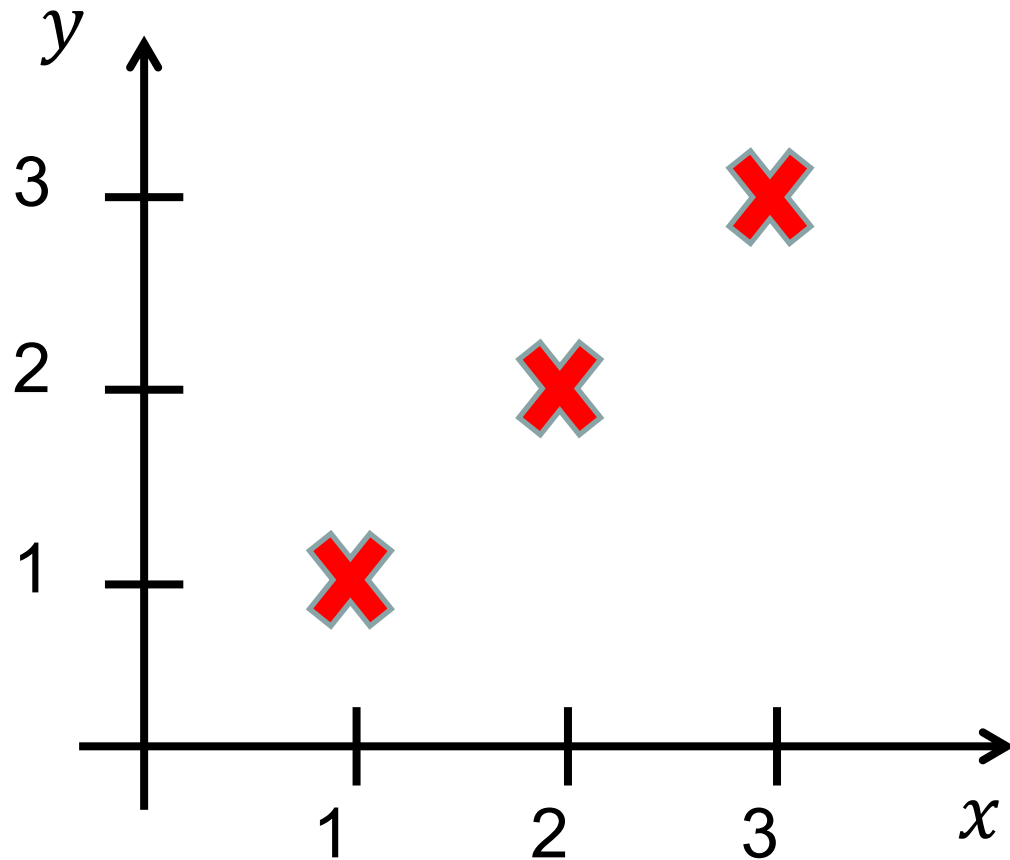


$J(\beta_1)$, function of β_1

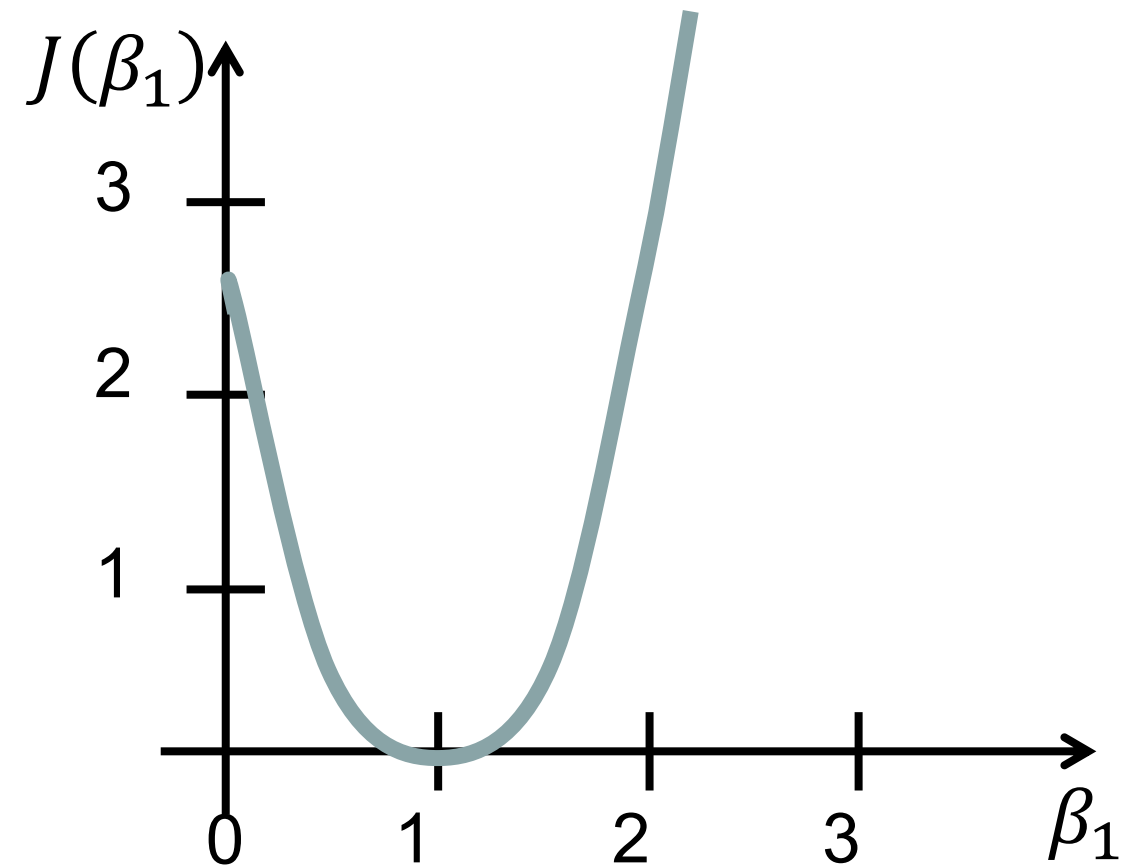


Simplified Hypothesis

$h_{\beta}(x)$, function of x



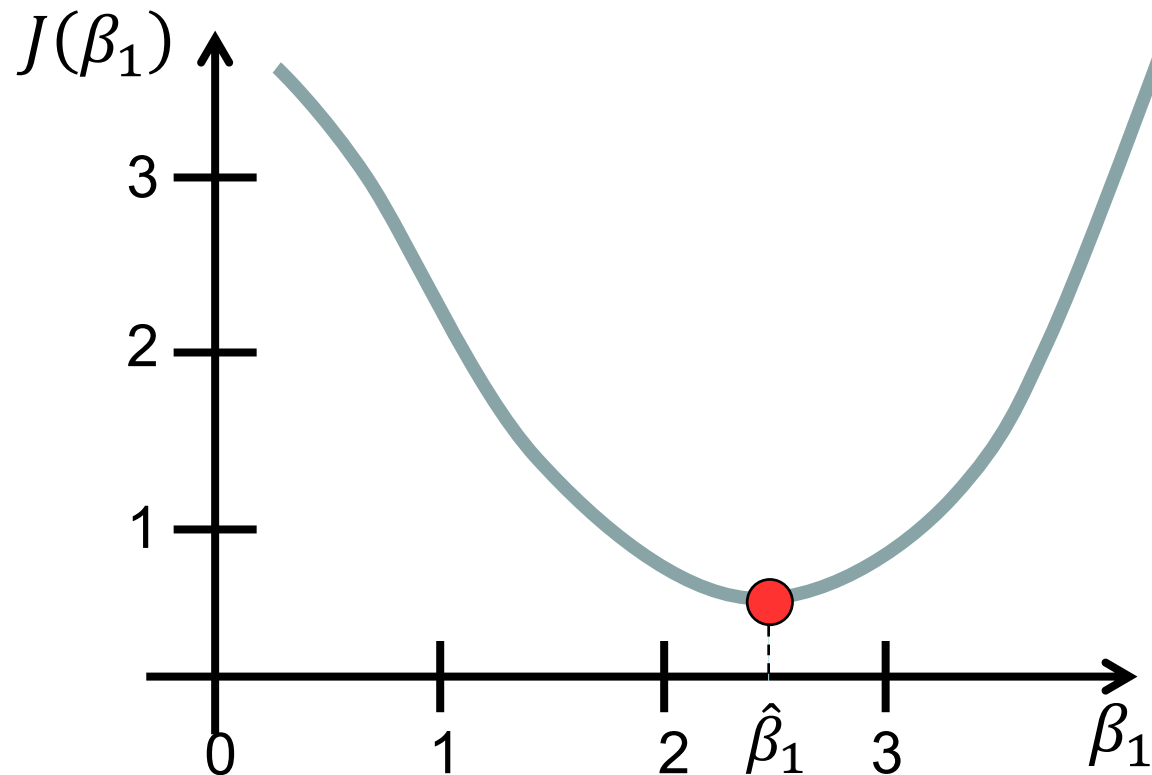
$J(\beta_1)$, function of β_1



Simplified Hypothesis

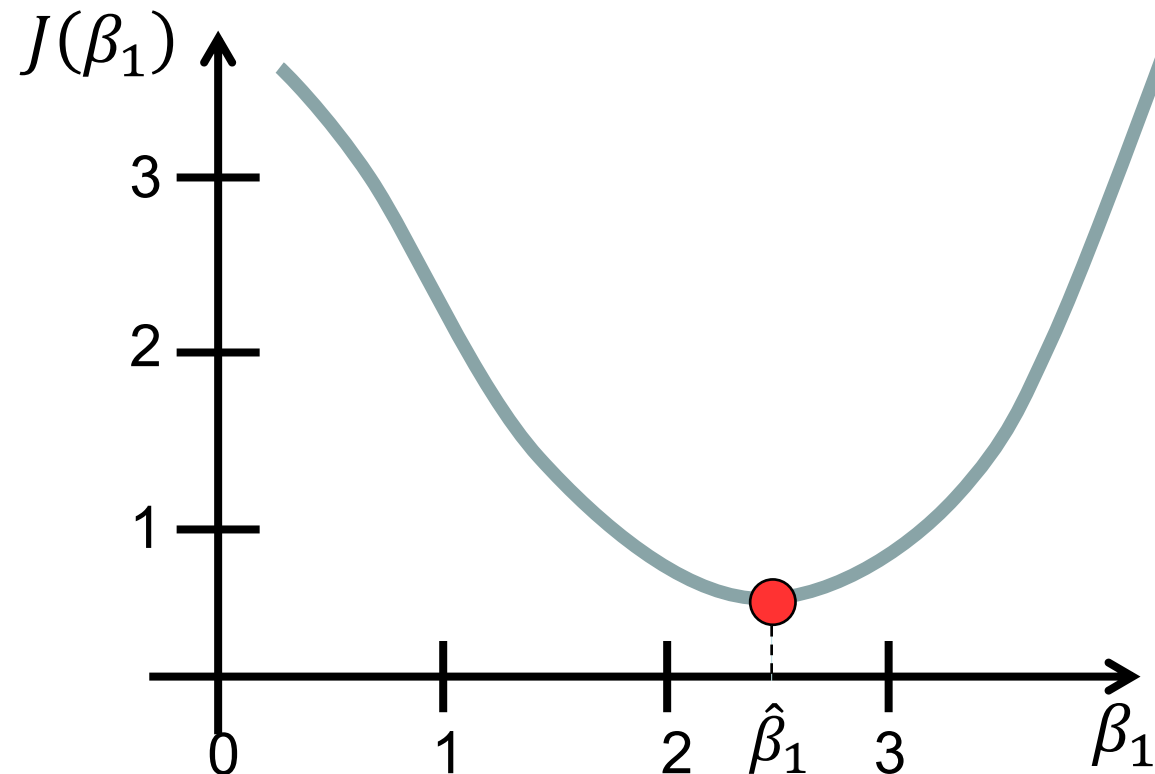
Minimize $J(\beta_1)$
 β_1

How to solve the problem?



Gradient Descent

Gradient descent is an iterative optimization algorithm for finding a **local minimum** of a **differentiable** function.



Gradient Descent

Start with some β_1

Repeat until convergence

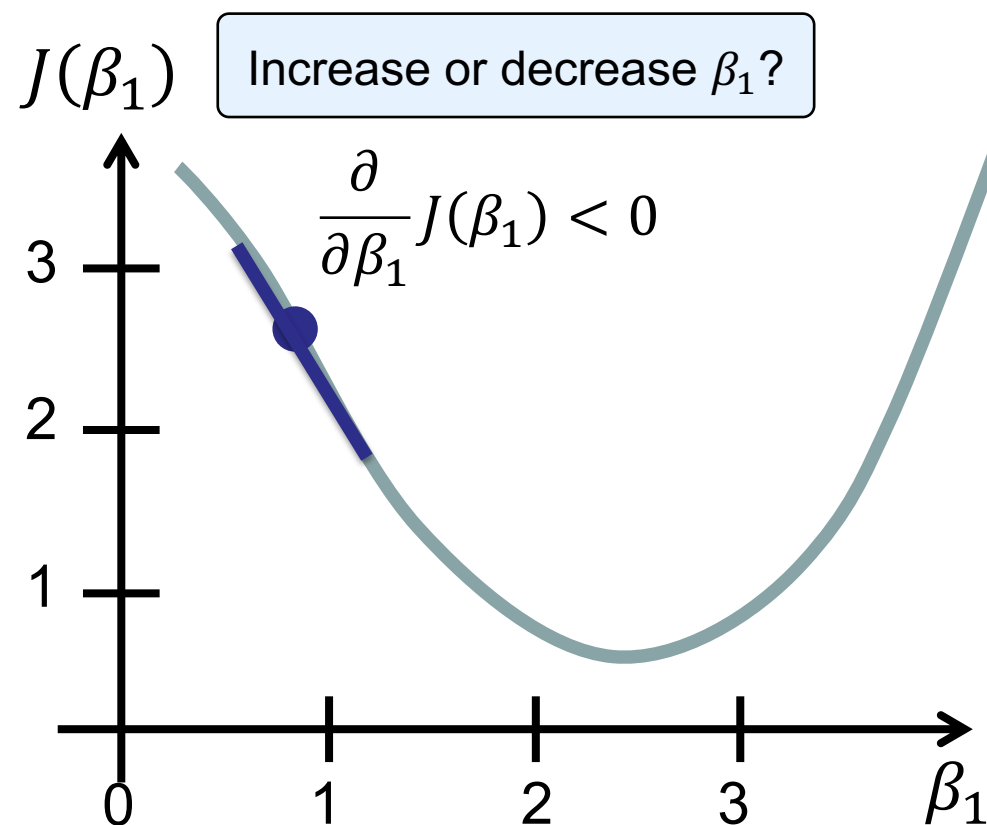
{

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_1)$$

}

α : Learning rate (step size)

$\frac{\partial}{\partial \beta_1} J(\beta_1)$: derivative



Gradient Descent

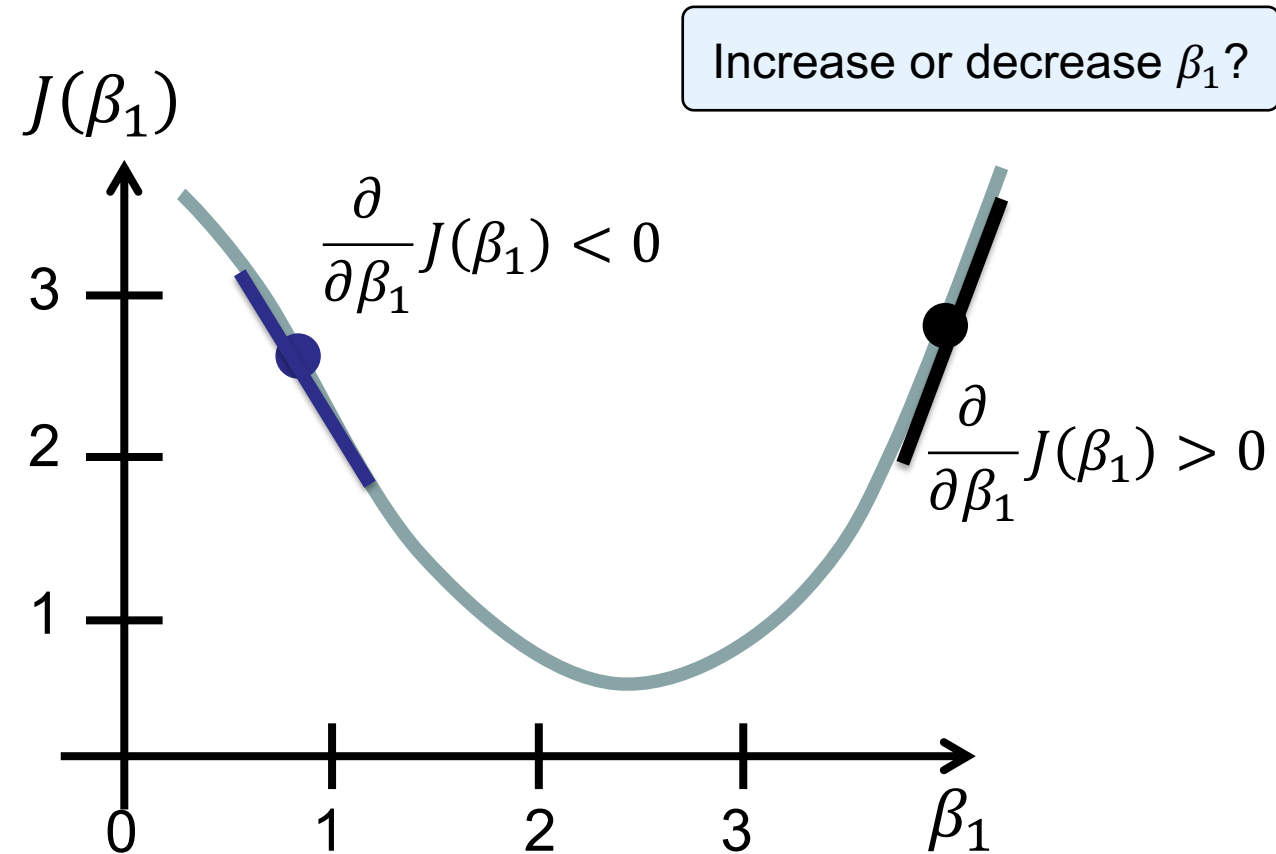
Start with some β_1

Repeat until convergence

$$\left\{ \begin{array}{l} \beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_1) \end{array} \right\}$$

α : Learning rate (step size)

$\frac{\partial}{\partial \beta_1} J(\beta_1)$: derivative



Gradient Descent

Start with some β_1

Repeat until convergence

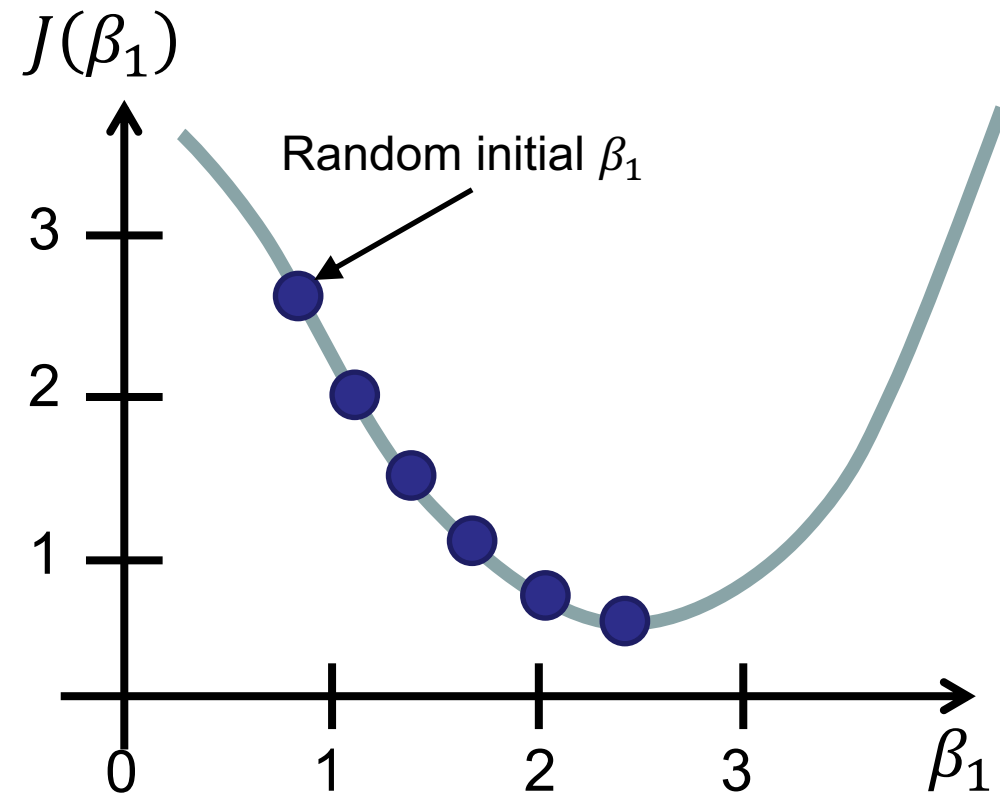
{

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_1)$$

}

α : Learning rate (step size)

$\frac{\partial}{\partial \beta_1} J(\beta_1)$: derivative



Gradient Descent

Start with some β_0, β_1

Repeat until convergence

{

$$\beta_0 := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1)$$

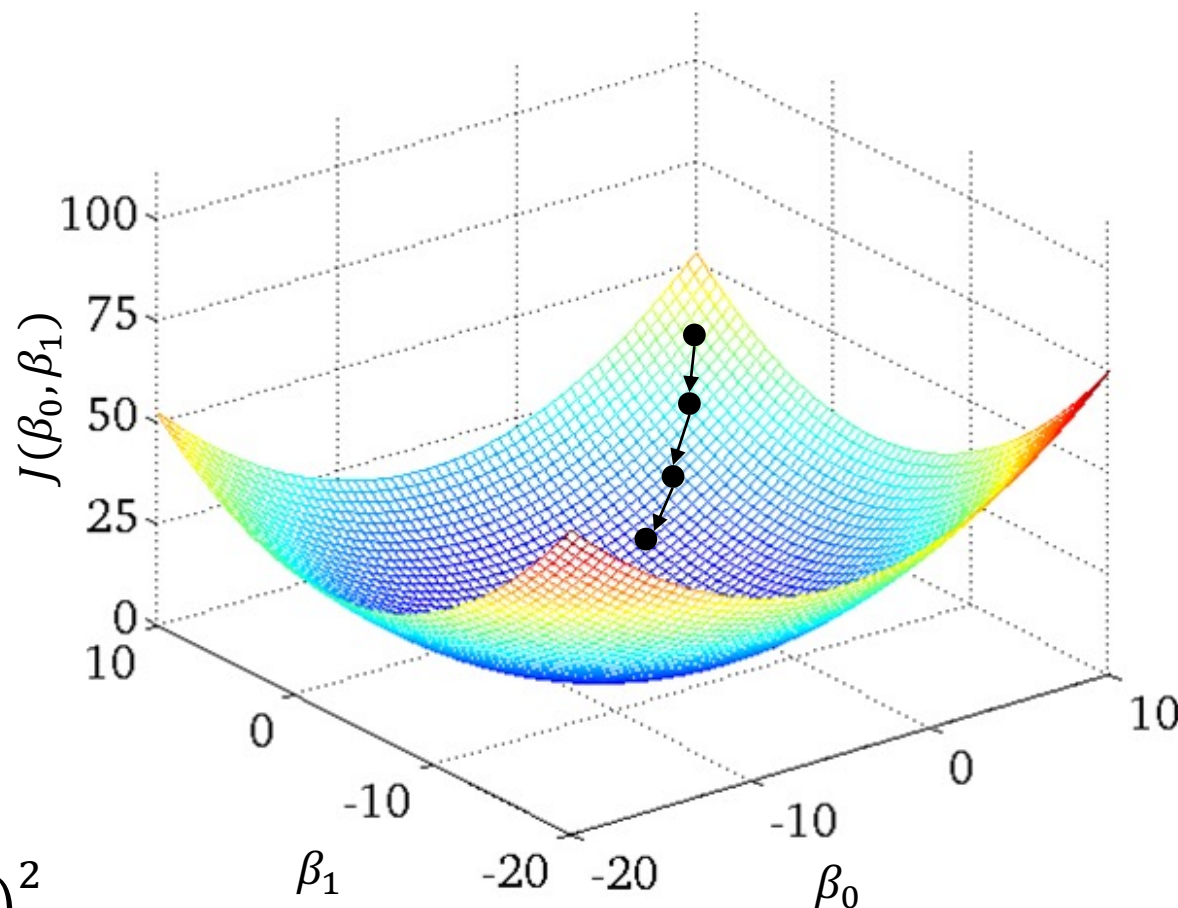
$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

}

Linear regression model

$$\text{Cost function: } J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$$

$$\text{Hypothesis: } h_{\beta}(x) = \beta_0 + \beta_1 x$$



Gradient Descent

Start with some β_0, β_1

Repeat until convergence

{

$$\beta_0 := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

}

How to compute the gradient $\frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1), \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$
(i.e., partial derivative)?

Linear regression model

Cost function: $J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$

Hypothesis: $h_{\beta}(x) = \beta_0 + \beta_1 x$

Computing Partial Derivative

- Differentiate equation $J(\beta_0, \beta_1)$ with respect to β_0

$$\begin{aligned}\frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1) &= \frac{\partial}{\partial \beta_0} \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2 \\ &= \frac{\partial}{\partial \beta_0} \frac{1}{2m} \sum_{i=1}^m (\beta_0 + \beta_1 x^i - y^i)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (\beta_0 + \beta_1 x^i - y^i) \\ &= \frac{1}{m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)\end{aligned}$$

Computing Partial Derivative

- Differentiate equation $J(\beta_0, \beta_1)$ with respect to β_1

$$\begin{aligned}\frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1) &= \frac{\partial}{\partial \beta_1} \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2 \\ &= \frac{\partial}{\partial \beta_1} \frac{1}{2m} \sum_{i=1}^m (\beta_0 + \beta_1 x^i - y^i)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (\beta_0 + \beta_1 x^i - y^i) x^i \\ &= \frac{1}{m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i) x^i\end{aligned}$$

Gradient Descent for Linear Regression

Start with some β_0, β_1

Repeat until convergence

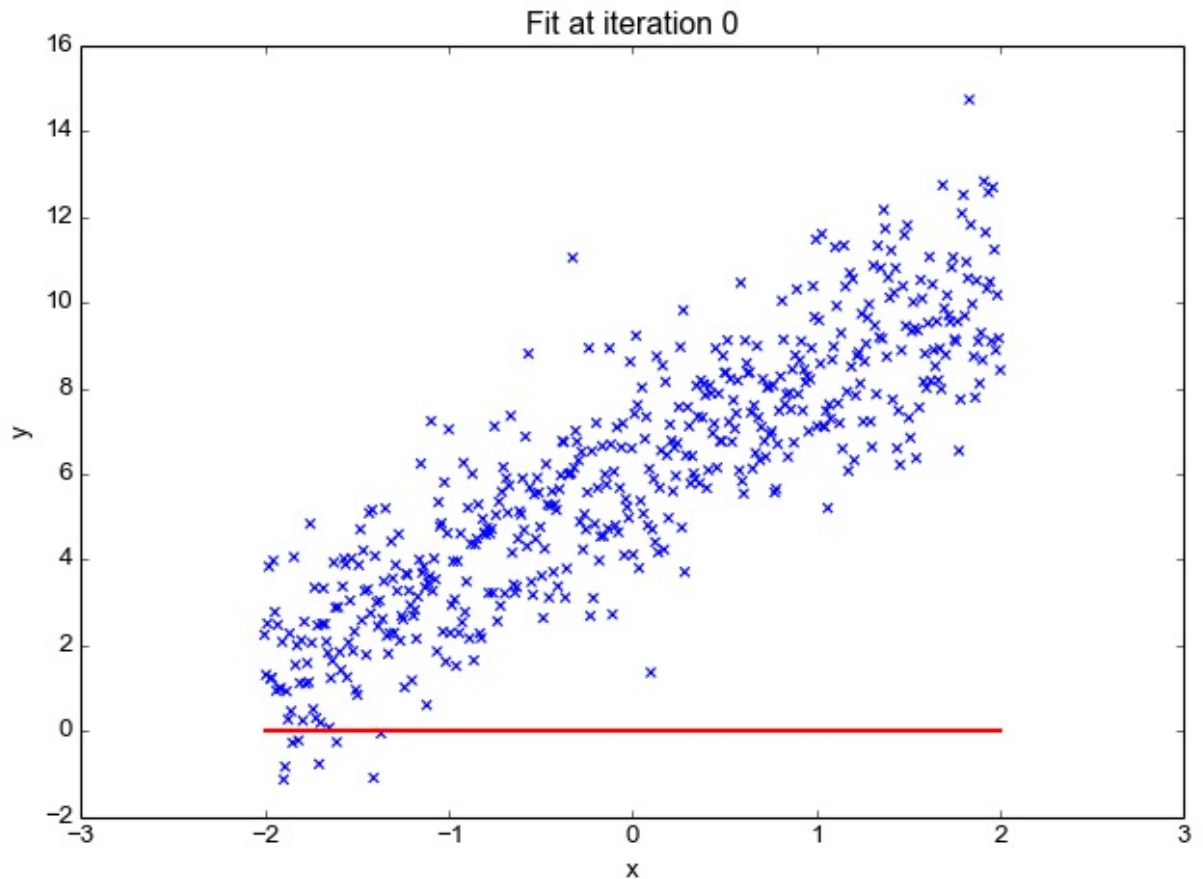
{

$$\beta_0 := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

}

Note: update β_0 and β_1 simultaneously



Linear Regression

- Hypothesis: $h_{\beta}(x) = \beta_0 + \beta_1 x$
- Cost function: $J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)^2$
- Gradient Descent:

$$\beta_0 := \beta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i)$$

$$\beta_1 := \beta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\beta}(x^i) - y^i) x^i$$

Can we write these equations in a more compact form?

Hypothesis

Hypothesis:

$$\begin{aligned}h_{\beta}(x) &= \beta_0 + \beta_1 x \\&= [1 \quad x] \times \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \\&= \mathbf{x} \times \mathbf{b}\end{aligned}$$

where $\mathbf{x} = [1 \quad x]$ $\mathbf{b} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$

Hypothesis

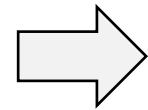
For m data points, $\hat{y}^i = h_{\beta}(x^i) = \beta_0 + \beta_1 x^i = \mathbf{x}^i \times \mathbf{b}$

$$\hat{y}^1 = \mathbf{x}^1 \mathbf{b}$$

$$\hat{y}^2 = \mathbf{x}^2 \mathbf{b}$$

...

$$\hat{y}^m = \mathbf{x}^m \mathbf{b}$$



$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^1 \\ \hat{y}^2 \\ \ddots \\ \hat{y}^m \end{bmatrix} = \begin{bmatrix} \mathbf{x}^1 \mathbf{b} \\ \mathbf{x}^2 \mathbf{b} \\ \ddots \\ \mathbf{x}^m \mathbf{b} \end{bmatrix} = \begin{bmatrix} \mathbf{x}^1 \\ \mathbf{x}^2 \\ \ddots \\ \mathbf{x}^m \end{bmatrix} \mathbf{b} = \mathbf{X} \times \mathbf{b}$$

$$\text{where } \mathbf{X} = \begin{bmatrix} \mathbf{x}^1 \\ \mathbf{x}^2 \\ \ddots \\ \mathbf{x}^m \end{bmatrix} = \begin{bmatrix} 1 & x^1 \\ 1 & x^2 \\ \ddots & \ddots \\ 1 & x^m \end{bmatrix}$$

Cost Function

Cost function $J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^i - y^i)^2$

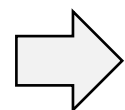
$$= \frac{1}{2m} \sum_{i=1}^m (\hat{y}^i - y^i) \times (\hat{y}^i - y^i)$$
$$= \frac{1}{2m} [(\hat{y}^1 - y^1) \quad \dots \quad (\hat{y}^m - y^m)] \times \begin{bmatrix} (\hat{y}^1 - y^1) \\ \vdots \\ (\hat{y}^m - y^m) \end{bmatrix}$$
$$= \frac{1}{2m} (\hat{\mathbf{y}} - \mathbf{y})^T \times (\hat{\mathbf{y}} - \mathbf{y})$$

where $\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^1 \\ \hat{y}^2 \\ \vdots \\ \hat{y}^m \end{bmatrix}$ $\mathbf{y} = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix}$

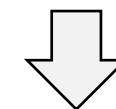
Gradient Descent

$$\beta_0 := \beta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)$$

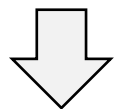
$$\beta_1 := \beta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i) x^i$$



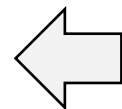
$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} - \alpha \frac{1}{m} \begin{bmatrix} \sum_{i=1}^m (\hat{y}^i - y^i) \\ \sum_{i=1}^m (\hat{y}^i - y^i) x^i \end{bmatrix}$$



$$\mathbf{b} = \mathbf{b} - \alpha \frac{1}{m} \mathbf{X}^T \times (\hat{\mathbf{y}} - \mathbf{y})$$



$$\mathbf{b} = \mathbf{b} - \alpha \frac{1}{m} \mathbf{X}^T \times (\mathbf{X} \times \mathbf{b} - \mathbf{y})$$



$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} - \alpha \frac{1}{m} \begin{bmatrix} 1 & \cdots & 1 \\ x^1 & \cdots & x^m \end{bmatrix} \begin{bmatrix} (\hat{y}^1 - y^1) \\ \vdots \\ (\hat{y}^m - y^m) \end{bmatrix}$$

Linear Regression

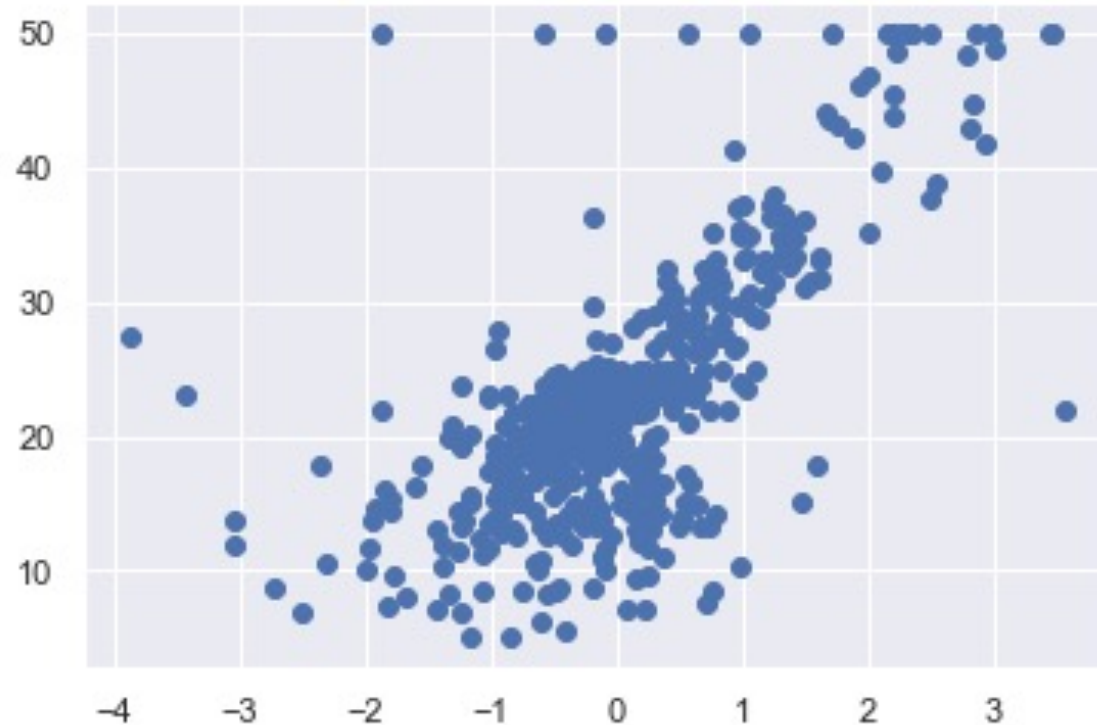
- Hypothesis: $\hat{\mathbf{y}} = \mathbf{X} \times \mathbf{b}$
- Cost Function: $J(\beta_0, \beta_1) = \frac{1}{2m} (\hat{\mathbf{y}} - \mathbf{y})^T \times (\hat{\mathbf{y}} - \mathbf{y})$
- Gradient Descent: $\mathbf{b} = \mathbf{b} - \alpha \frac{1}{m} \mathbf{X}^T \times (\mathbf{X} \times \mathbf{b} - \mathbf{y})$

where

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^1 \\ \hat{y}^2 \\ \hat{y}^3 \\ \vdots \\ \hat{y}^m \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x^1 \\ 1 & x^2 \\ \ddots & \ddots \\ 1 & x^m \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y^1 \\ y^2 \\ y^3 \\ \vdots \\ y^m \end{bmatrix}$$

Cohort Problem CS0

CS0. *Plot:* Read data for Boston Housing Prices and write a function `get_features_targets()` to get the columns for the features and the targets from the input argument data frame.



Cohort Problem CS1

CS1. Cost Function: Write `def compute_cost(X, y, beta)` to compute the cost function of a linear regression model.

$$\begin{array}{ll} \text{Cost Function:} & J(\beta_0, \beta_1) = \frac{1}{2m} (\hat{\mathbf{y}} - \mathbf{y})^T \times (\hat{\mathbf{y}} - \mathbf{y}) \\ \text{Hypothesis:} & \hat{\mathbf{y}} = \mathbf{X} \times \mathbf{b} \end{array}$$

Note: m is the number of data points (i.e., number of rows in \mathbf{X})

Cohort Problem CS2

CS2. *Gradient Descent:* Write a function called

```
def gradient_descent(X, y, beta, alpha, num_iters):
```

- `X`: is a 2-D numpy array for the features
- `y`: is a vector array for the target
- `beta`: is a column vector for the initial guess of the parameters
- `alpha`: is the learning rate
- `num_iters`: is the number of iteration to perform

$$\text{Gradient Descent: } \mathbf{b} = \mathbf{b} - \alpha \frac{1}{m} \mathbf{X}^T \times (\mathbf{X} \times \mathbf{b} - \mathbf{y})$$

Thank You!