

Gradient Descent

Multiple Regression using Gradient Descent

Rahul Singh
rsingh@arrsingh.com

Simple Linear Regression

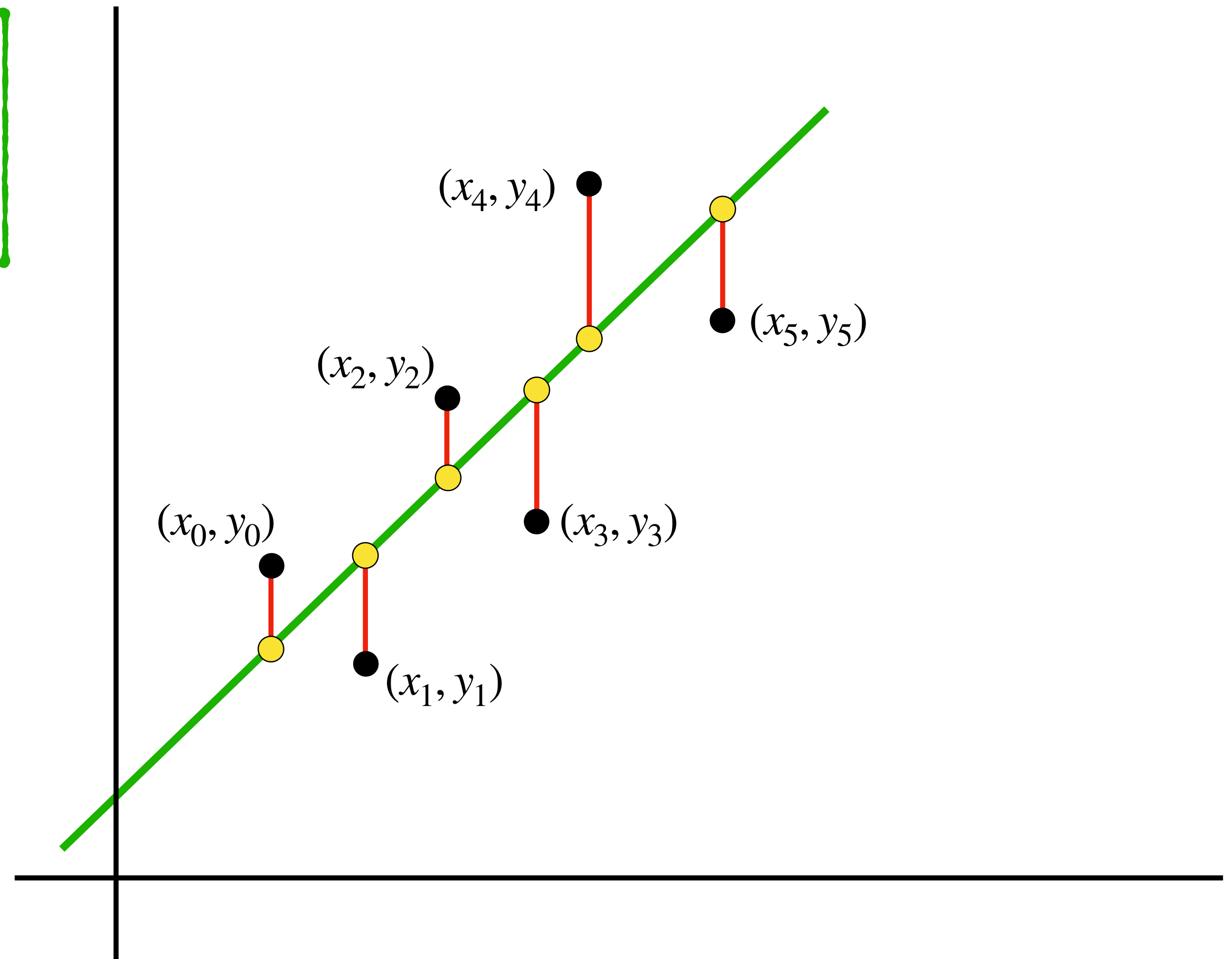
The Problem Statement:

Simple Linear Regression: Find the values of β_0 and β_1 such that the **Mean Squared Error (MSE)** is minimized.

The line of best fit is $\hat{y} = \beta_0 + \beta_1 \hat{x}$

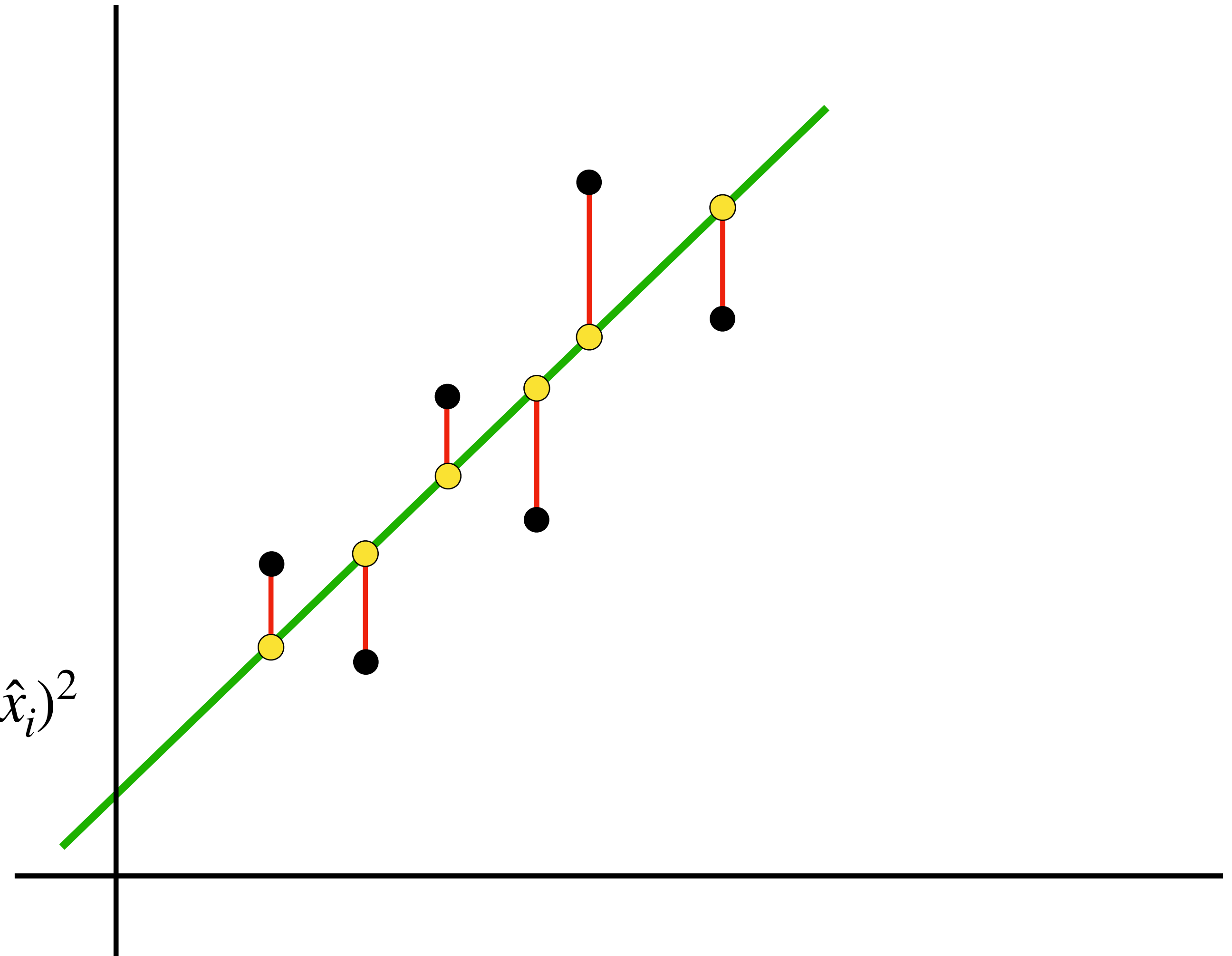
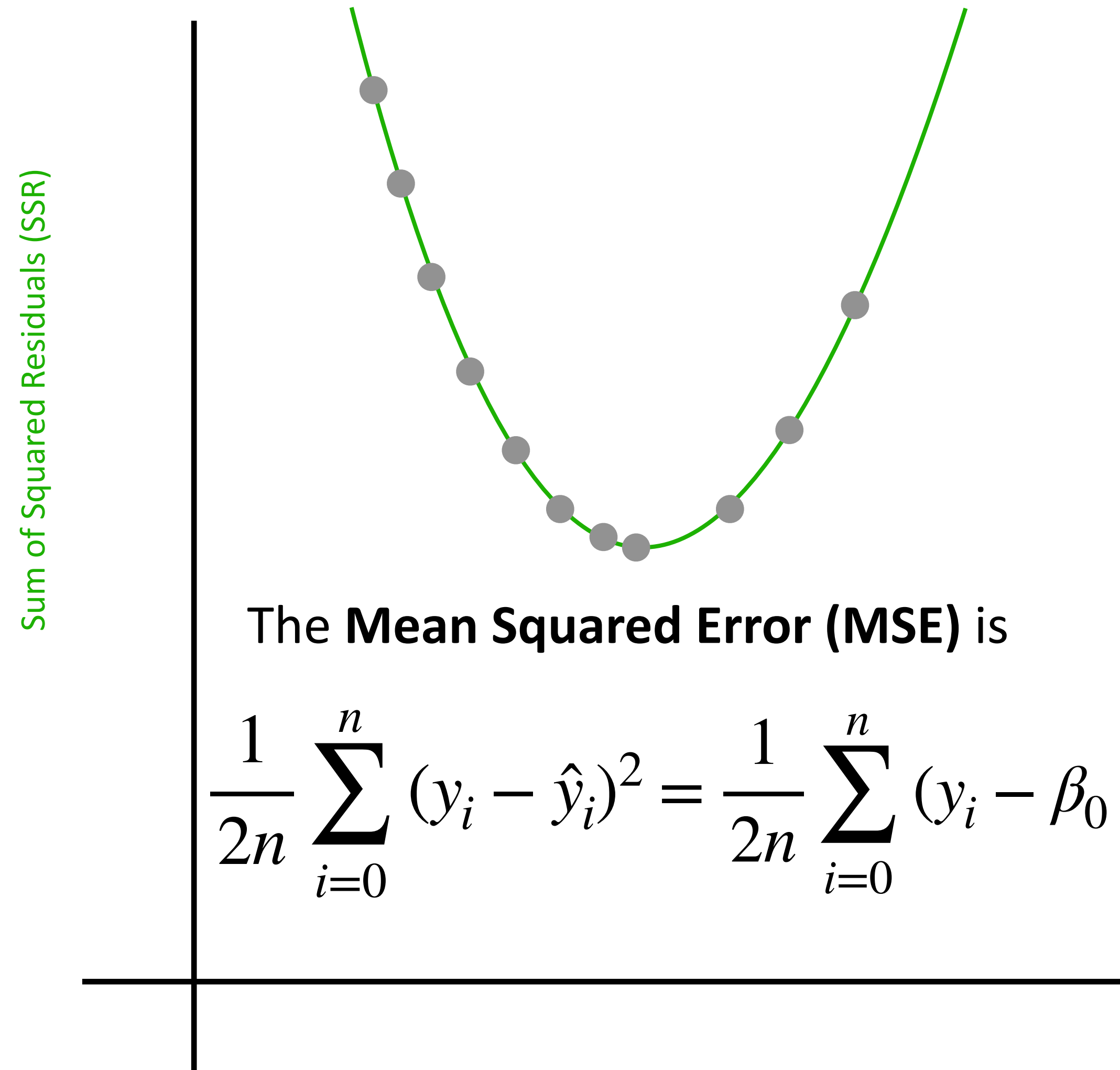
Mean Squared Error (MSE)

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \hat{y}_i)^2 = \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_i)^2$$

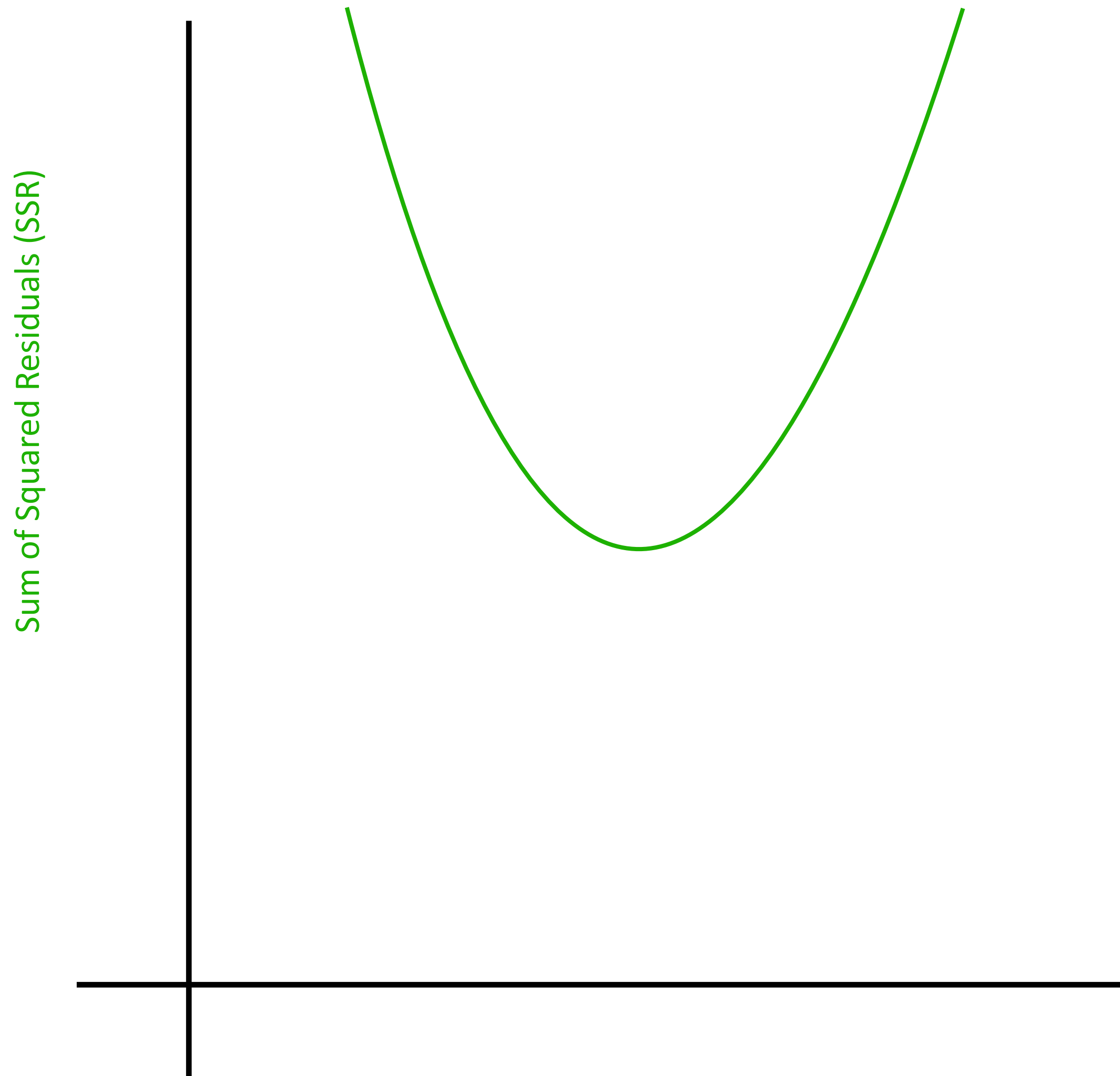


Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve

Simple Linear Regression



Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve



Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β_0 and β_1

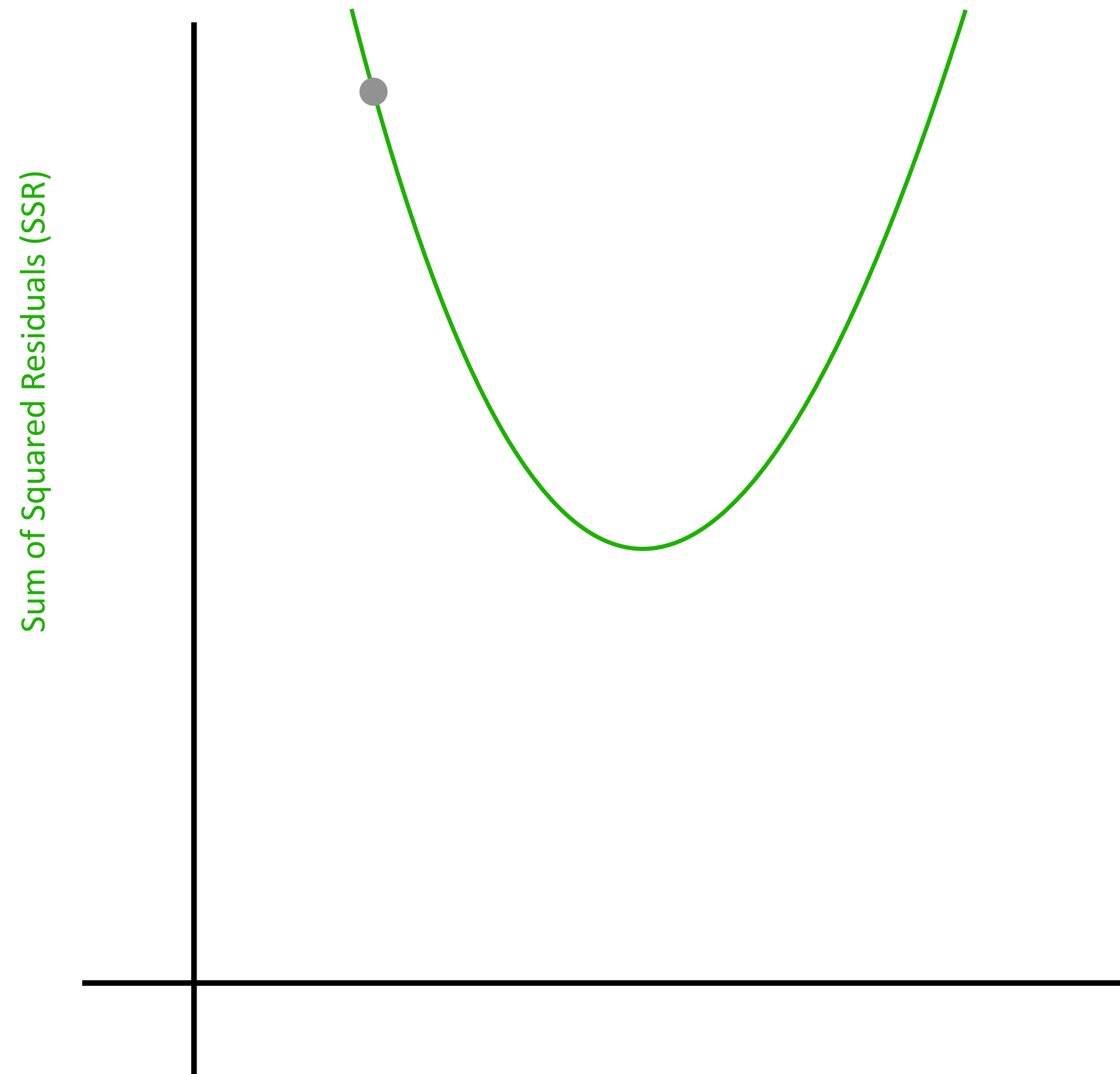
Step 2: Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 and β_1 by subtracting the step size

Step 5: Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve



Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β_0 and β_1

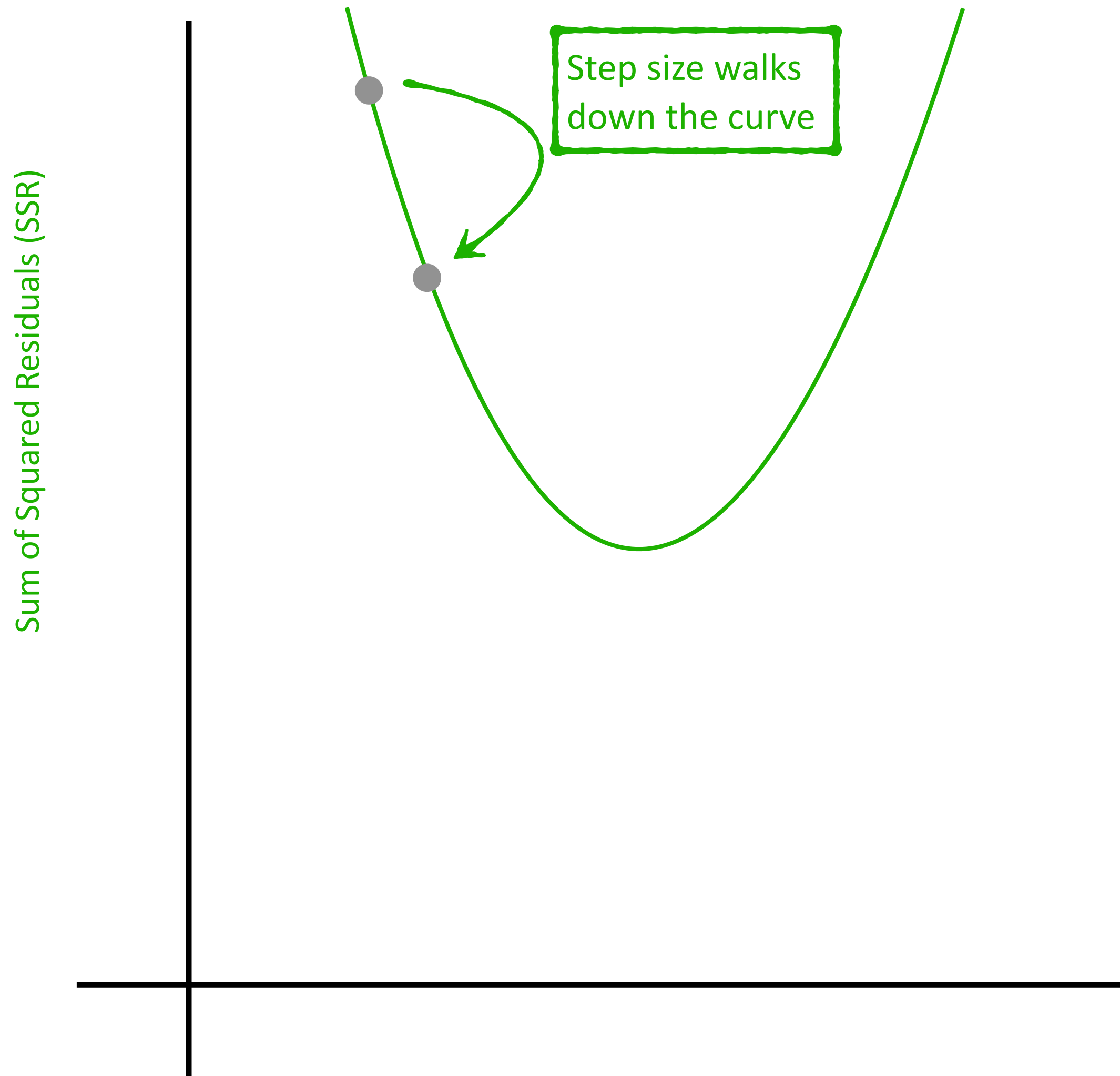
Step 2: Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 and β_1 by subtracting the step size

Step 5: Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve

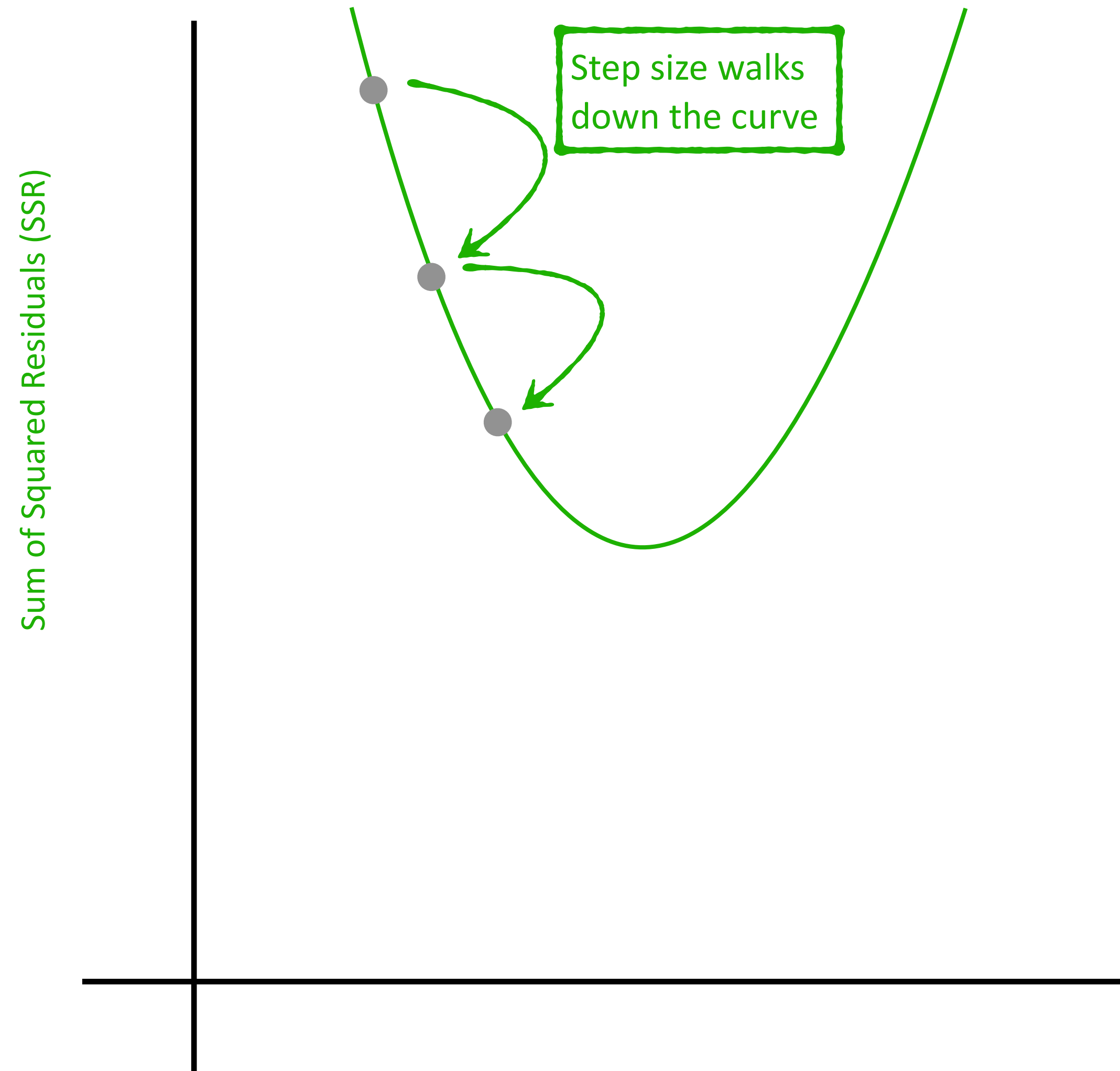


Gradient Descent

Gradient Descent: Basic Concept

- Step 1:** Start with random values for β_0 and β_1
- Step 2:** Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point
- Step 3:** Calculate a step size that is proportional to the slope
- Step 4:** Calculate new values for β_0 and β_1 by subtracting the step size
- Step 5:** Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve

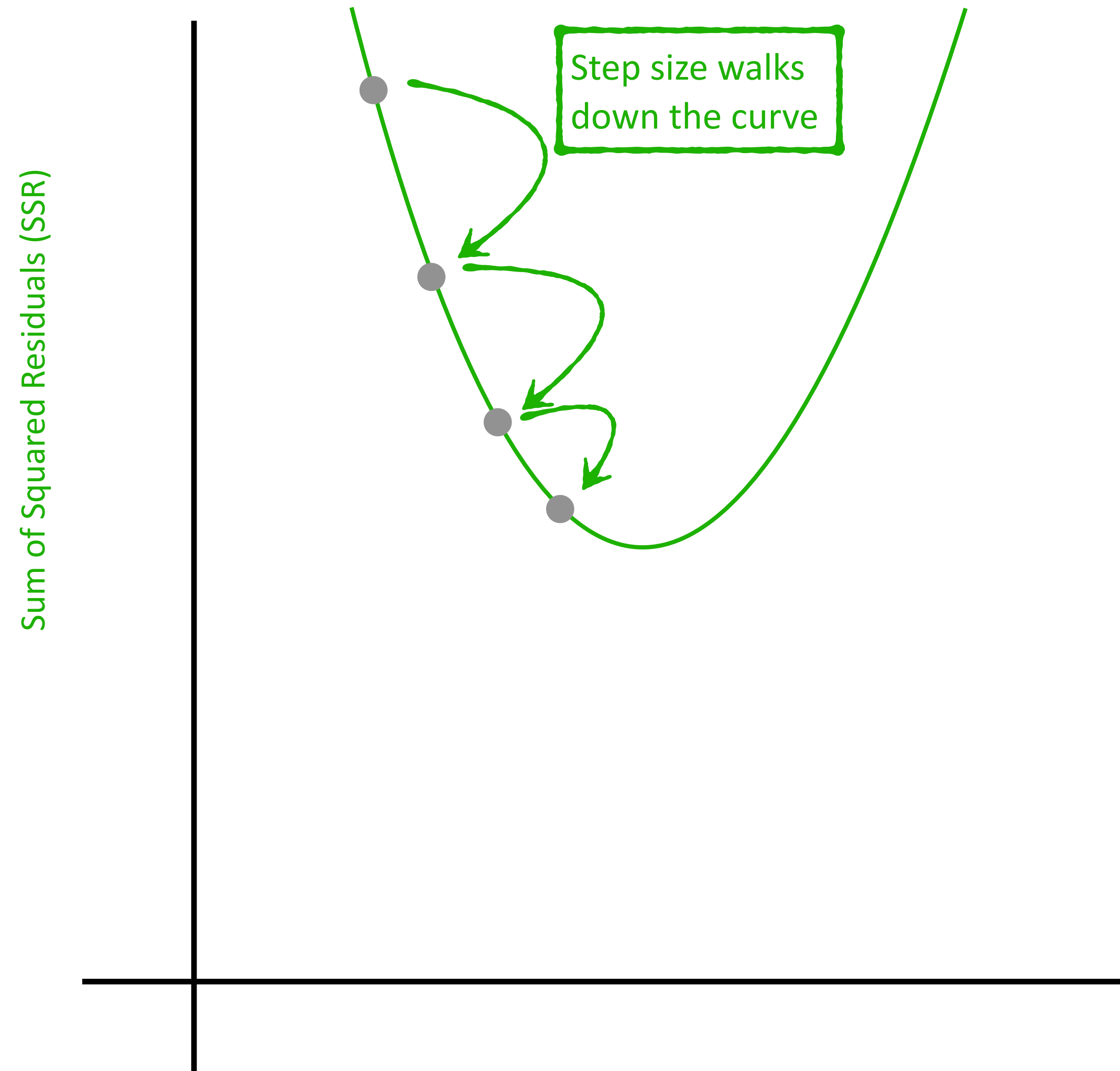


Gradient Descent

Gradient Descent: Basic Concept

- Step 1:** Start with random values for β_0 and β_1
- Step 2:** Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point
- Step 3:** Calculate a step size that is proportional to the slope
- Step 4:** Calculate new values for β_0 and β_1 by subtracting the step size
- Step 5:** Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve

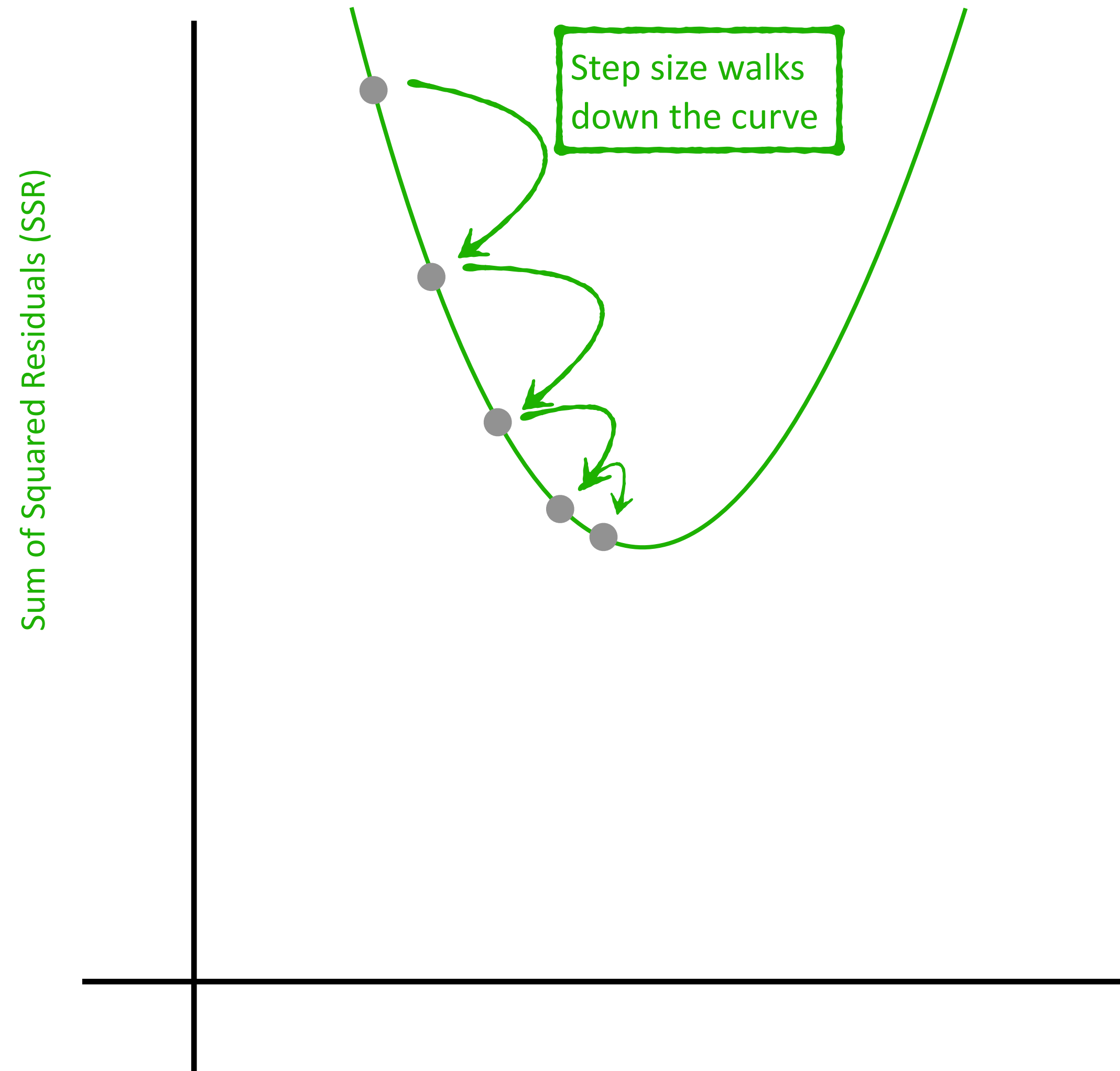


Gradient Descent

Gradient Descent: Basic Concept

- Step 1:** Start with random values for β_0 and β_1
- Step 2:** Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point
- Step 3:** Calculate a step size that is proportional to the slope
- Step 4:** Calculate new values for β_0 and β_1 by subtracting the step size
- Step 5:** Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve

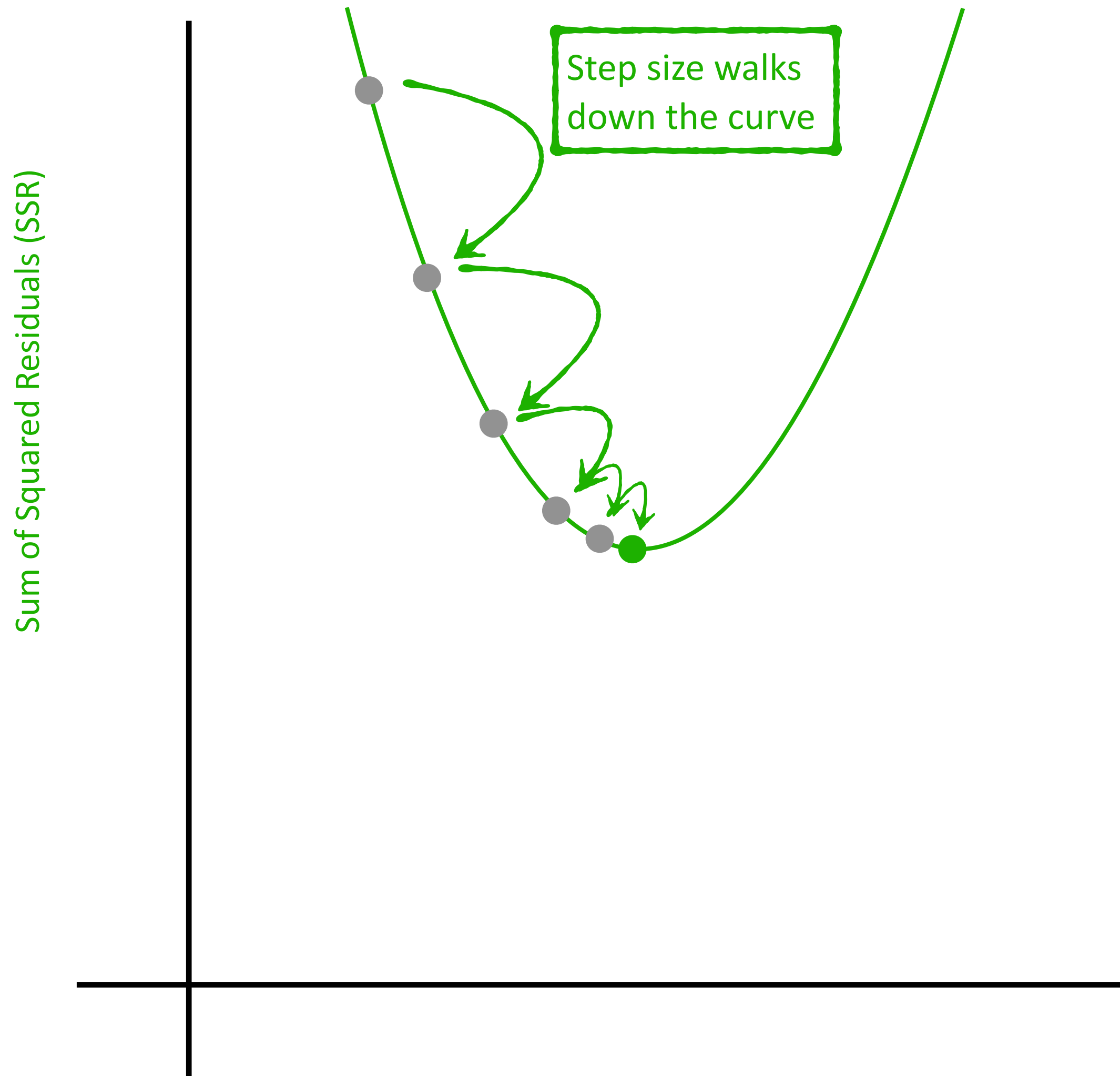


Gradient Descent

Gradient Descent: Basic Concept

- Step 1:** Start with random values for β_0 and β_1
- Step 2:** Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point
- Step 3:** Calculate a step size that is proportional to the slope
- Step 4:** Calculate new values for β_0 and β_1 by subtracting the step size
- Step 5:** Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve

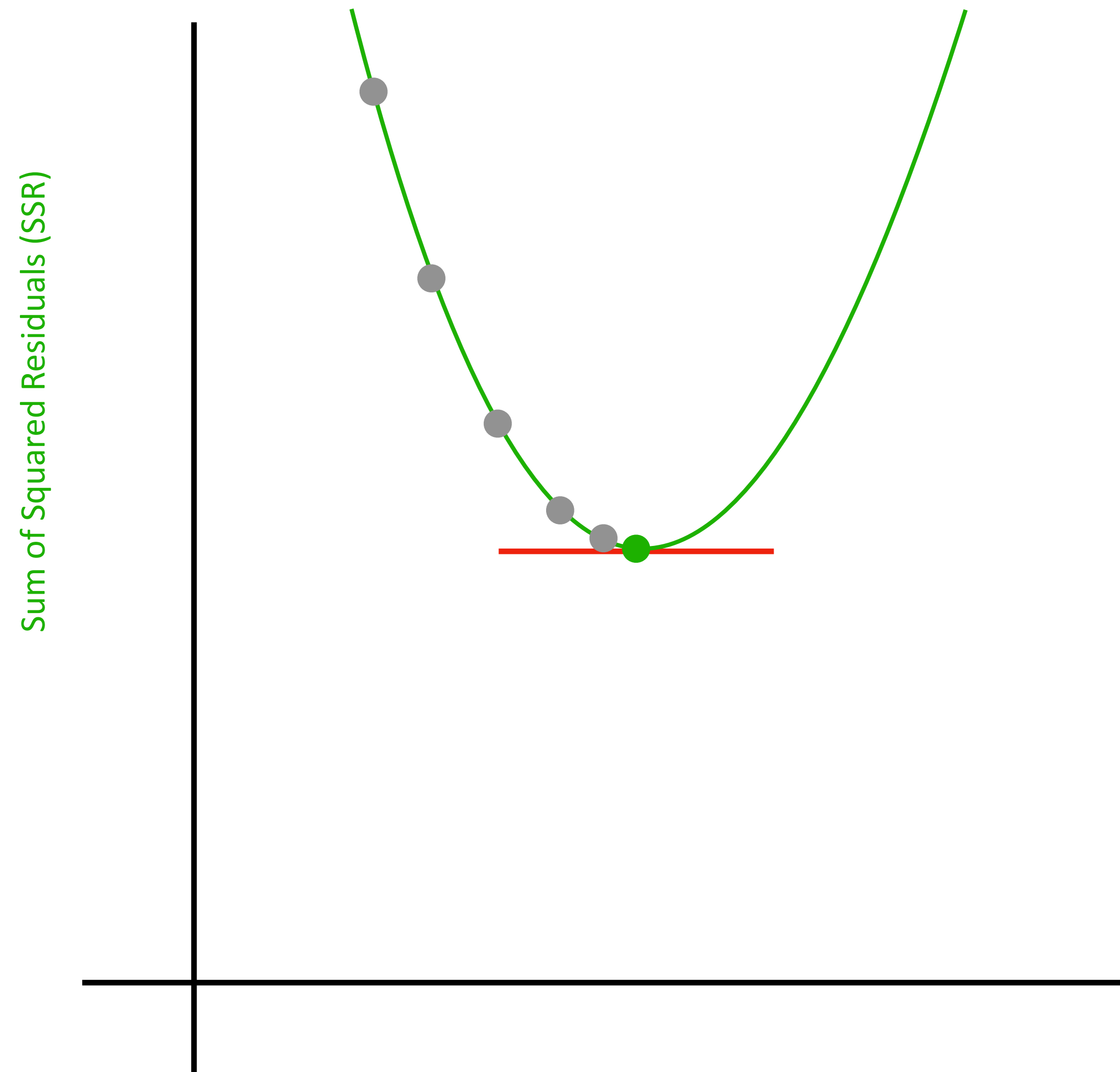


Gradient Descent

Gradient Descent: Basic Concept

- Step 1:** Start with random values for β_0 and β_1
- Step 2:** Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point
- Step 3:** Calculate a step size that is proportional to the slope
- Step 4:** Calculate new values for β_0 and β_1 by subtracting the step size
- Step 5:** Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve



Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β_0 and β_1

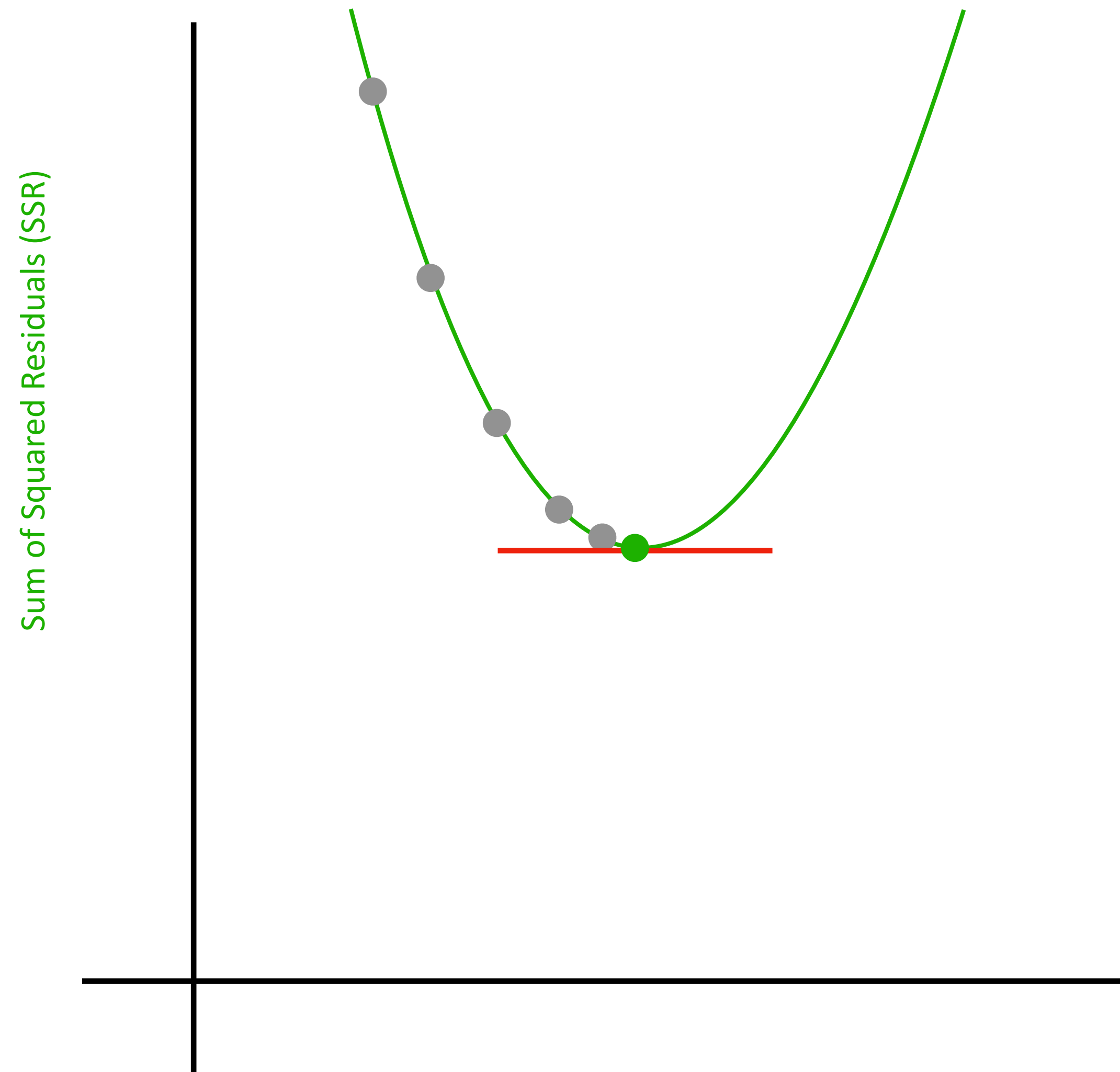
Step 2: Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 and β_1 by subtracting the step size

Step 5: Go to step 2 and repeat

Mean Squared Error (MSE) for various values of β_0 and β_1 follows this curve



Gradient Descent

Gradient Descent: Basic Concept

Gradient Descent continues in this manner until the step size is close to zero or a fixed number of iterations

A linear model in 2 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1$$

Has 2 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i})^2$$

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 and β_1

Step 2: Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 and β_1 by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in 2 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1$$

Has 2 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i})^2$$

Compute 2
partial derivatives

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 and β_1

Step 2: Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$\frac{\partial}{\partial \beta_0} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i})^2$$

$$\frac{\partial}{\partial \beta_1} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i})^2$$

A linear model in 2 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1$$

Has 2 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i})^2$$

Compute 2 step sizes

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 and β_1

Step 2: Compute the partial derivative of the MSE w.r.t β_0 and β_1 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$step_size_{\beta_0} = \frac{\partial}{\partial \beta_0} MSE \times learning_rate$$

$$step_size_{\beta_1} = \frac{\partial}{\partial \beta_1} MSE \times learning_rate$$

A linear model in 3 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2$$

Has 3 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i})^2$$

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 , β_1 and β_2

Step 2: Compute the partial derivative of the MSE w.r.t β_0 , β_1 and β_2 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 , β_1 and β_2 by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in 3 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2$$

Has 3 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i})^2$$

Compute 3
partial derivatives

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 , β_1 and β_2

Step 2: Compute the partial derivative of the MSE w.r.t β_0 , β_1 and β_2 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$\frac{\partial}{\partial \beta_0} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i})^2$$

$$\frac{\partial}{\partial \beta_1} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i})^2$$

$$\frac{\partial}{\partial \beta_2} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i})^2$$

A linear model in 3 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2$$

Has 3 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i})^2$$

Compute 3 step sizes

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 , β_1 and β_2

Step 2: Compute the partial derivative of the MSE w.r.t β_0 , β_1 and β_2 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$step_size_{\beta_0} = \frac{\partial}{\partial \beta_0} MSE \times learning_rate$$

$$step_size_{\beta_1} = \frac{\partial}{\partial \beta_1} MSE \times learning_rate$$

$$step_size_{\beta_2} = \frac{\partial}{\partial \beta_2} MSE \times learning_rate$$

A linear model in 4 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3$$

Has 4 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 , β_1 , β_2 and β_3

Step 2: Compute the partial derivative of the MSE w.r.t β_0 , β_1 , β_2 and β_3 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 , β_1 , β_2 and β_3 by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in 4 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3$$

Has 4 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

Compute 4
partial derivatives

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 , β_1 , β_2 and β_3

Step 2: Compute the partial derivative of the MSE w.r.t β_0 , β_1 , β_2 and β_3 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$\frac{\partial}{\partial \beta_0} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

$$\frac{\partial}{\partial \beta_1} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

$$\frac{\partial}{\partial \beta_2} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

$$\frac{\partial}{\partial \beta_3} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

A linear model in 4 dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3$$

Has 4 parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i})^2$$

Compute 4 step sizes

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for β_0 , β_1 , β_2 and β_3

Step 2: Compute the partial derivative of the MSE w.r.t β_0 , β_1 , β_2 and β_3 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$step_size_{\beta_0} = \frac{\partial}{\partial \beta_0} MSE \times learning_rate$$

$$step_size_{\beta_1} = \frac{\partial}{\partial \beta_1} MSE \times learning_rate$$

$$step_size_{\beta_2} = \frac{\partial}{\partial \beta_2} MSE \times learning_rate$$

$$step_size_{\beta_3} = \frac{\partial}{\partial \beta_3} MSE \times learning_rate$$

A linear model in $k+1$ dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

Has $k + 1$ parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for $\beta_0, \beta_1, \beta_2$ and β_3

Step 2: Compute the partial derivative of the MSE w.r.t $\beta_0, \beta_1, \beta_2$ and β_3 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for $\beta_0, \beta_1, \beta_2$ and β_3 by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in $k+1$ dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

Has $k + 1$ parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

Compute $k + 1$
partial derivatives

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for $\beta_0, \beta_1, \beta_2$ and β_3

Step 2: Compute the partial derivative of the MSE w.r.t $\beta_0, \beta_1, \beta_2$ and β_3 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$\frac{\partial}{\partial \beta_0} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

$$\frac{\partial}{\partial \beta_1} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

$$\frac{\partial}{\partial \beta_2} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

⋮

$$\frac{\partial}{\partial \beta_n} \frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

A linear model in $k+1$ dimensions...

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

Has $k + 1$ parameters

And a cost function...

$$\frac{1}{2n} \sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

Compute $k + 1$
step sizes

Gradient Descent

Gradient Descent Algorithm

Step 1: Start with random values for $\beta_0, \beta_1, \beta_2$ and β_3

Step 2: Compute the partial derivative of the MSE w.r.t $\beta_0, \beta_1, \beta_2$ and β_3 - this is the slope at that point

Step 3: Calculate a step size that is proportional to the slope

$$step_size_{\beta_0} = \frac{\partial}{\partial \beta_0} MSE \times learning_rate$$

$$step_size_{\beta_1} = \frac{\partial}{\partial \beta_1} MSE \times learning_rate$$

$$step_size_{\beta_2} = \frac{\partial}{\partial \beta_2} MSE \times learning_rate$$

⋮

$$step_size_{\beta_k} = \frac{\partial}{\partial \beta_k} MSE \times learning_rate$$

A linear model in k+1 dimensions...

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for $\beta_0, \beta_1, \beta_2$ and β_3

Step 2: Compute the partial derivative of the SSR

at that point

proportional

Computing $k + 1$ partial derivatives isn't practical

Has $k + 1$

And a c

$$\sum_{i=0}^n (y_i - \beta_0 - \beta_1 \hat{x}_{1i} - \beta_2 \hat{x}_{2i} - \beta_3 \hat{x}_{3i} - \dots - \beta_k \hat{x}_{ki})^2$$

Compute k+ 1
step sizes

$$step_size_{\beta_2} = \frac{\partial}{\partial \beta_2} SSR \times learning_rate$$

⋮

$$step_size_{\beta_k} = \frac{\partial}{\partial \beta_k} SSR \times learning_rate$$

Lets use a Matrix

Multiple Regression

Linear Model in
 $k + 1$ Dimensions

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

$$\hat{y}_n = 1 \times \beta_0 + \hat{x}_{1n} \times \beta_1 + \hat{x}_{2n} \times \beta_2 + \hat{x}_{3n} \times \beta_3 + \dots + \hat{x}_{kn} \times \beta_k$$
$$\begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \cdot \\ \cdot \\ \cdot \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & \hat{x}_{10} & \hat{x}_{20} & \hat{x}_{30} & \cdot & \cdot & \cdot & \hat{x}_{k0} \\ 1 & \hat{x}_{11} & \hat{x}_{21} & \hat{x}_{31} & \cdot & \cdot & \cdot & \hat{x}_{k1} \\ 1 & \hat{x}_{12} & \hat{x}_{22} & \hat{x}_{32} & \cdot & \cdot & \cdot & \hat{x}_{k2} \\ 1 & \hat{x}_{13} & \hat{x}_{23} & \hat{x}_{33} & \cdot & \cdot & \cdot & \hat{x}_{k3} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & \hat{x}_{1n} & \hat{x}_{2n} & \hat{x}_{3n} & \cdot & \cdot & \cdot & \hat{x}_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \cdot \\ \cdot \\ \cdot \\ \beta_k \end{bmatrix}$$

- 1 dependent variable \hat{y}
- k independent variables $\hat{x}_1, \hat{x}_2, \hat{x}_3 \dots \hat{x}_k$
- $k + 1$ parameters - $\beta_0, \beta_1, \beta_2, \beta_3 \dots \beta_k$

Multiple Regression

Linear Model in
 $k + 1$ Dimensions

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

$$\hat{Y} = \hat{X}\beta$$

\hat{Y} and \hat{X} are matrices

$$\hat{Y} = \hat{X}\beta$$

$$\begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & \hat{x}_{10} & \hat{x}_{20} & \hat{x}_{30} & \cdot & \cdot & \cdot & \hat{x}_{k0} \\ 1 & \hat{x}_{11} & \hat{x}_{21} & \hat{x}_{21} & \cdot & \cdot & \cdot & \hat{x}_{k1} \\ 1 & \hat{x}_{12} & \hat{x}_{22} & \hat{x}_{22} & \cdot & \cdot & \cdot & \hat{x}_{k2} \\ 1 & \hat{x}_{13} & \hat{x}_{23} & \hat{x}_{23} & \cdot & \cdot & \cdot & \hat{x}_{k3} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & \hat{x}_{1n} & \hat{x}_{2n} & \hat{x}_{3n} & \cdot & \cdot & \cdot & \hat{x}_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}$$

- 1 dependent variable \hat{y}
- k independent variables $\hat{x}_1, \hat{x}_2, \hat{x}_3 \dots \hat{x}_k$
- $k + 1$ parameters - $\beta_0, \beta_1, \beta_2, \beta_3 \dots \beta_k$

Multiple Regression

Linear Model in
 $k + 1$ Dimensions

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

$$\hat{Y} = \hat{X}\beta$$

The Mean Squared Error (MSE):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Partial Derivative w.r.t β :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2$$

$$\begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \hat{Y} \\ \cdot \\ \cdot \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & \hat{x}_{10} & \hat{x}_{20} & \hat{x}_{30} & \cdot & \cdot & \cdot & \hat{x}_{k0} \\ 1 & \hat{x}_{11} & \hat{x}_{21} & \hat{x}_{21} & \cdot & \cdot & \cdot & \hat{x}_{k1} \\ 1 & \hat{x}_{12} & \hat{x}_{22} & \hat{x}_{22} & \cdot & \cdot & \cdot & \hat{x}_{k2} \\ 1 & \hat{x}_{13} & \hat{x}_{23} & \hat{X} & \cdot & \cdot & \cdot & \hat{x}_{k3} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & \hat{x}_{1n} & \hat{x}_{2n} & \hat{x}_{3n} & \cdot & \cdot & \cdot & \hat{x}_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta \\ \cdot \\ \cdot \\ \beta_k \end{bmatrix}$$

Lets compute this
matrix derivative

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{\partial}{\partial \beta} \frac{1}{2n} \left(\sqrt{(Y - X\beta)^T (Y - X\beta)} \right)^2 \\ &= \frac{1}{2n} \frac{\partial}{\partial \beta} (Y - X\beta)^T (Y - X\beta) \end{aligned}$$

$$\text{let } A = (Y - X\beta)$$

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \\ &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \end{aligned}$$

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{\partial}{\partial \beta} \frac{1}{2n} \left(\sqrt{(Y - X\beta)^T (Y - X\beta)} \right)^2 \\ &= \frac{1}{2n} \frac{\partial}{\partial \beta} (Y - X\beta)^T (Y - X\beta) \end{aligned}$$

$$\text{let } A = (Y - X\beta)$$

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \\ &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \end{aligned}$$

Euclidean norm of a matrix:
 $\|A\| = \sqrt{A^T A}$

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{\partial}{\partial \beta} \frac{1}{2n} \left(\sqrt{(Y - X\beta)^T (Y - X\beta)} \right)^2 \\ &= \frac{1}{2n} \frac{\partial}{\partial \beta} (Y - X\beta)^T (Y - X\beta) \end{aligned}$$

$$\text{let } A = (Y - X\beta)$$

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \\ &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \end{aligned}$$

Euclidean norm of a matrix:
 $\|A\| = \sqrt{A^T A}$

Chain Rule for Derivative

[See Tutorial on Differential Calculus](#)

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \\ &= \frac{2}{2n} A^T \frac{\partial}{\partial \beta} (Y - X\beta) \\ &= \frac{1}{n} A^T (-X) \\ &= \frac{1}{n} (Y - X\beta)^T (-X) \\ &= -\frac{1}{n} (Y - X\beta)^T (X) \end{aligned}$$

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \\ &= \frac{2}{2n} A^T \frac{\partial}{\partial \beta} (Y - X\beta) \\ &= \frac{1}{n} A^T (-X) \\ &= \frac{1}{n} (Y - X\beta)^T (-X) \\ &= -\frac{1}{n} (Y - X\beta)^T (X) \end{aligned}$$

Chain Rule for Derivative

[See Tutorial on Differential Calculus](#)

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \\ &= \frac{2}{2n} A^T \frac{\partial}{\partial \beta} (Y - X\beta) \\ &= \frac{1}{n} A^T (-X) \\ &= \frac{1}{n} (Y - X\beta)^T (-X) \\ &= -\frac{1}{n} (Y - X\beta)^T (X) \end{aligned}$$

Chain Rule for Derivative

[See Tutorial on Differential Calculus](#)

$$\frac{\partial}{\partial A} A^T A = 2A^T$$

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\begin{aligned} \frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 &= \frac{1}{2n} \frac{\partial}{\partial \beta} (A)^T (A) \frac{\partial}{\partial \beta} A \\ &= \frac{2}{2n} A^T \frac{\partial}{\partial \beta} (Y - X\beta) \\ &= \frac{1}{n} A^T (-X) \\ &= \frac{1}{n} (Y - X\beta)^T (-X) \\ &= -\frac{1}{n} (Y - X\beta)^T (X) \end{aligned}$$

Chain Rule for Derivative

[See Tutorial on Differential Calculus](#)

$$\frac{\partial}{\partial A} A^T A = 2A^T$$

$$\frac{\partial}{\partial \beta} (Y - X\beta) = -X$$

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 = -\frac{1}{n} (Y - X\beta)^T (X)$$

Gradient Vector: $1 \times (k + 1)$ row vector

because...

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Multiple Regression

Partial Derivative w.r.t β :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 = -\frac{1}{n} (Y - X\beta)^T (X)$$

Gradient Vector: $1 \times (k + 1)$ row vector

because...

Y is a $(n + 1) \times 1$ column vector

$X\beta$ is a $(n + 1) \times 1$ column vector

$(Y - X\beta)$ is a $(n + 1) \times 1$ column vector

$(Y - X\beta)^T$ is a $1 \times (n + 1)$ row vector

X is a $(n + 1) \times (k + 1)$ matrix

$(Y - X\beta)^T (X)$ is a $1 \times (k + 1)$ row vector

Multiple Regression

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Partial Derivative w.r.t β :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 = -\frac{1}{n} (Y - X\beta)^T (X)$$

Gradient Vector: $1 \times (k + 1)$ row vector

Transpose it to get the column vector

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 = \left(-\frac{1}{n} (Y - X\beta)^T (X)\right)^T$$

$$= -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Vector: $(k + 1) \times 1$ column vector

Multiple Regression

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

Partial Derivative w.r.t β :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 = -\frac{1}{n} (Y - X\beta)^T (X)$$

Gradient Vector: $1 \times (k + 1)$ row vector

Transpose it to get the column vector

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \|Y - X\beta\|^2 = \left(-\frac{1}{n} (Y - X\beta)^T (X)\right)^T$$

$$(AB)^T = B^T A^T$$

$$= -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Vector: $(k + 1) \times 1$ column vector

Multiple Regression

$$\hat{Y} = \hat{X}\beta$$

Linear Model in
 $k + 1$ Dimensions

The **Mean Squared Error (MSE)**:

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Cost Function

Partial Derivative w.r.t β :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Partial Derivative w.r.t β

Lets walk through gradient descent using this matrix representation of the Cost Function and its partial derivative (the gradient vector)

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

$$d_cost = -\frac{1}{n} X^T (Y - X\beta)$$

Step 3: Calculate a step size that is proportional to the slope

Step 4: Calculate new values for β_0 and β_1 by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

$$d_cost = -\frac{1}{n} X^T (Y - X\beta)$$

Step 3: Calculate a step size that is proportional to the slope

$$step_size = d_cost \times learning_rate$$

Step 4: Calculate new values for β by subtracting the step size

Step 5: Go to step 2 and repeat

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

$$d_cost = -\frac{1}{n} X^T (Y - X\beta)$$

Step 3: Calculate a step size that is proportional to the slope

$$step_size = d_cost \times learning_rate$$

Step 4: Calculate new values for β by subtracting the step size

$$\beta = \beta - step_size$$

Step 5: Go to step 2 and repeat

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

$$d_cost = -\frac{1}{n} X^T (Y - X\beta)$$

Step 3: Calculate a step size that is proportional to the slope

$$step_size = d_cost \times learning_rate$$

Step 4: Calculate new values for β by subtracting the step size

$$\beta = \beta - step_size$$

Step 5: Go to step 2 and repeat

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

Gradient Descent continues in this manner until the step size is close to zero or a fixed number of iterations

Step 4: Calculate new values for β by subtracting the step size

$$\beta = \beta - step_size$$

Step 5: Go to step 2 and repeat

A linear model in $k + 1$ dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t β):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Matrix algebra allows us to compute gradients and step sizes in a single computation

Gradient Descent

Gradient Descent: Basic Concept

Step 1: Start with random values for β

Step 2: Compute the partial derivative of the cost function w.r.t β

Gradient Descent continues in this manner until the step size is close to zero or a fixed number of iterations

Step 4: Calculate new values for β by subtracting the step size

$$\beta = \beta - step_size$$

Step 5: Go to step 2 and repeat

Related Tutorials & Textbooks

Multiple Regression ↗

Multiple regression extends the two dimensional linear model introduced in Simple Linear Regression to $k + 1$ dimensions with one dependent variable, k independent variables and $k+1$ parameters.

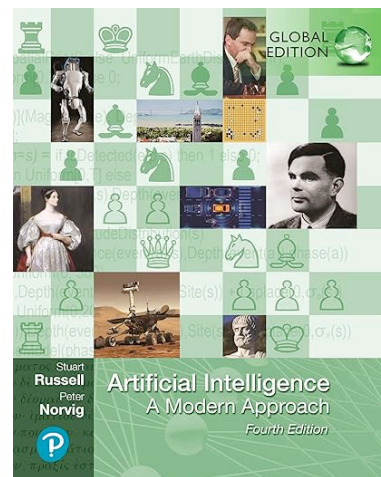
Gradient Descent for Simple Linear Regression ↗

Gradient Descent algorithm for multiple regression and how it can be used to optimize $k + 1$ parameters for a Linear model in multiple dimensions.

Logistic Regression ↗

An introduction to Logistic Regression. A Logistic Regression model use used to predict a binary value (the dependent variable) for one or more independent variables using a threshold to classify a probability.

Recommended Textbooks



Artificial Intelligence: A Modern Approach

by Peter Norvig, Stuart Russell

For a complete list of tutorials see:

<https://arrsingh.com/ai-tutorials>