



# Introduction to Machine Learning

## Week 9

Olivier JAYLET

School of Information Technology and Engineering

# Introduction

A machine learning algorithm can be viewed as an optimization program.

Today, we will have a look at a common algorithm used to find the parameters that minimize a known cost function  $f(\cdot)$  : **the gradient descent algorithm.**

- **vanilla** version of the gradient descent algorithm.
- stochastic gradient
- mini-batch gradient

Let us consider a very general model:

$$y = f(X) + \varepsilon,$$

where  $y$  is a variable to predict (i.e. target variable or response variable),  $f(\cdot)$  is an unknown model,  $X$  is a set of  $p$  predictors (i.e. features, or inputs, or explanatory variables) and  $\varepsilon$  is an error term.

For the example, we assume that the response variable is linearly dependent on the set of explanatory variables:

$$y = X\beta + \varepsilon.$$

We do not know the true generating data process and only observe some realizations of  $y$  and  $X$  for  $n$  observations.

When estimating a linear regression, we assume that the error term is normally distributed with zero mean and variance  $\sigma^2$ .

The vector of coefficients  $\beta$  can be estimated with **Ordinary Least Squares** (OLS).

The problem is to estimate the coefficients of vector  $\beta$  which **minimize an objective function**:

$$\arg \min_{\beta} \sum_{i=1}^n \mathcal{L}(y_i, f(X_i)),$$

Here, with OLS, an analytical solution exists:

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

# Time complexity

Time complexity is a measure of the time used by an algorithm, expressed as a function of the size of the input. Time counts the number of calculation steps before arriving at a result.

# OLS Complexity

Each steps of the solution are :

- $X^T X : O(nm^2)$ ,
- $(X^T X)^{-1} : O(m^3)$ ,
- $X^T : O(m^2n)$ ,
- Final multiplication :  $O(mn)$ .

## Dominant Term:

- For small  $n$  (small dataset):  $O(m^3)$  (matrix inversion dominates).
- For large  $n$  (large dataset):  $O(nm^2)$  (matrix multiplication dominates).

## Effect of large dataset : $n \gg m$

When the number of samples (rows)  $n$  is much greater than the number of features  $m$ :

- The dominant time complexity is  $O(nm^2)$ , which comes from computing  $X^T X$ .
- This dependency on  $n$  makes OLS expensive for **large datasets**, as increasing  $n$  leads to a quadratic increase in computational cost.

## Effect of High Dimensionality : $m \gg n$

When the number of features (columns  $m$ ) is much greater than the number of samples (rows  $n$ ):

- The dominant time complexity is  $O(m^3)$ , which comes from matrix inversion of  $(X^T X)$ , an  $m \times m$  matrix.
- This makes dimensionality (number of features) a critical factor for computational feasibility.

# Need of Numerical Solutions

Numerical optimization techniques are often used instead of analytical solutions :

## **Expensive Calculations of Analytical Solutions:**

- Computing the closed-form solution involving operations like matrix multiplication and matrix inversion can be computationally expensive for large datasets or high-dimensional data.
- The complexity of these operations makes them infeasible for large-scale problems.

# Need of Numerical Solutions

## Applicability to Non-Linear Models:

- Analytical solutions are typically restricted to linear models. For non-linear models or cases where the loss function is not quadratic, numerical optimization methods, are required to minimize the loss function.
- These techniques are more flexible and can handle the iterative process needed to converge on a solution for non-linear and complex models.

Thus, numerical solutions provide scalability and adaptability for solving real-world machine learning problems efficiently.

# Vanilla gradient descent

Let us illustrate this with a simple example.

Let us consider the following function:

$$f(x) = (x + 3) \times (x - 2)^2 \times (x + 1).$$

The global minimum of that function is reached in

$$x = -1 - \sqrt{\frac{3}{2}}.$$

Let us generate some values from this process, for  $x \in [-3, 3]$ .

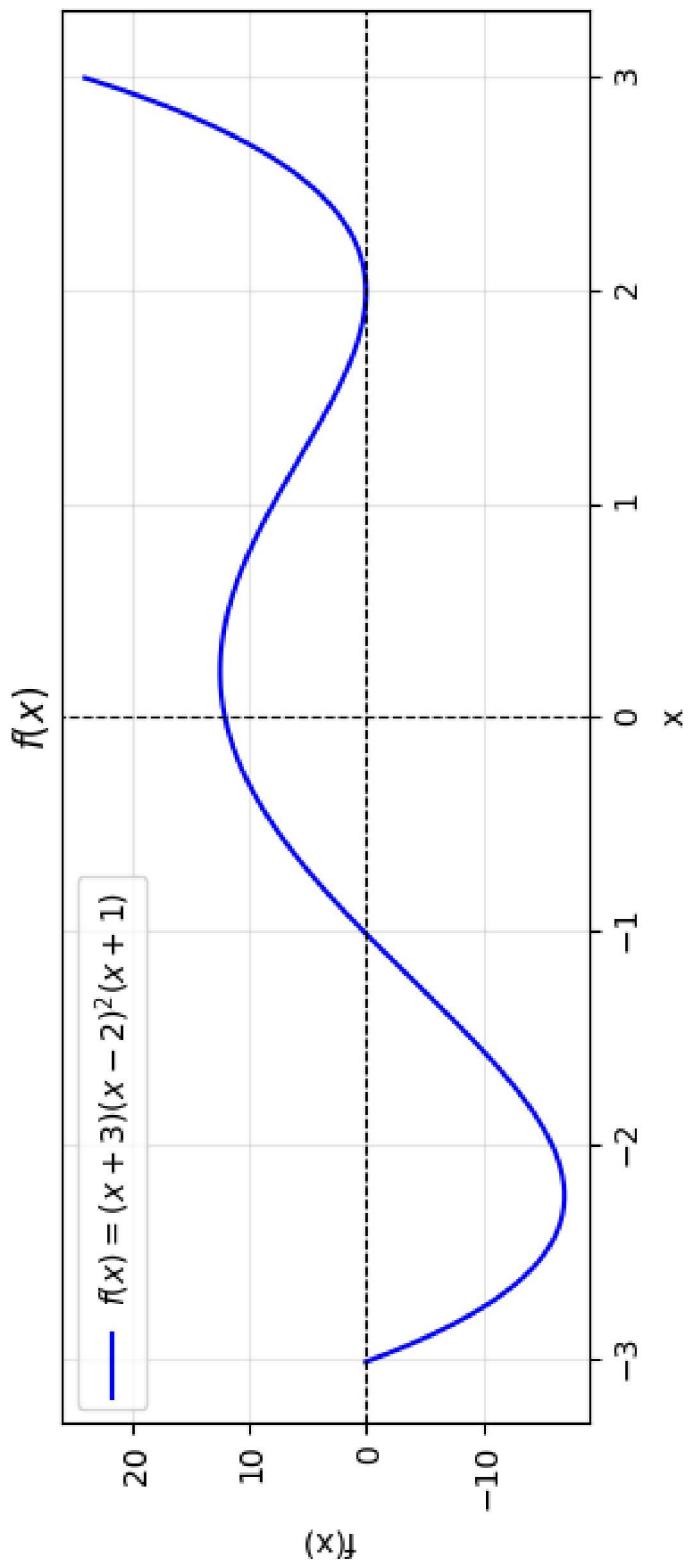


Figure: Function with a single input: a more complex function

## Vanilla gradient descent

If we want to minimize this function using gradient descent, we can proceed as follows :

# Initialization

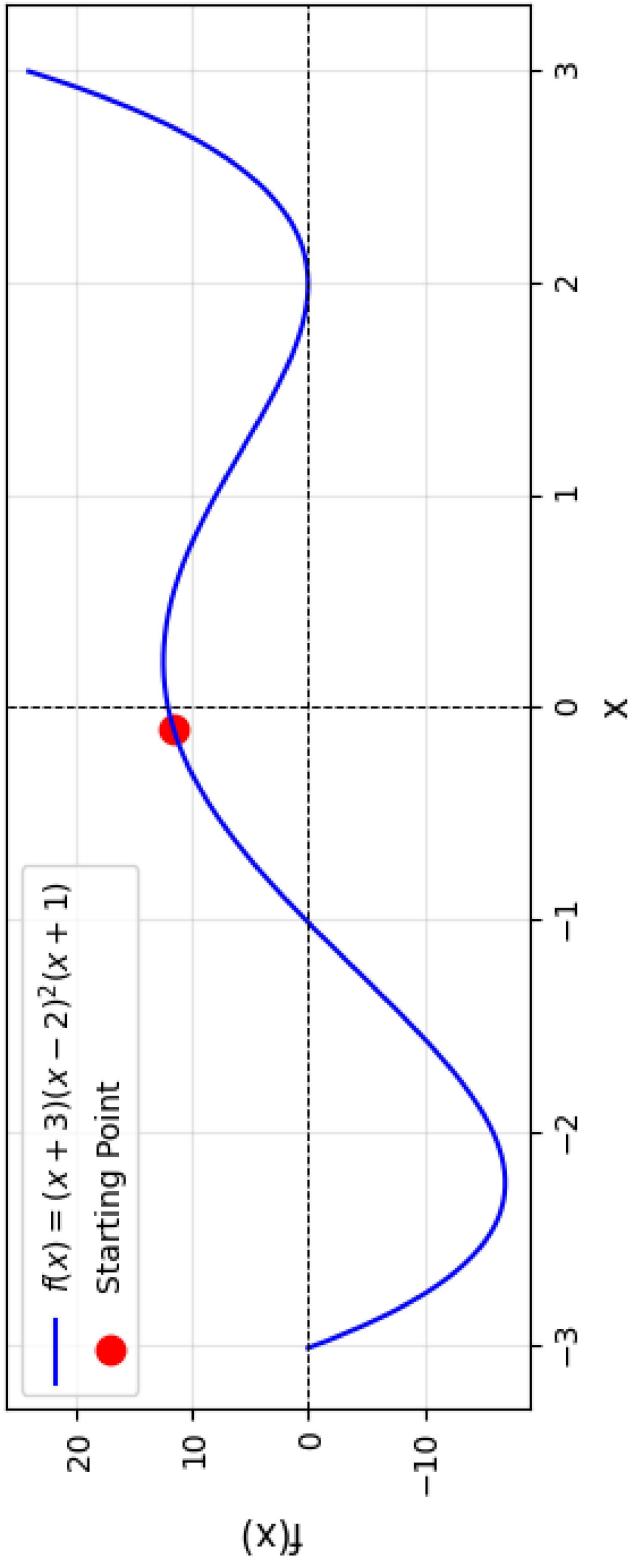


Figure: Initialize random point on the function

# Vanilla gradient descent

Then, from that point, we need to decide on two things so as to reduce the objective function:

- 1 in which direction to go next (left or right)
- 2 and how far we want to go.

# Direction

To decide the direction, we can compute the derivative of the function at this specific point of interest. The slope of the derivative will guide us:

- **if it is positive:** we need to shift to the left
- **if it is negative:** we need to shift to the right.

# Direction

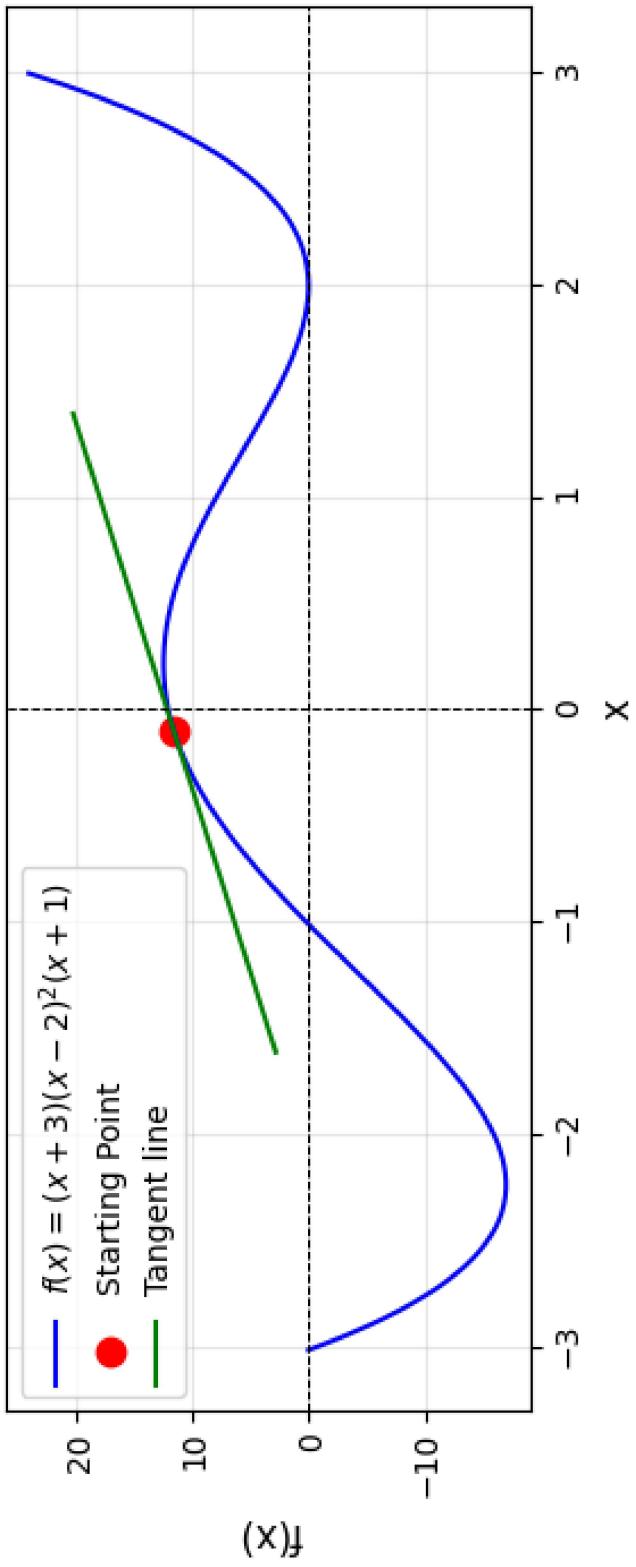


Figure: The derivative of  $f(x) = (x+3) \times (x-2)^2 \times (x+1)$  at  $x = -0.5$  is 12.5, i.e. the slope of the tangent is positive.

- Here, the slope is positive. We thus need to go left.
- We still need to decide how far we want to go, *i.e.*, we must decide the size of the step we will take.
- This step is called the **learning rate**:
  - On the one hand, if this learning rate is too small, we increase the risk of ending up in a local minimum.
  - On the other hand, if we pick a too large value for the learning rate, we face a risk of overshooting the minimum and keeping bouncing around a (local) minimum forever.

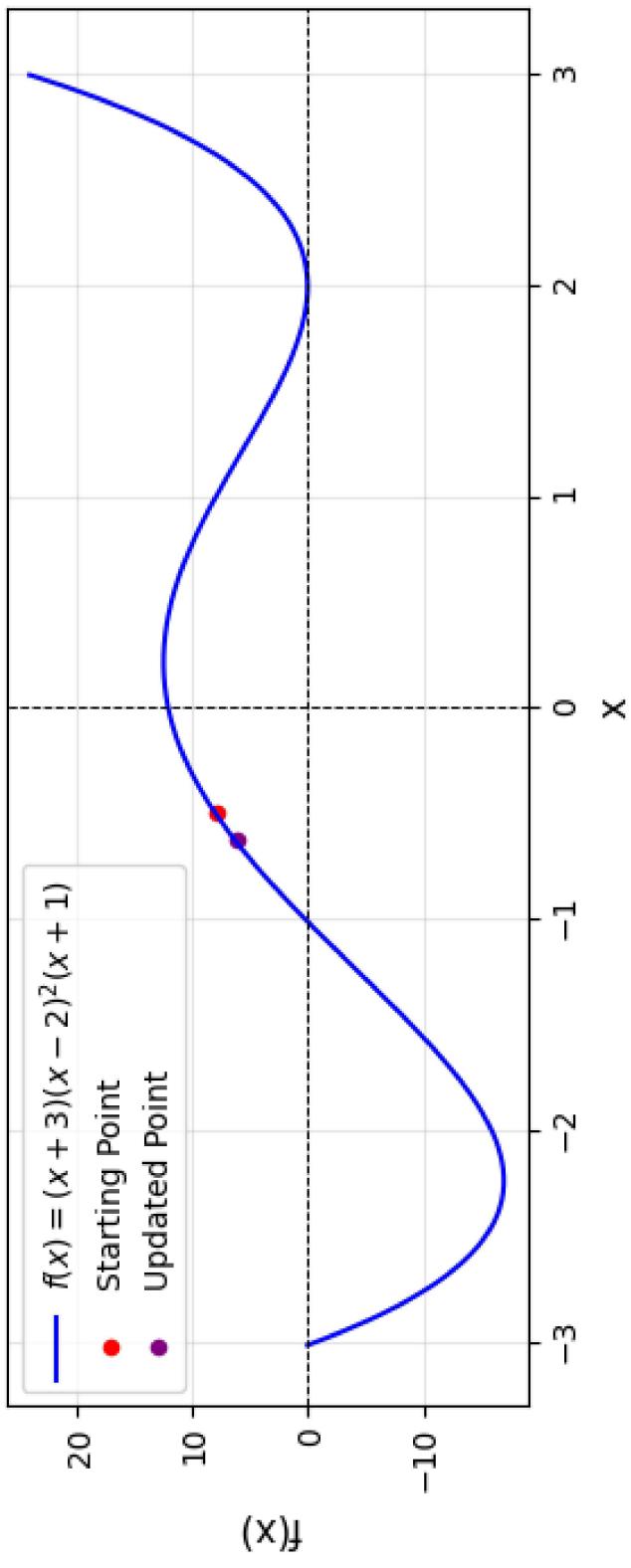


Figure: Create a new point

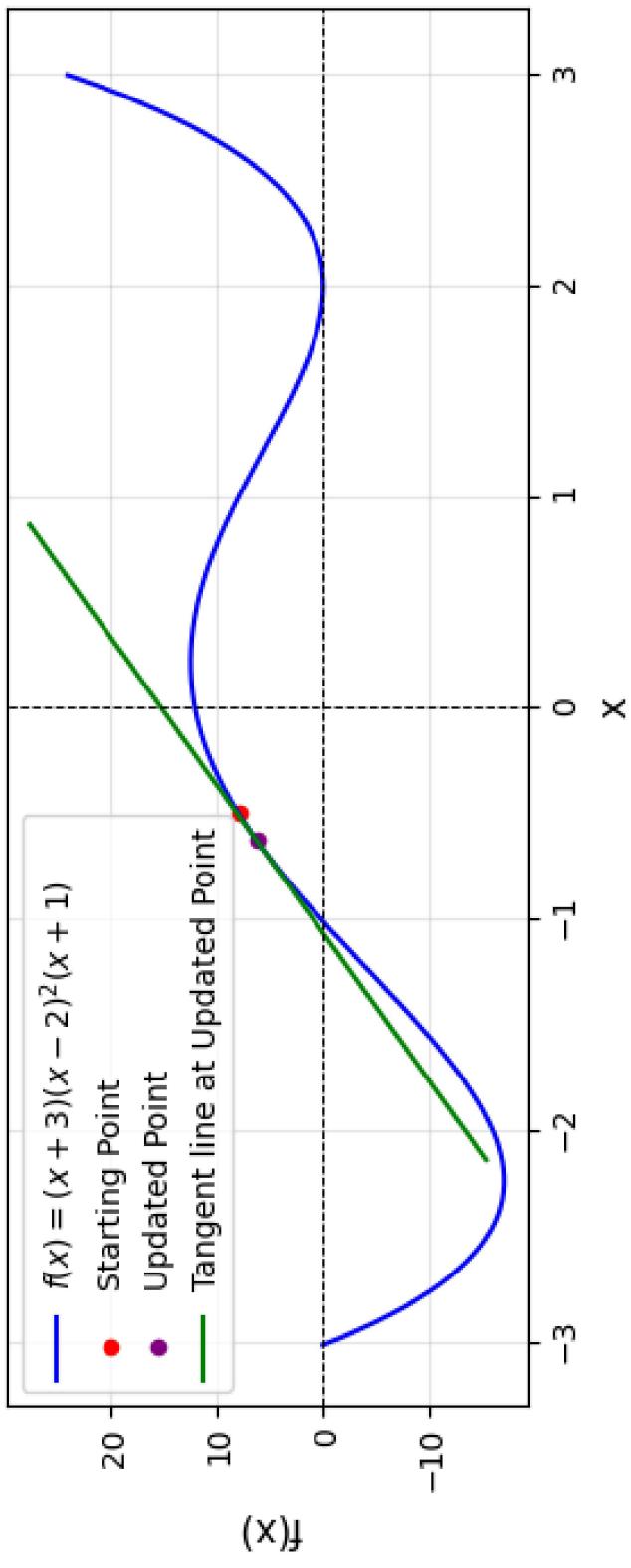


Figure: Update the derivative

## Iterative Procedure

- Then, we can repeat the procedure multiple times through a loop.
- We will update our parameter from one iteration to the other and will stop either:
  - When a maximum number of iterations is reached
  - When the improvement (reduction in the objective function from one step to the next) is too small (below a threshold we will call **tolerance**).

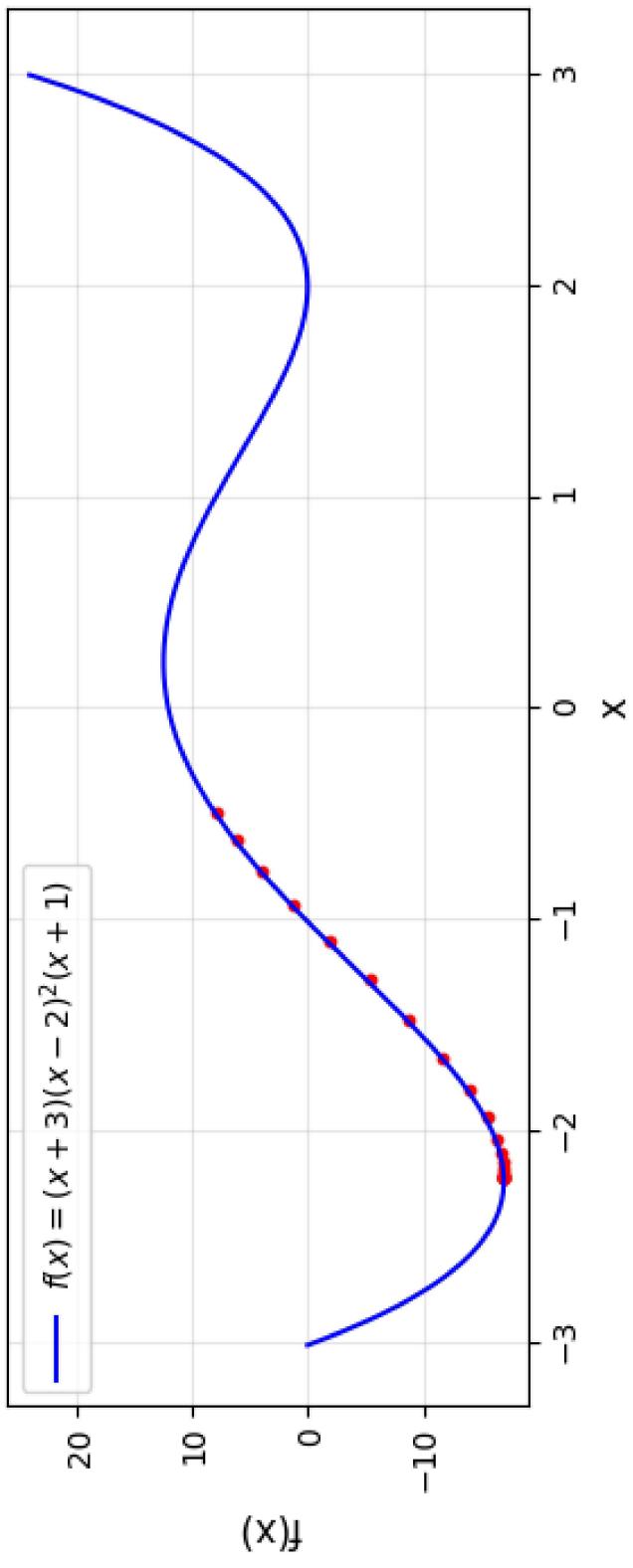


Figure: Process repeated into a loop

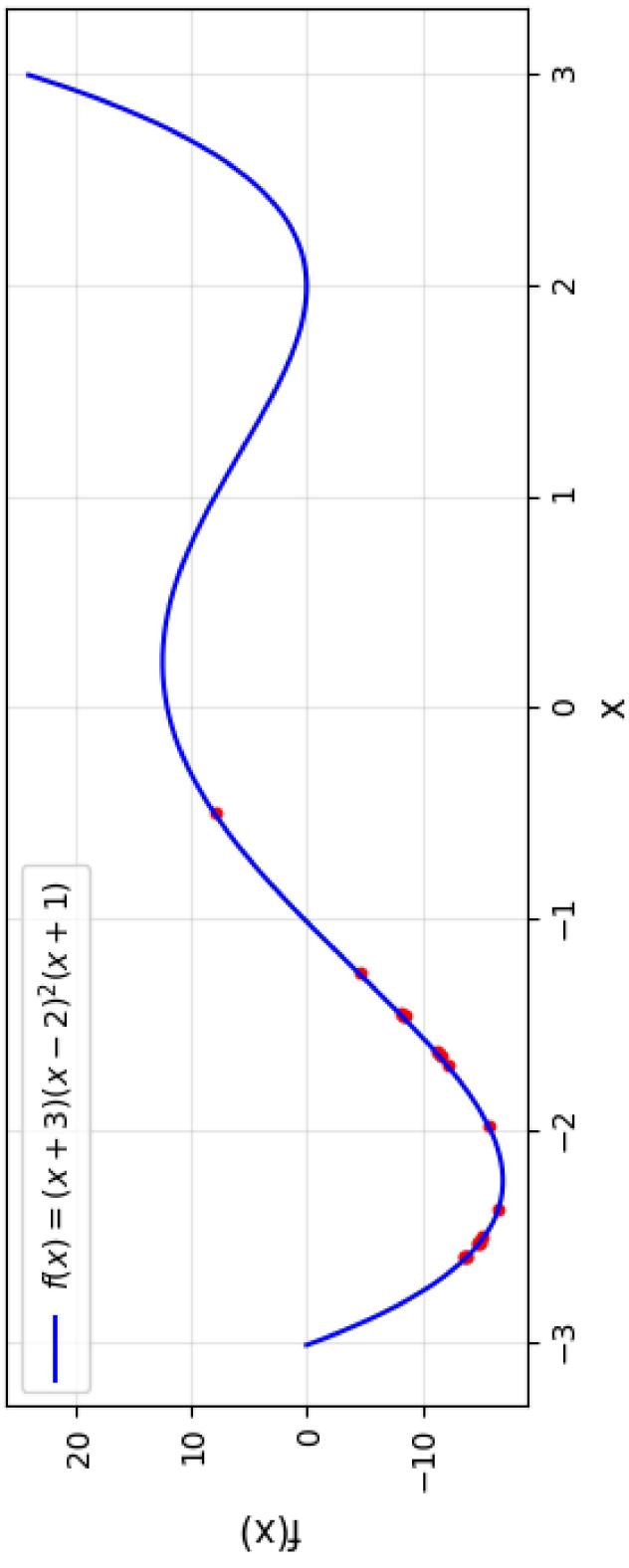


Figure: Too high learning rate (1)

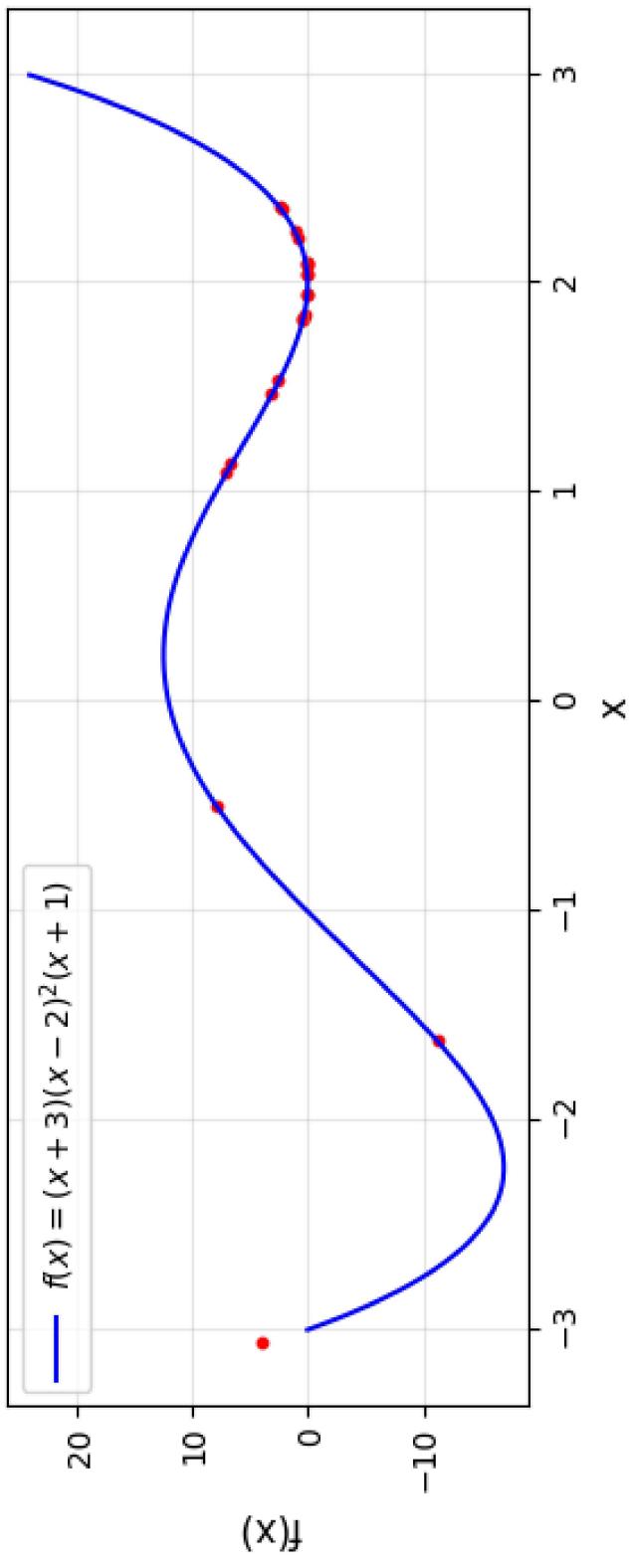


Figure: Too high learning rate (2)

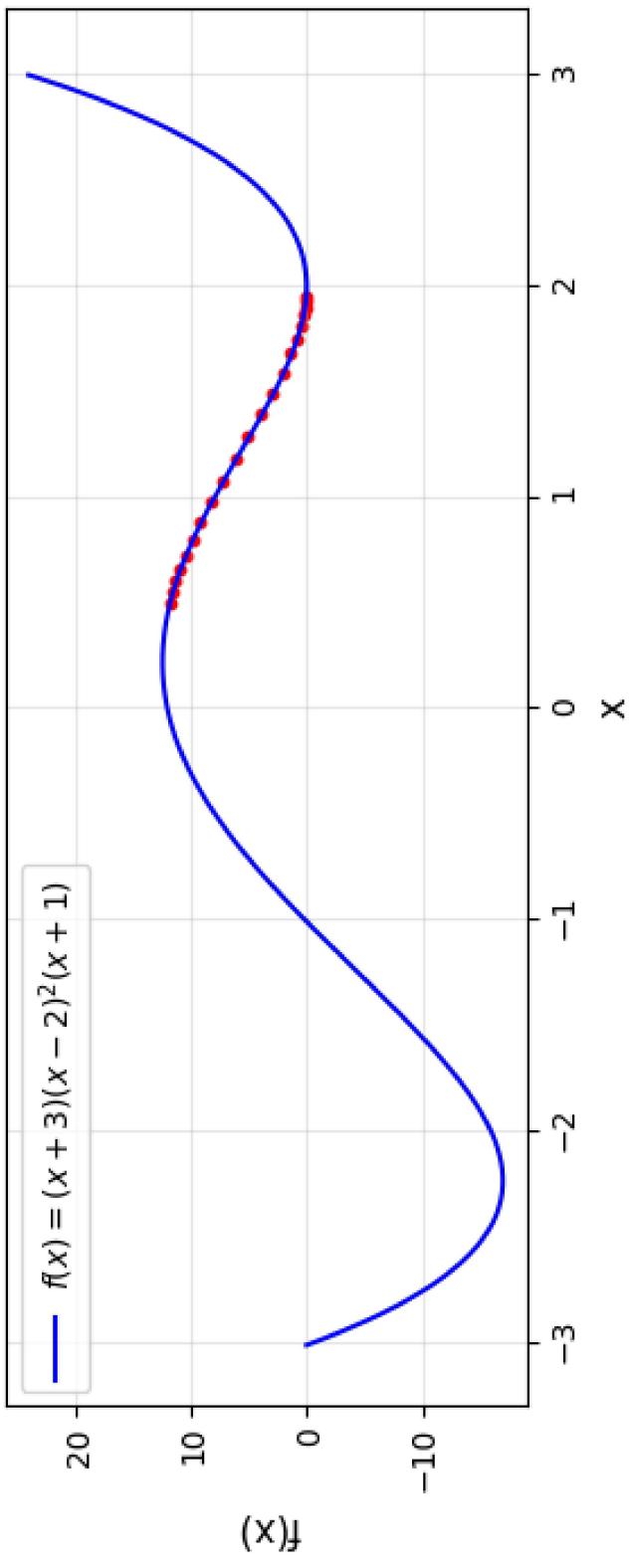


Figure: Dependence on the initial values

# Higher Dimensions Optimization Problems

Let us consider the following data generating process:

$$f(x_1, x_2) = x_1^2 + x_2^2$$

Surface of  $f(x_1, x_2)$  with Initial Point

Initial Point

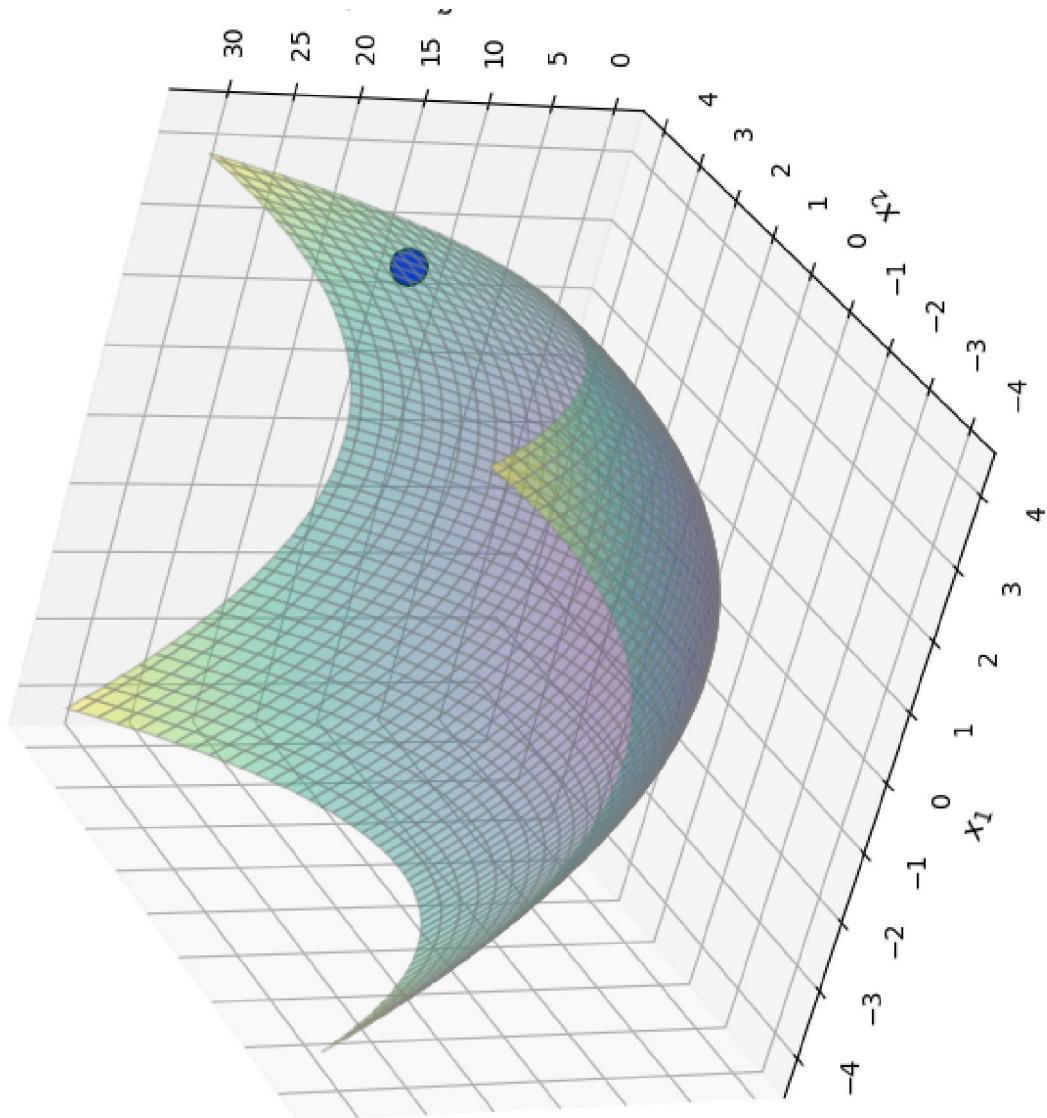


Figure: Finding a local minima

# Intuition

From that point, we need to decide:

- **the direction** to go to
- **and the magnitude of the step** to take in that direction.

# Computing the gradient

The direction is obtained by computing the first derivative of the objective function  $f(\cdot)$  with respect to each argument  $x_1$  and  $x_2$ , at point  $\theta$ . In other words, we need to evaluate the gradient of the function at point  $\theta$ .

$$\nabla f(\theta) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\theta) \\ \frac{\partial f}{\partial x_2}(\theta) \end{bmatrix}$$

The values will give us the steepest ascent.

The updated value of the parameters after the end of the  $t$  th step will be:

$$\begin{bmatrix} x_1^{(t+1)} \\ x_2^{(t+1)} \end{bmatrix} = \begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \end{bmatrix} - \eta \cdot \begin{bmatrix} \frac{\partial f}{\partial x_1}(x_1^{(t)}, x_2^{(t)}) \\ \frac{\partial f}{\partial x_2}(x_1^{(t)}, x_2^{(t)}) \end{bmatrix},$$

where

$$\begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \end{bmatrix}$$

is the current value of the vector of parameters,

$\eta \in \mathbb{R}^+$  is the learning rate, and

$$\begin{bmatrix} \frac{\partial f}{\partial x_1}(x_1^{(t)}, x_2^{(t)}) \\ \frac{\partial f}{\partial x_2}(x_1^{(t)}, x_2^{(t)}) \end{bmatrix}$$
 is the gradient of the function at point  $\theta = (x_1^{(t)}, x_2^{(t)})$ .

In a more general context, the update rule becomes:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \cdot \nabla \mathcal{L}(\theta^{(t)}),$$

where:

- $\eta$  is the learning rate,
- $\nabla \mathcal{L}(\theta^{(t)})$  is the gradient of the loss function  $\mathcal{L}$  at  $\theta^{(t)}$ , and
- $\theta$  denotes the vector of parameters being optimized.

### Surface of $f(x_1, x_2)$ with Gradient Descent Path

Gradient Descent Path

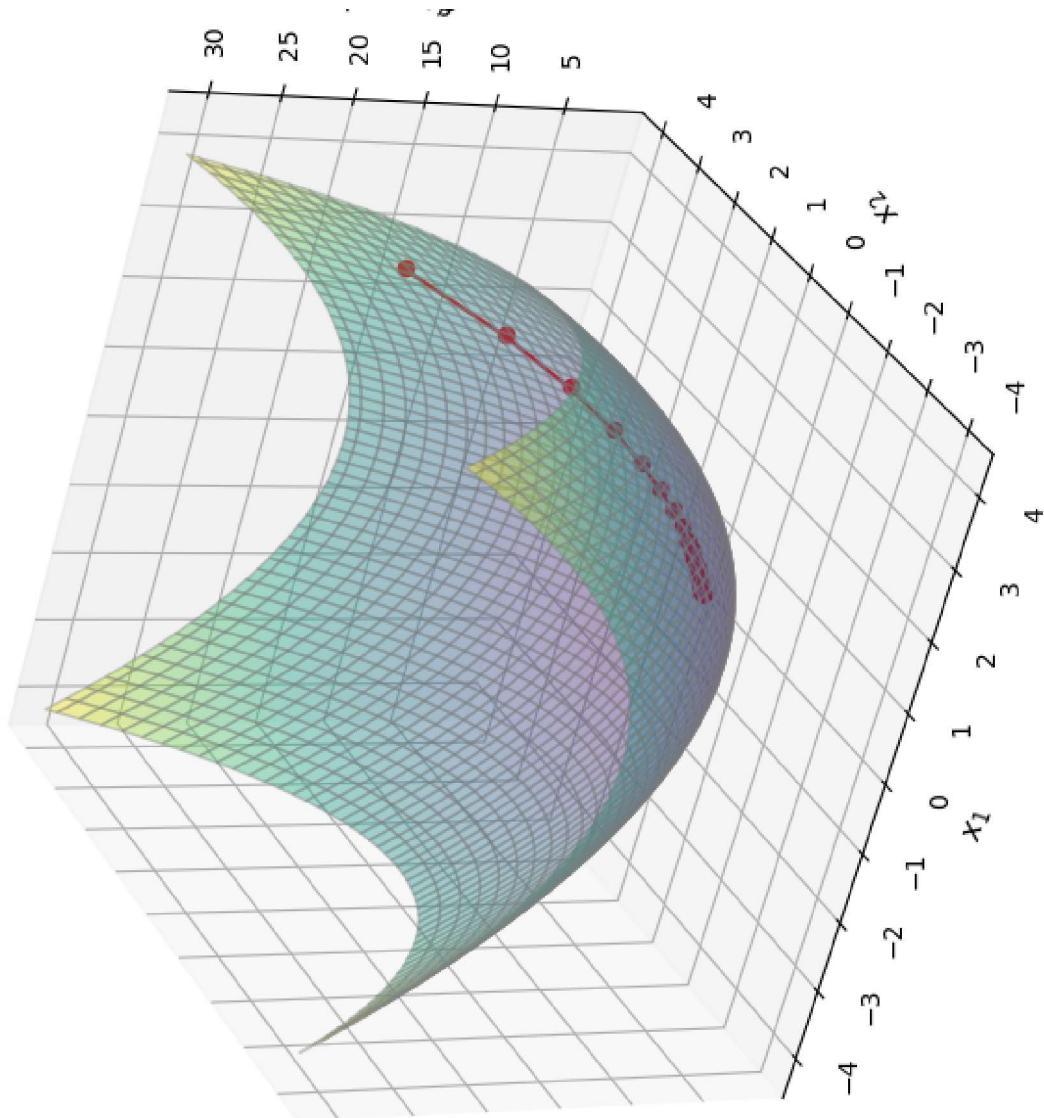


Figure: Finding a local minima

The gradient descent algorithm can be written as:

---

### Algorithm Vanilla Gradient Descent

---

```
1: Input:  $\mathcal{L}(x)$ ,  $\eta$ ,  $\theta(0)$ ,  $T$ ,  $\epsilon$ 
2: Initialize:  $t \leftarrow 0$ 
3: repeat
4:    $g \leftarrow \nabla \mathcal{L}(\theta(t))$ 
5:    $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot g$ 
6:    $t \leftarrow t + 1$ 
7: until  $\|g\| < \epsilon$  OR  $t > T$ 
8: Return:  $\theta^{(t)}$ 
```

---

where,  $\mathcal{L}(x)$  is a loss function,  $\eta$  is a learning rate,  $\theta(0)$  is the set of initial parameters,  $T$  is the maximum number of iterations, and  $\epsilon$  is the tolerance.

So far, we have estimated the  $\rho$  parameters that minimize an objective function  $\mathcal{L}(\theta)$ , where  $\theta$  is a vector of the  $\rho$  parameters to be estimated.

And the gradient descent algorithm updates  $\theta$  using the following rule:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \cdot \nabla \mathcal{L}(\theta^{(t)})$$

In the previous example, to compute the gradient of the objective function  $\mathcal{L}$ , we have used the **whole dataset**.

The learning rate  $\eta$  was a constant.

Now, we will consider different ways of updating the parameters.

- Focus on the frequency of updates and on the samples used to update the parameters.
- Make the learning rate vary along the iteration process.

## Algorithm Stochastic Gradient Descent (SGD)

---

```
1: Input:  $\mathcal{L}(x^{(i)}, y^{(i)})$ ,  $\eta$ ,  $\theta^{(0)}$ ,  $T$ ,  $\epsilon$ 
2: Initialize:  $t \leftarrow 0$ 
3: repeat
4:   Randomly shuffle the dataset
5:   for each example  $(x^{(i)}, y^{(i)})$  in the dataset do
6:      $g \leftarrow \nabla \mathcal{L}(\theta^{(t)}; x^{(i)}, y^{(i)})$ 
7:      $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot g$ 
8:      $t \leftarrow t + 1$ 
9:   end for
10:  until  $\|g\| < \epsilon$  OR  $t > T$ 
11: Return:  $\theta^{(t)}$ 
```

---

where,  $\mathcal{L}(x^{(i)}, y^{(i)})$  is the loss for a single example,  $\eta$  is the learning rate,  $\theta^{(0)}$  are the initial parameters,  $T$  is the maximum number of iterations, and  $\epsilon$  is the tolerance.

---

## Algorithm Mini-Batch Gradient Descent

---

```
1: Input:  $\mathcal{L}(X^{(b)}, Y^{(b)})$ ,  $\eta$ ,  $\theta^{(0)}$ ,  $T$ ,  $\epsilon$ , batch size  $B$ 
2: Initialize:  $t \leftarrow 0$ 
3: repeat
4:   Randomly shuffle the dataset
5:   Partition the dataset into mini-batches of size  $B$ 
6:   for each mini-batch  $(X^{(b)}, Y^{(b)})$  do
7:      $g \leftarrow \nabla \mathcal{L}(\theta^{(t)}; X^{(b)}, Y^{(b)})$ 
8:      $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot g$ 
9:      $t \leftarrow t + 1$ 
10:  end for
11: until  $\|g\| < \epsilon$  OR  $t > T$ 
12: Return:  $\theta^{(t)}$ 
```

---

where,  $\mathcal{L}(X^{(b)}, Y^{(b)})$  is the loss for a mini-batch,  $\eta$  is the learning rate,  $\theta^{(0)}$  are the initial parameters,  $T$  is the maximum number of iterations,  $\epsilon$  is the tolerance, and  $B$  is the mini-batch size.

# Linear Regression with OLS (sklearn)

```
1 from sklearn.linear_model import LinearRegression  
  
1 ols_model = LinearRegression()  
2 ols_model.fit(X_train, y_train)  
  
1 print("OLS Coefficients:", ols_model.coef_)  
2 print("OLS Intercept:", ols_model.intercept_)
```

# Linear Regression with Vanilla Gradient Descent (sklearn)

```
1 from sklearn.linear_model import SGDRegressor
```

```
1 vanilla_gd_model = SGDRegressor(           % Maximum number of iterations
2     max_iter=1000,                         % Tolerance
3     tol=1e-3,                             % Constant learning rate
4     learning_rate='constant',             % Initial learning rate value
5     eta0=0.01,                            %
6     random_state=42,                      %
7     penalty=None,                         % Disable shuffling
8     shuffle=False)                       %
9 )
```

```
1 vanilla_gd_model.fit(X_train_scaled, y_train)
```

# Linear Regression with Stochastic Gradient Descent (sklearn)

```
1 from sklearn.linear_model import SGDRegressor
```

```
1 sgd_model = SGDRegressor(                                % Maximum number of iterations  
2     max_iter=1000,                                     % Tolerance  
3     tol=1e-3,                                         % Constant learning rate  
4     learning_rate='constant',                         % Initial learning rate value  
5     eta0=0.01,                                         % No regularization  
6     random_state=42,                                    % shuffle is True by default  
7     penalty=None,                                     % shuffle is True by default  
8     shuffle=True  
9 )  
10 sgd_model.fit(X_train, Y_train)
```

```
1 print("SGD Coefficients:", sgd_model.coef_)  
2 print("SGD Intercept:", sgd_model.intercept_)
```

# Mini-Batch Gradient Descent (sklearn) - Setup

```
1 from sklearn.linear_model import SGDRegressor  
2 import numpy as np  
  
3  
4 batch_size = 32          % Mini-batch size  
5 n_epochs = 100           % Number of epochs  
6  
7 mini_batch_gd_model = SGDRegressor(  
8     learning_rate='constant',  
9     eta0=0.01,  
10    random_state=42  
11 )
```

# Mini-Batch Gradient Descent (sklearn) - Training Loop

```
1 for epoch in range(n_epochs):
2     # Shuffle data at the start of each epoch
3     indices = np.random.permutation(len(X_train))
4     X_shuffled = X_train[indices]
5     y_shuffled = y_train[indices]
6
7     # Iterate over mini-batches
8     for i in range(0, len(X_train), batch_size):
9         X_batch = X_shuffled[i:i + batch_size]
10        y_batch = y_shuffled[i:i + batch_size]
11
12        # Update model parameters with the current mini-
13        batch
mini_batch_gd_model.partial_fit(X_batch, y_batch)
```