

Basic Machine Learning

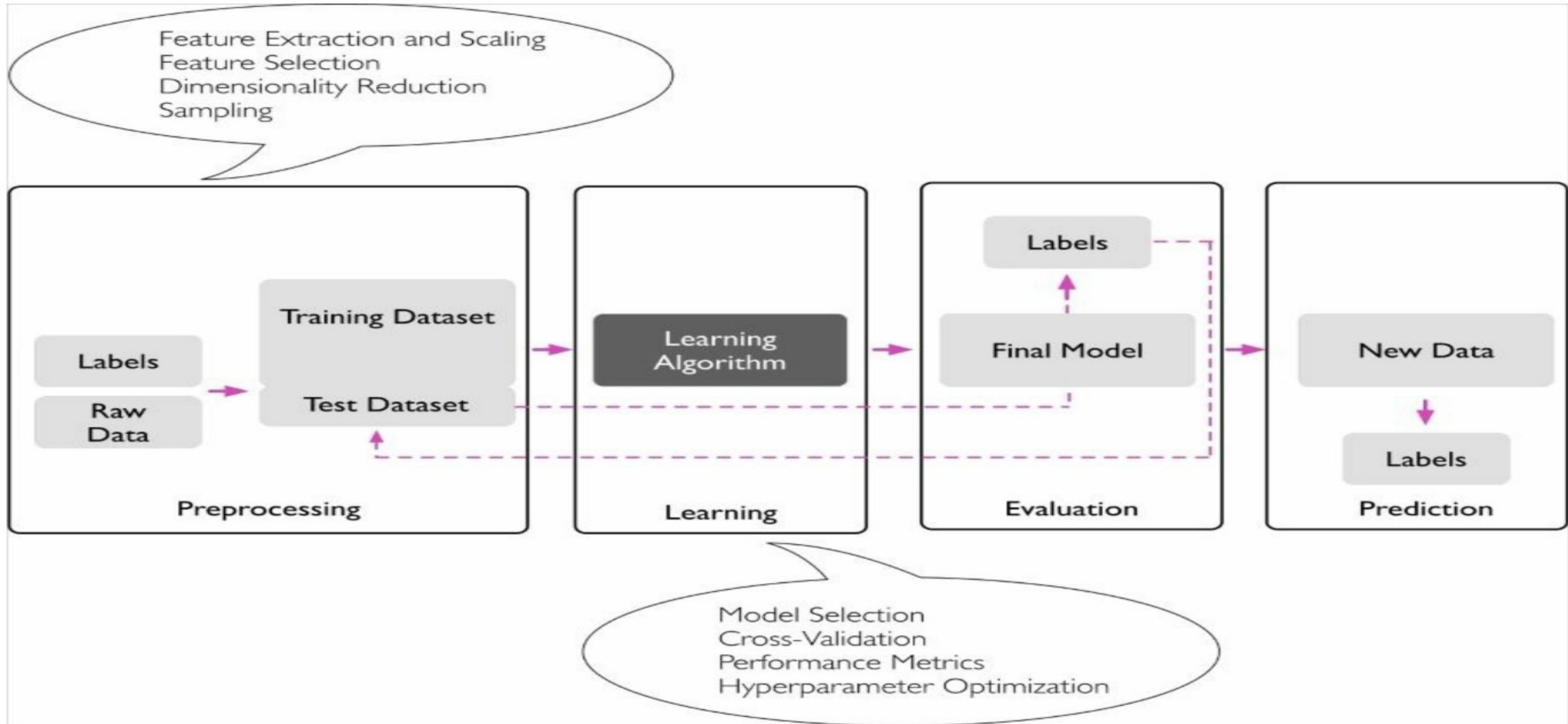
Understanding data, basic terminology, notations model building process

Kyi Thar
kyi.thar@pipe.bz

Introduction to the basic terminology and notations

Machine Learning Pipeline

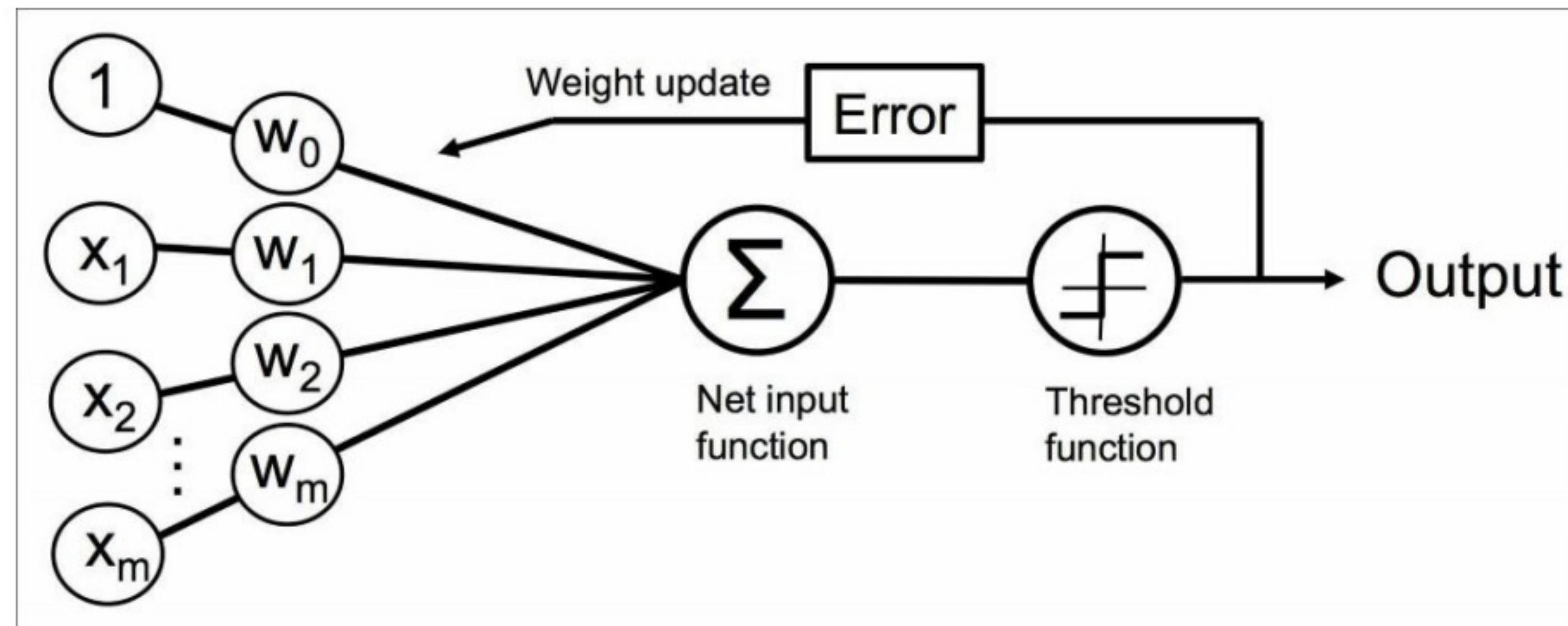
2



Introduction to the basic terminology and notations

Model

3



$$\Delta w_j = \eta(-1 - -1)x_j^{(i)} = 0$$

$$\Delta w_j = \eta(1 - 1)x_j^{(i)} = 0$$

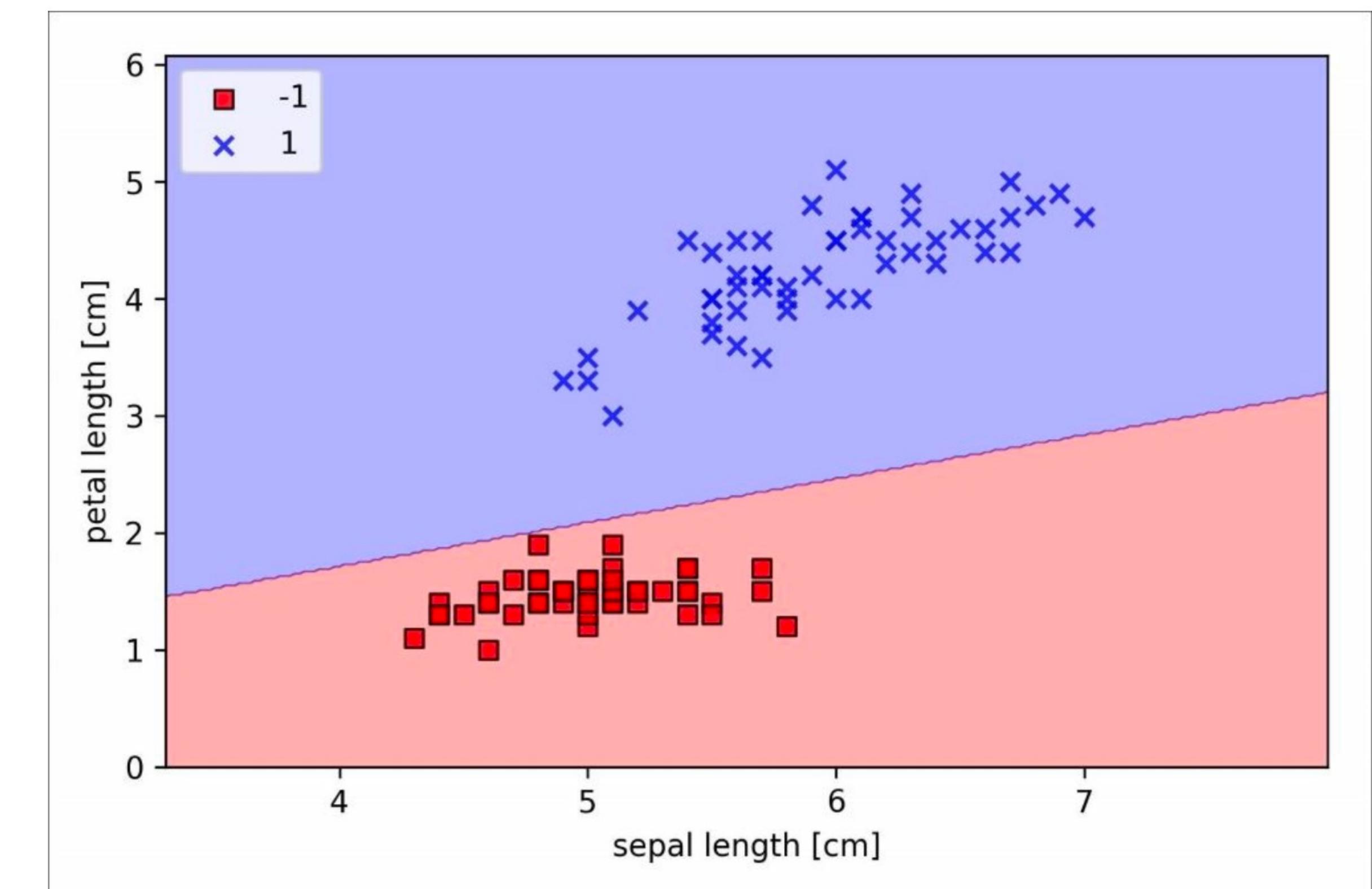
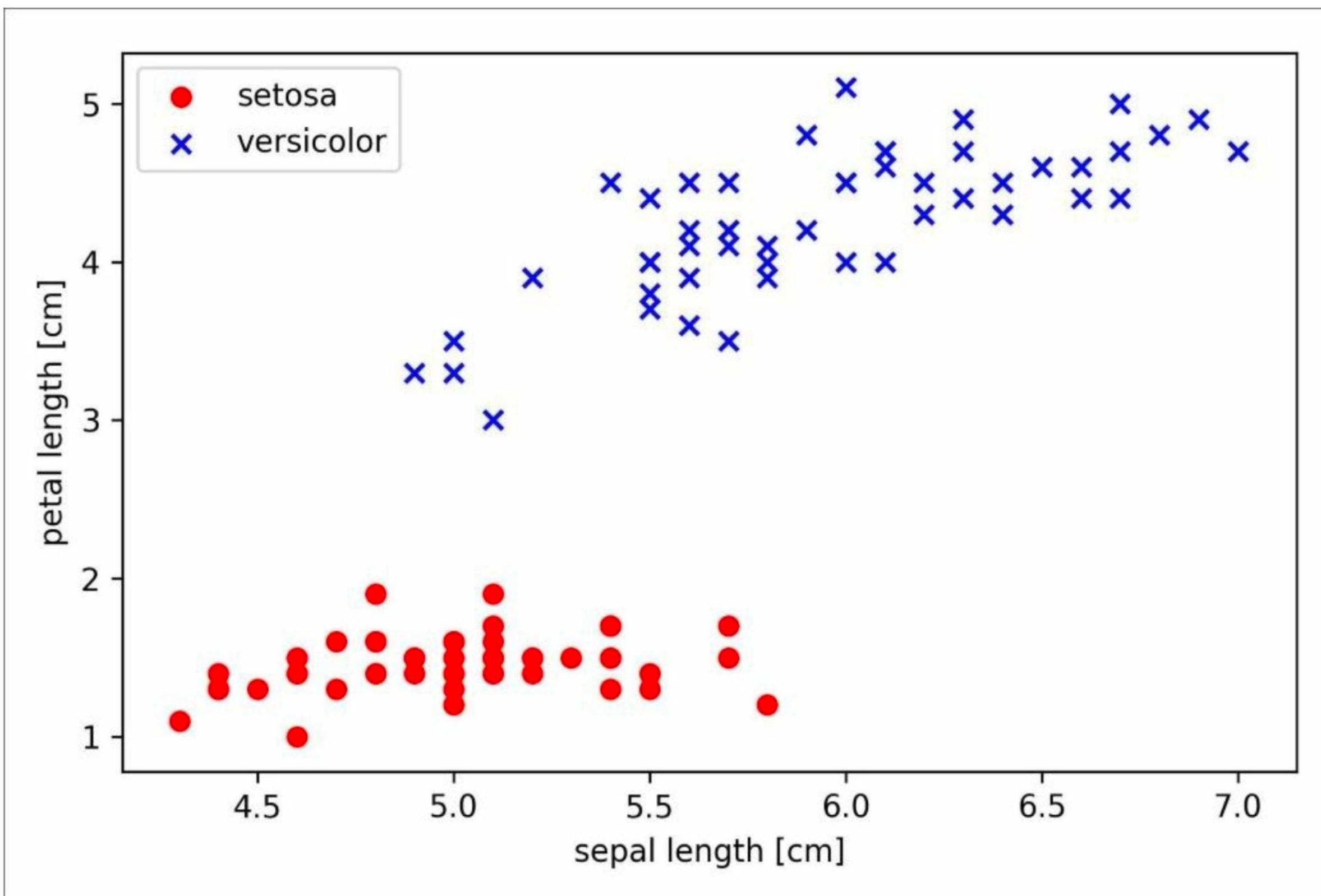
$$\Delta w_j = \eta(1 - -1)x_j^{(i)} = \eta(2)x_j^{(i)}$$

$$\Delta w_j = \eta(-1 - 1)x_j^{(i)} = \eta(-2)x_j^{(i)}$$

Introduction to the basic terminology and notations

Model

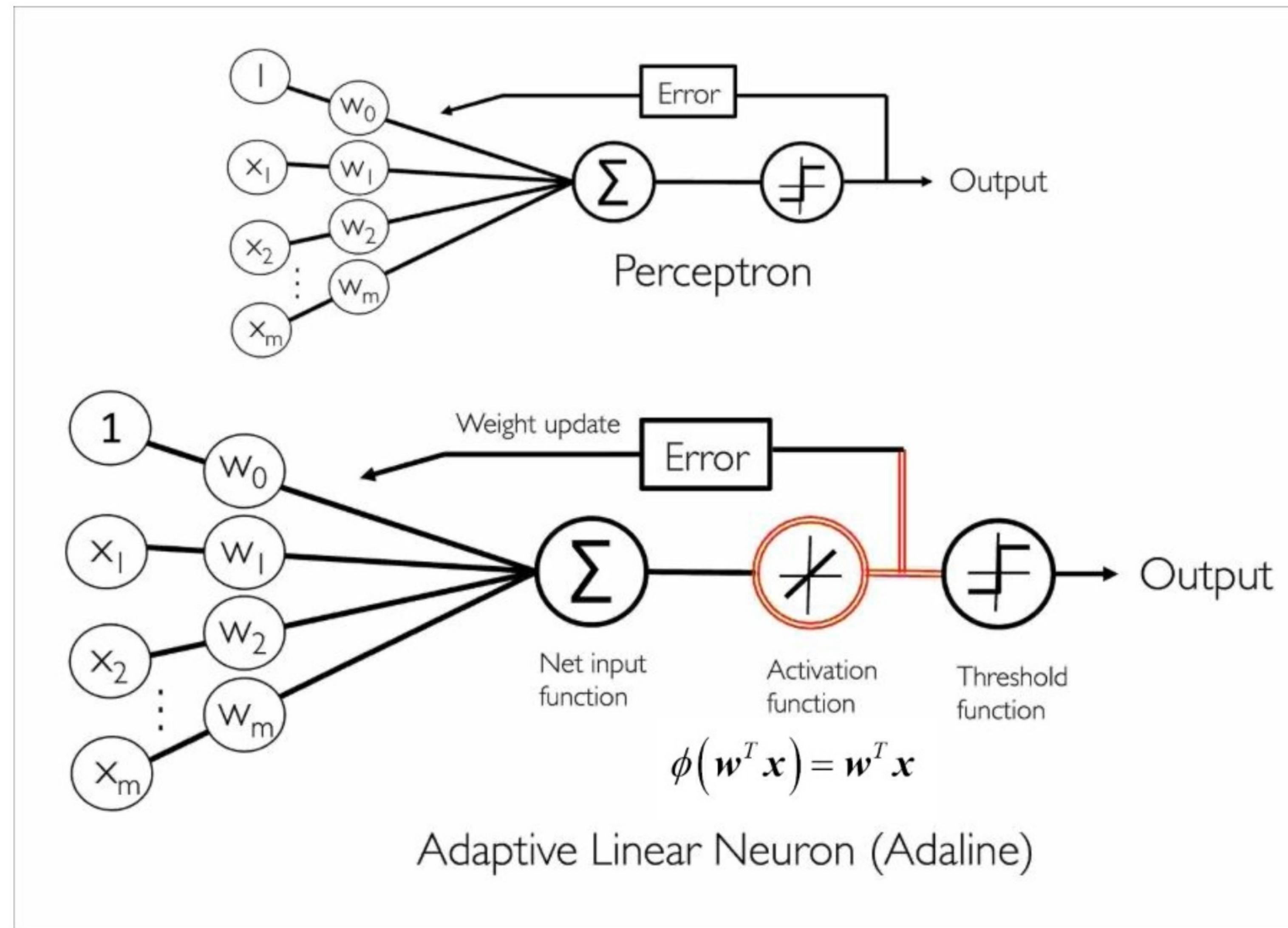
4



Introduction to the basic terminology and notations

Adaptive Linear Neuron

5



Introduction to the basic terminology and notations

Minimizing cost functions with gradient descent

6

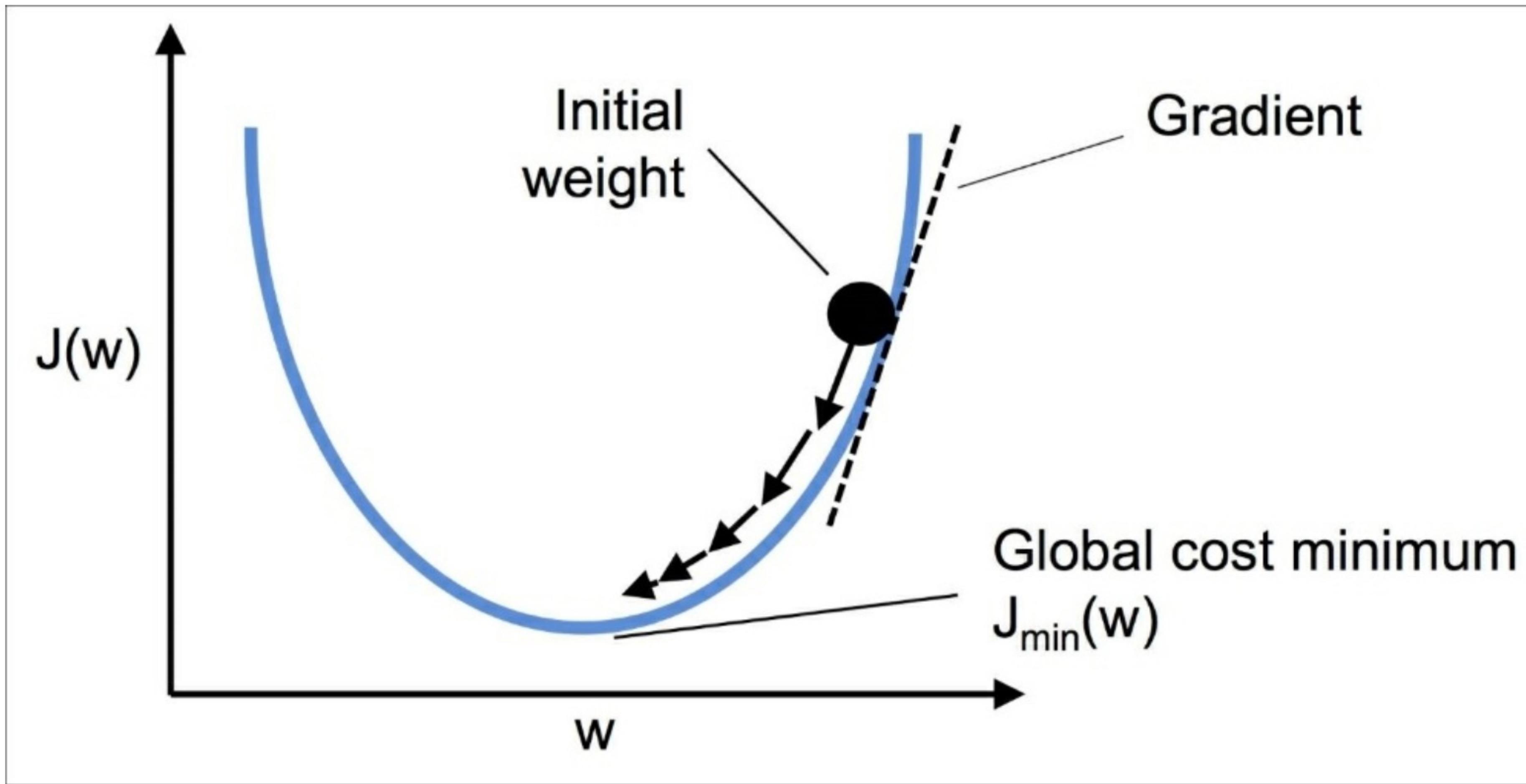
One of the key ingredients of supervised machine learning algorithms is a defined **objective function** that is to be optimized during the learning process. This objective function is often a cost function that we want to minimize. In the case of Adaline, we can define the cost function J to learn the weights as the **Sum of Squared Errors (SSE)** between the calculated outcome and the true class label:

$$J(\mathbf{w}) = \frac{1}{2} \sum_i (y^{(i)} - \phi(z^{(i)}))^2$$

Introduction to the basic terminology and notations

Minimizing cost functions with gradient descent

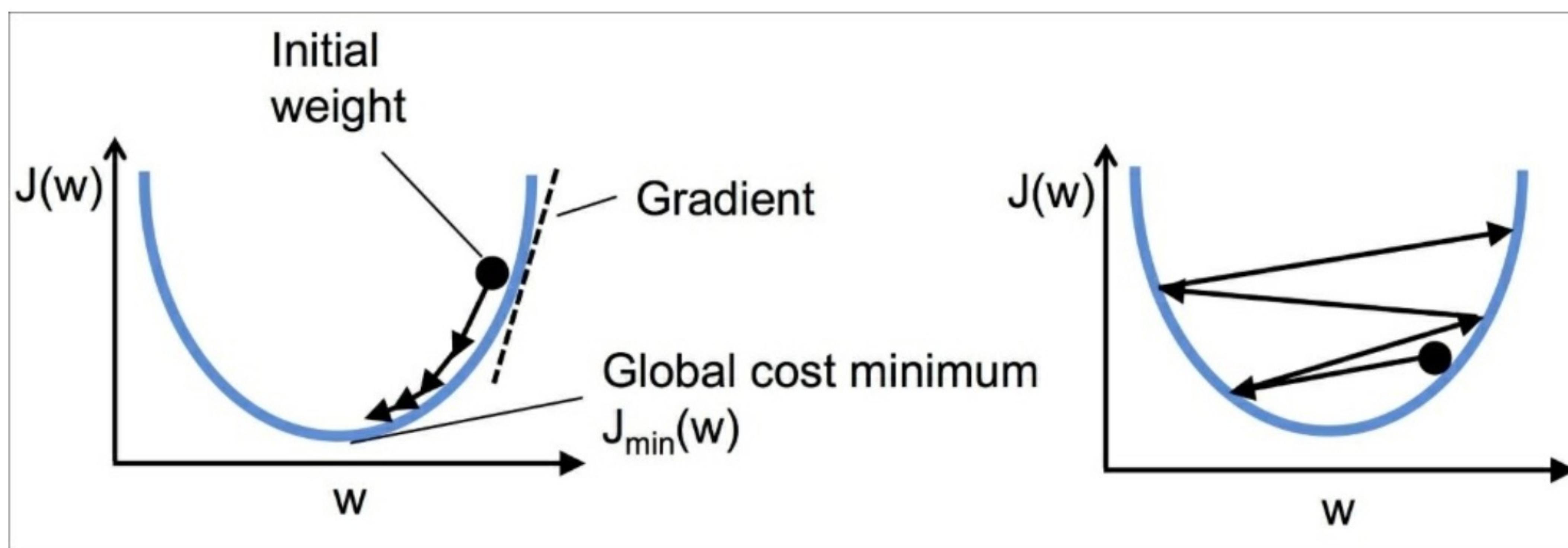
7



Introduction to the basic terminology and notations

Learning rate

8



Introduction to the basic terminology and notations

Learning rate

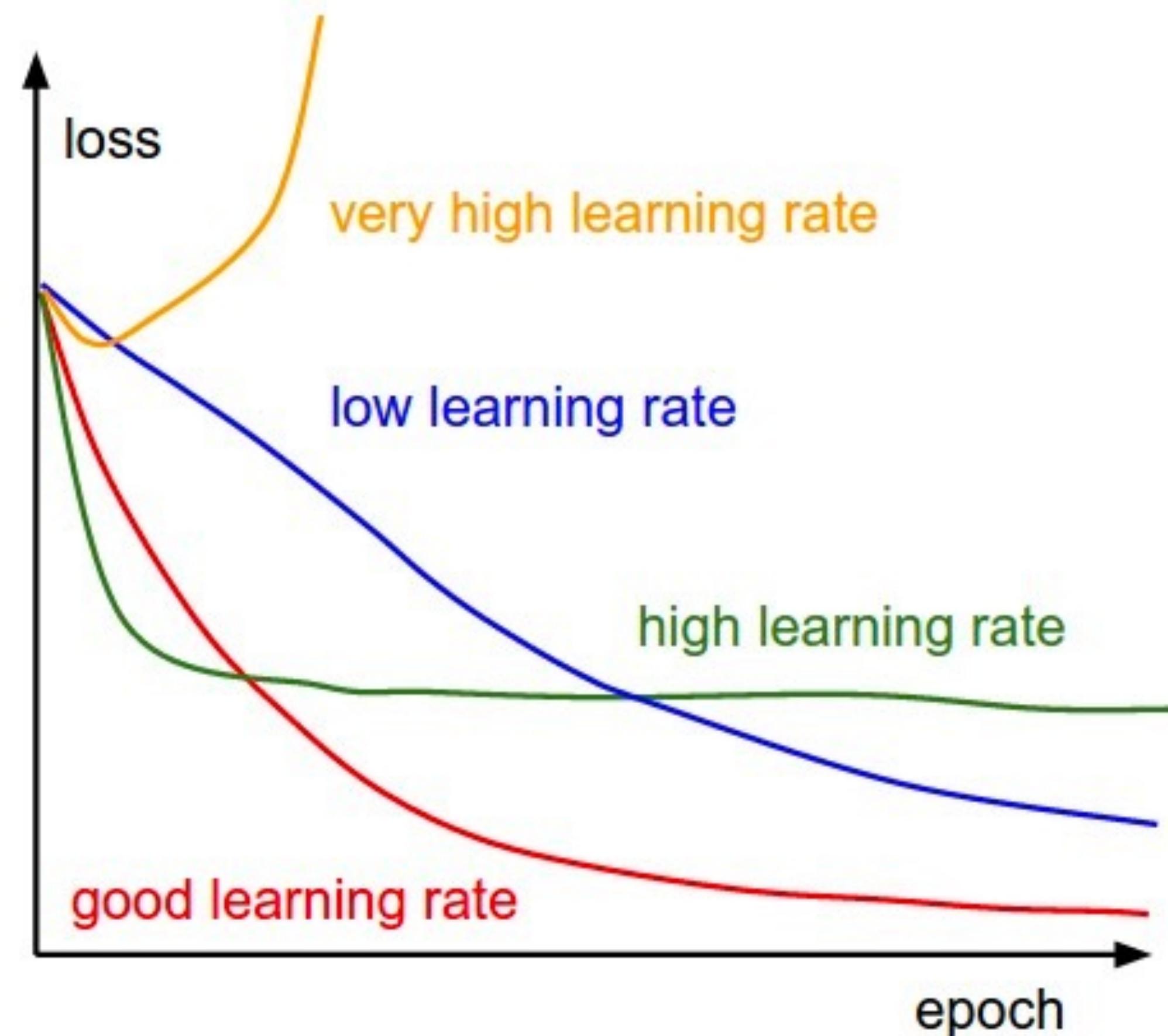
9

<https://the-learning-machine.com/article/optimization/gradient-descent>

Introduction to the basic terminology and notations

Learning rate

10



Introduction to the basic terminology and notations

Minimizing cost functions with gradient descent

11

Using gradient descent, we can now update the weights by taking a step in the opposite direction of the gradient $\nabla J(\mathbf{w})$ of our cost function $J(\mathbf{w})$:

$$\mathbf{w} := \mathbf{w} + \Delta\mathbf{w}$$

Where the weight change $\Delta\mathbf{w}$ is defined as the negative gradient multiplied by the learning rate η :

$$\Delta\mathbf{w} = -\eta \nabla J(\mathbf{w})$$

Introduction to the basic terminology and notations

Minimizing cost functions with gradient descent

12

To compute the gradient of the cost function, we need to compute the partial derivative of the cost function with respect to each weight w_j :

$$\frac{\partial J}{\partial w_j} = -\sum_i (y^{(i)} - \phi(z^{(i)}))x_j^{(i)}$$

So that we can write the update of weight w_j as:

$$\Delta w_j = -\eta \frac{\partial J}{\partial w_j} = \eta \sum_i (y^{(i)} - \phi(z^{(i)}))x_j^{(i)}$$

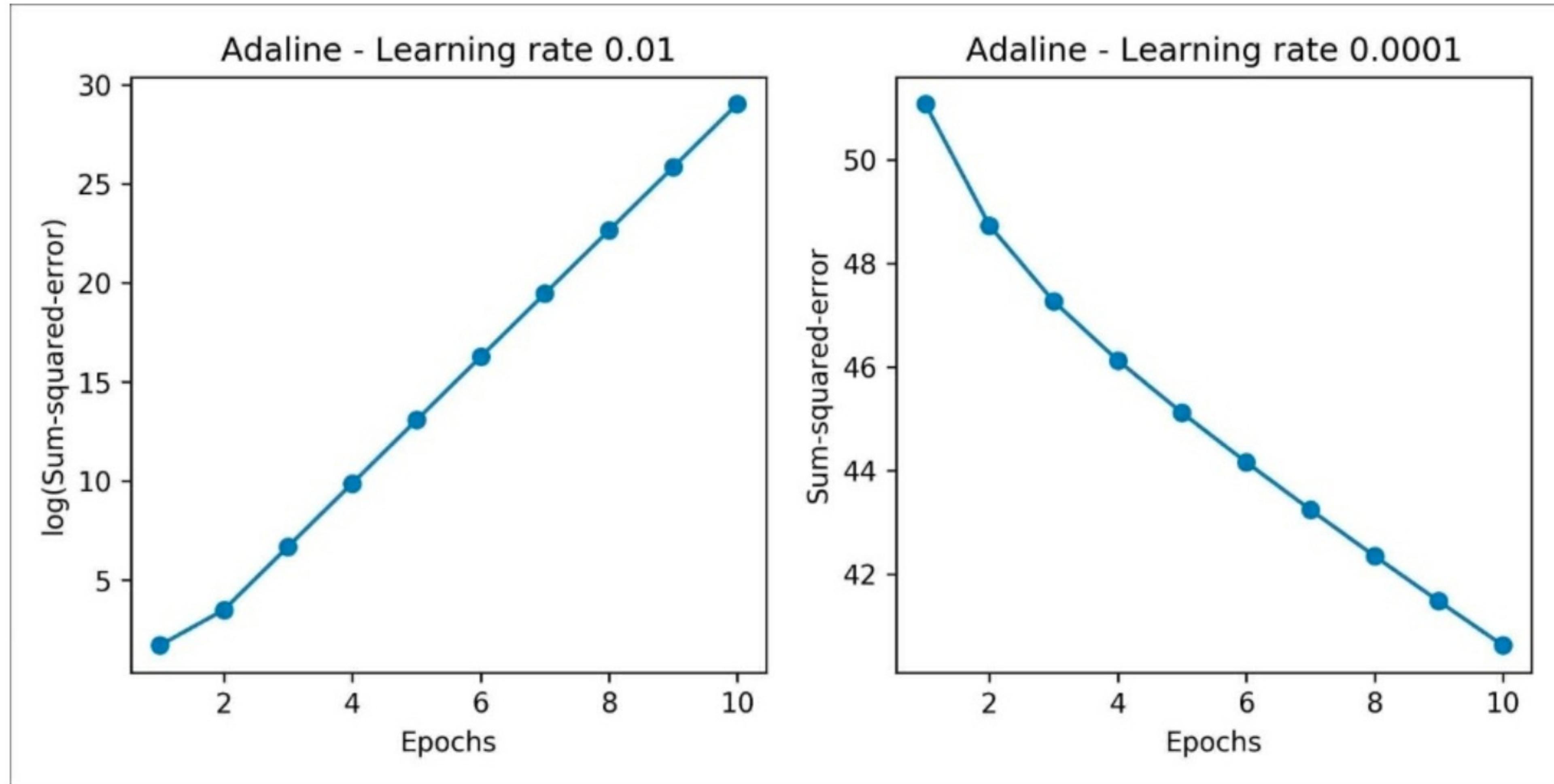
Since we update all weights simultaneously, our Adaline learning rule becomes:

$$\mathbf{w} := \mathbf{w} + \Delta \mathbf{w}$$

Introduction to the basic terminology and notations

Minimizing cost functions with gradient descent

13



Introduction to the basic terminology and notations

Improving gradient descent through feature scaling

14

Gradient descent is one of the many algorithms that benefit from feature scaling. In this section, we will use a feature scaling method called **standardization**, which gives our data the property of a standard normal distribution, which helps gradient descent learning to converge more quickly. Standardization shifts the mean of each feature so that it is centered at zero and each feature has a standard deviation of 1. For instance, to standardize the j th feature, we can simply subtract the sample mean

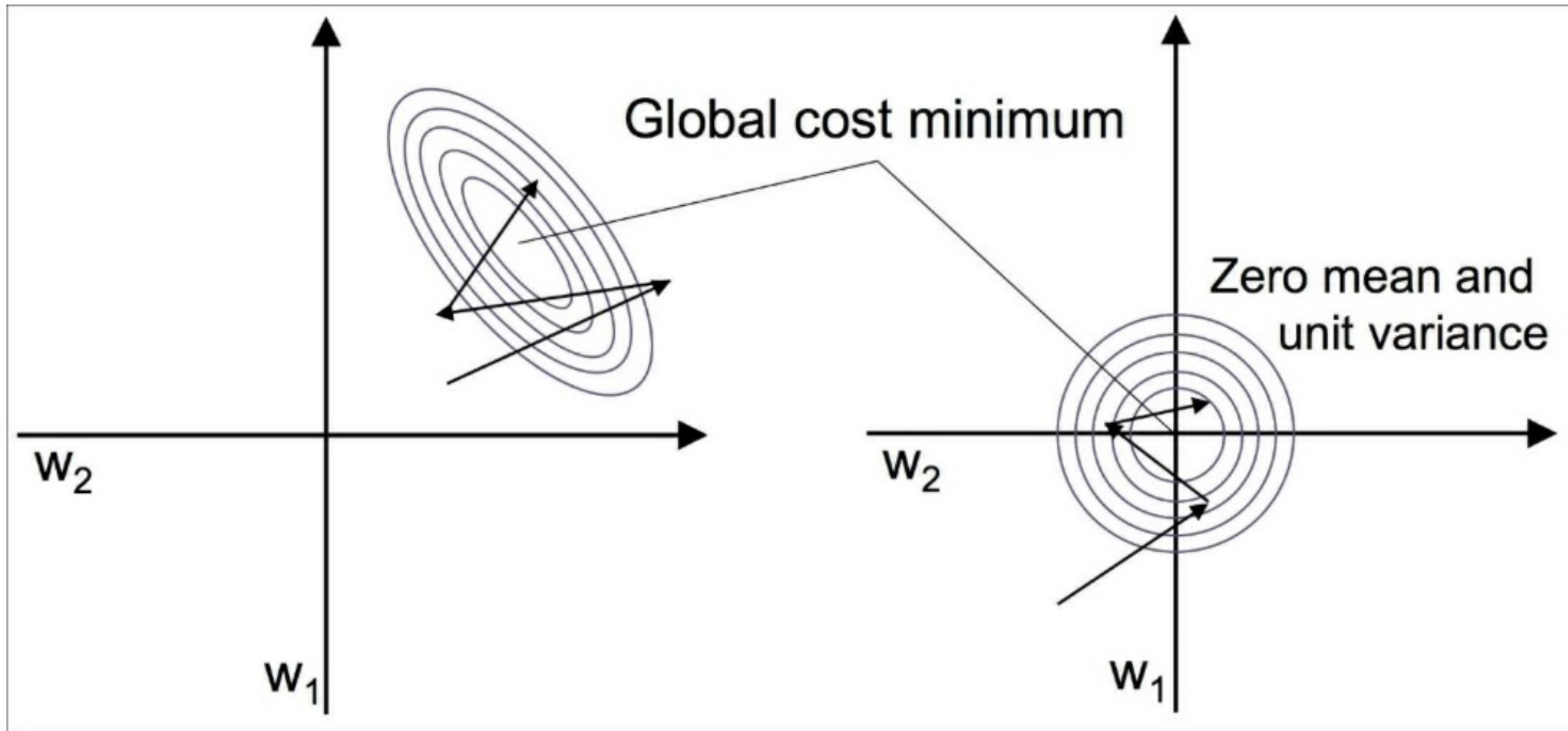
μ_j from every training sample and divide it by its standard deviation σ_j :

$$\mathbf{x}'_j = \frac{\mathbf{x}_j - \mu_j}{\sigma_j}$$

Introduction to the basic terminology and notations

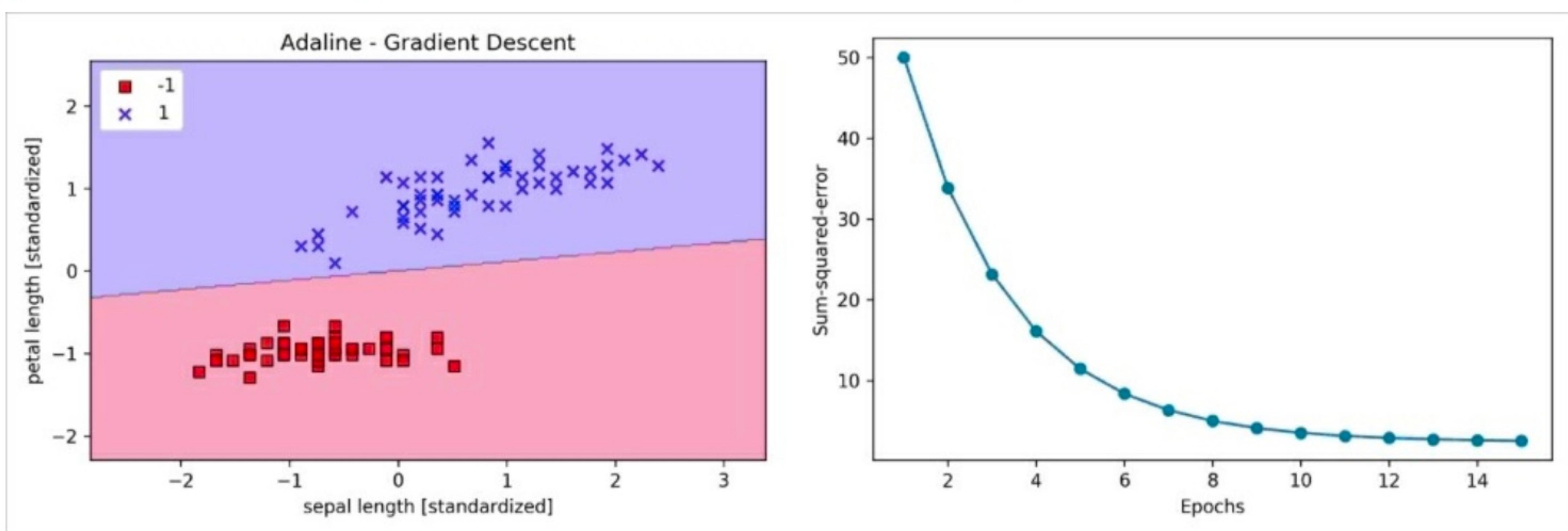
Improving gradient descent through feature scaling

15



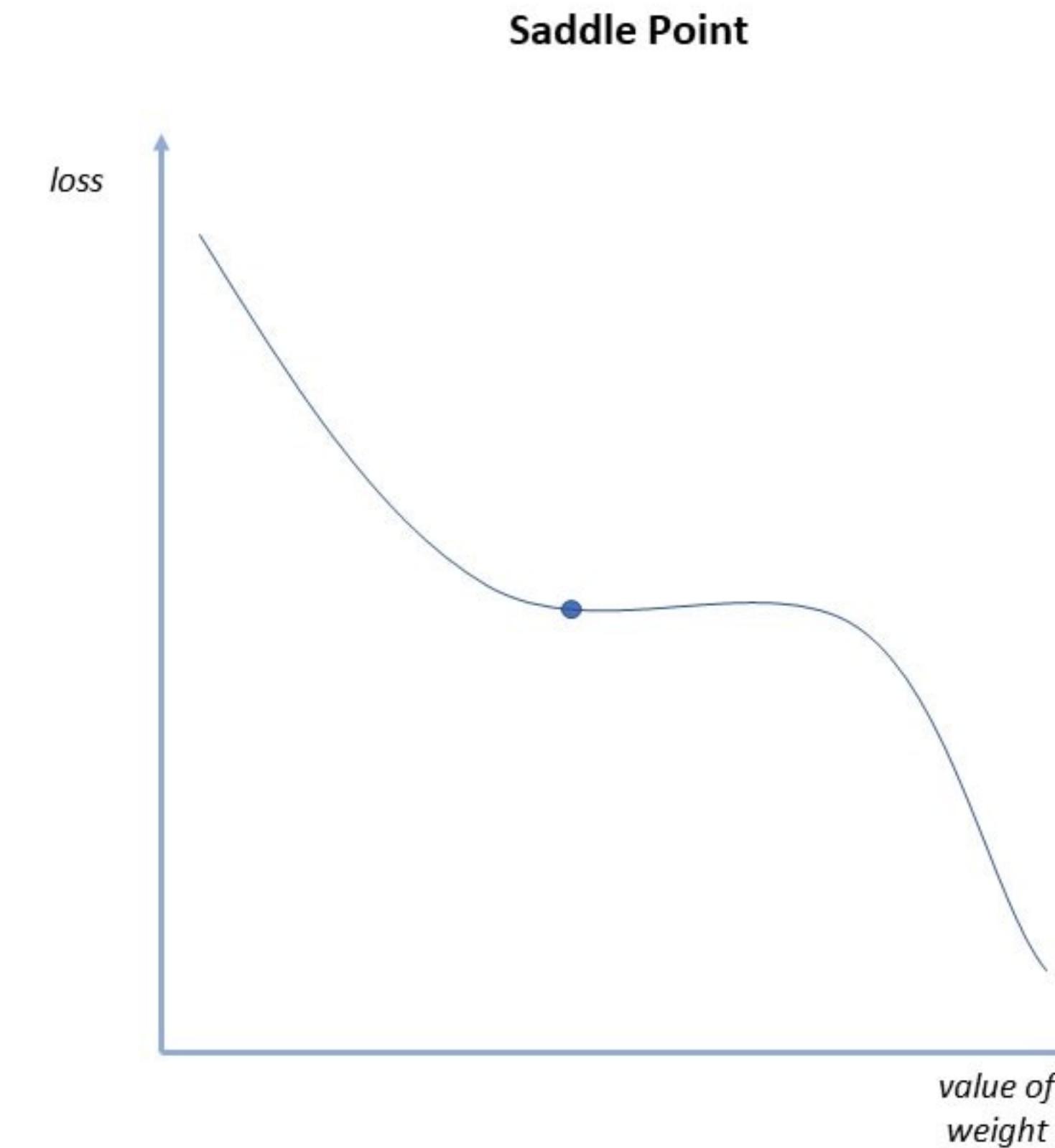
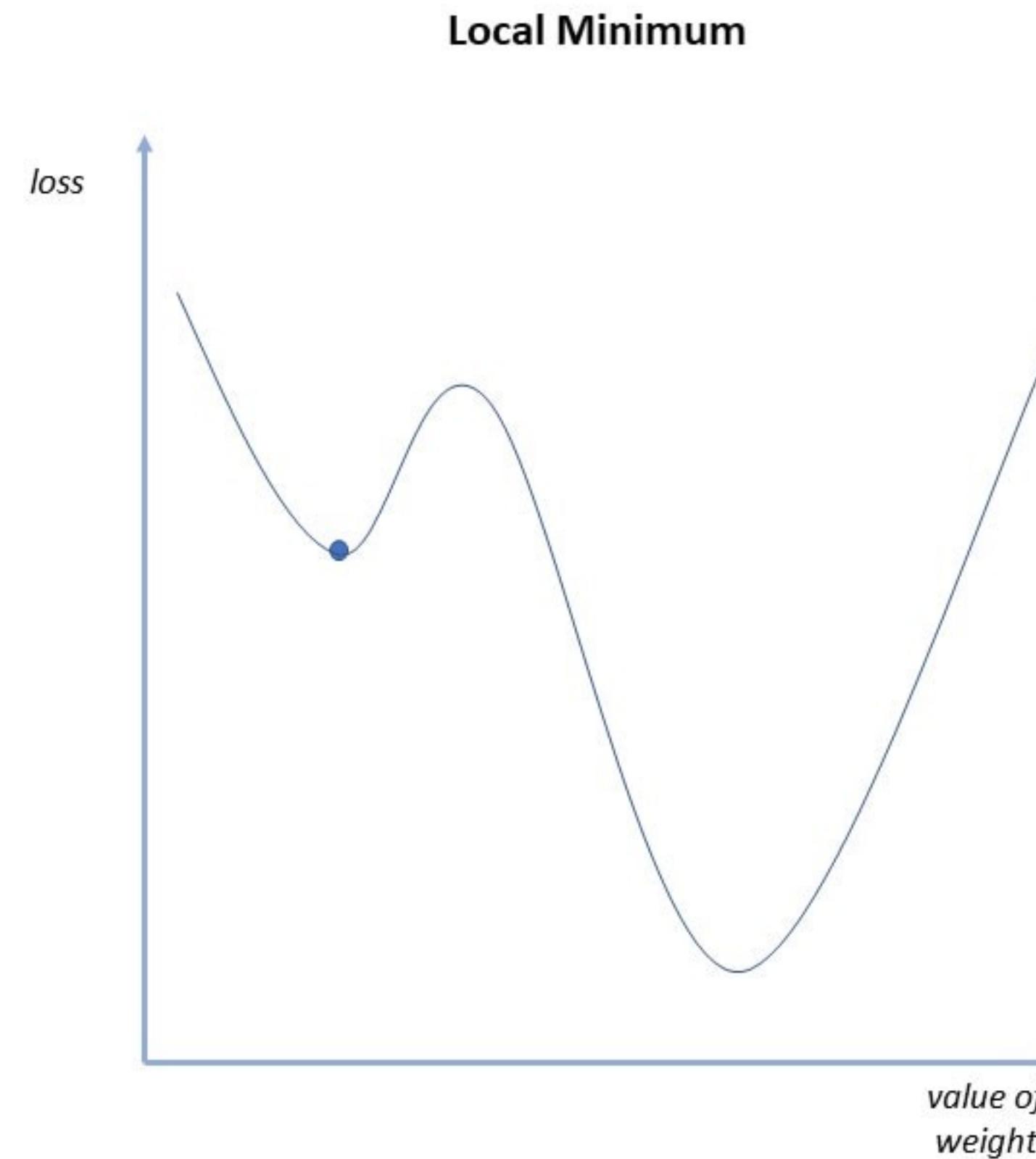
Introduction to the basic terminology and notations

Improving gradient descent through feature scaling



Introduction to the basic terminology and notations

Challenges with gradient descent



Introduction to the basic terminology and notations

Large-scale machine learning and stochastic gradient descent

18

In the previous section, we learned how to minimize a cost function by taking a step in the opposite direction of a cost gradient that is calculated from the whole training set; this is why this approach is sometimes also referred to as **batch gradient descent**. Now imagine we have a very large dataset with millions of data points, which is not uncommon in many machine learning applications. Running batch gradient descent can be computationally quite costly in such scenarios since we need to reevaluate the whole training dataset each time we take one step towards the global minimum.

A popular alternative to the batch gradient descent algorithm is **stochastic gradient descent**, sometimes also called iterative or online gradient descent. Instead of updating the weights based on the sum of the accumulated errors over all samples

$\mathbf{x}^{(i)}$:

$$\Delta \mathbf{w} = \eta \sum_i (y^{(i)} - \phi(z^{(i)})) \mathbf{x}^{(i)}$$

We update the weights incrementally for each training sample:

$$\eta (y^{(i)} - \phi(z^{(i)})) \mathbf{x}^{(i)}$$

Introduction to the basic terminology and notations

Large-scale machine learning and stochastic gradient descent

19

In the previous section, we learned how to minimize a cost function by taking a step in the opposite direction of a cost gradient that is calculated from the whole training set; this is why this approach is sometimes also referred to as **batch gradient descent**. Now imagine we have a very large dataset with millions of data points, which is not uncommon in many machine learning applications. Running batch gradient descent can be computationally quite costly in such scenarios since we need to reevaluate the whole training dataset each time we take one step towards the global minimum.

A popular alternative to the batch gradient descent algorithm is **stochastic gradient descent**, sometimes also called iterative or online gradient descent. Instead of updating the weights based on the sum of the accumulated errors over all samples

$\mathbf{x}^{(i)}$:

$$\Delta \mathbf{w} = \eta \sum_i (y^{(i)} - \phi(z^{(i)})) \mathbf{x}^{(i)}$$

We update the weights incrementally for each training sample:

$$\eta (y^{(i)} - \phi(z^{(i)})) \mathbf{x}^{(i)}$$

Introduction to the basic terminology and notations

Large-scale machine learning and stochastic gradient descent – adaptive learning rate

20

In stochastic gradient descent implementations, the fixed learning rate η is often replaced by an adaptive learning rate that decreases over time, for example:

$$\frac{c_1}{[\text{number of iterations}] + c_2}$$

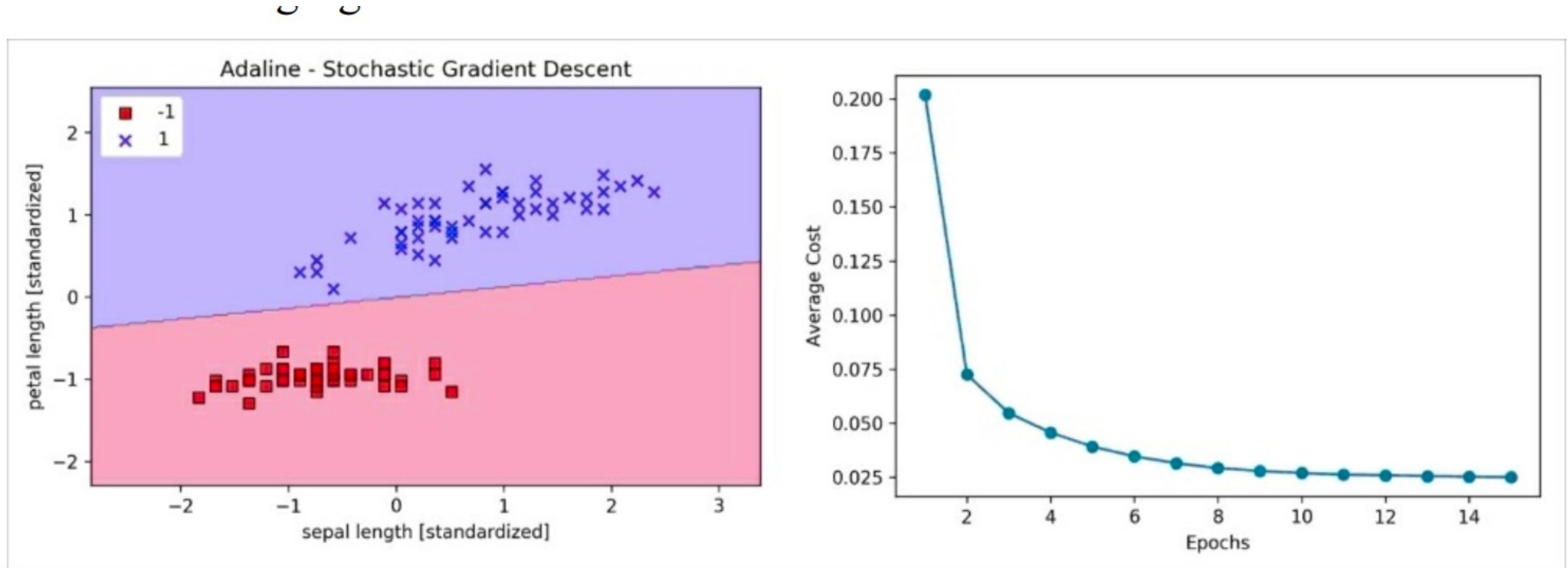
Where c_1 and c_2 are constants. We shall note that stochastic gradient descent does not reach the global minimum, but an area very close to it. And using an adaptive learning rate, we can achieve further annealing to the cost minimum.

Another advantage of stochastic gradient descent is that we can use it for **online learning**. In online learning, our model is trained on the fly as new training data arrives. This is especially useful if we are accumulating large amounts of data, for example, customer data in web applications. Using online learning, the system can immediately adapt to changes and the training data can be discarded after updating the model if storage space is an issue.

Introduction to the basic terminology and notations

Large-scale machine learning and stochastic gradient descent – adaptive learning rate

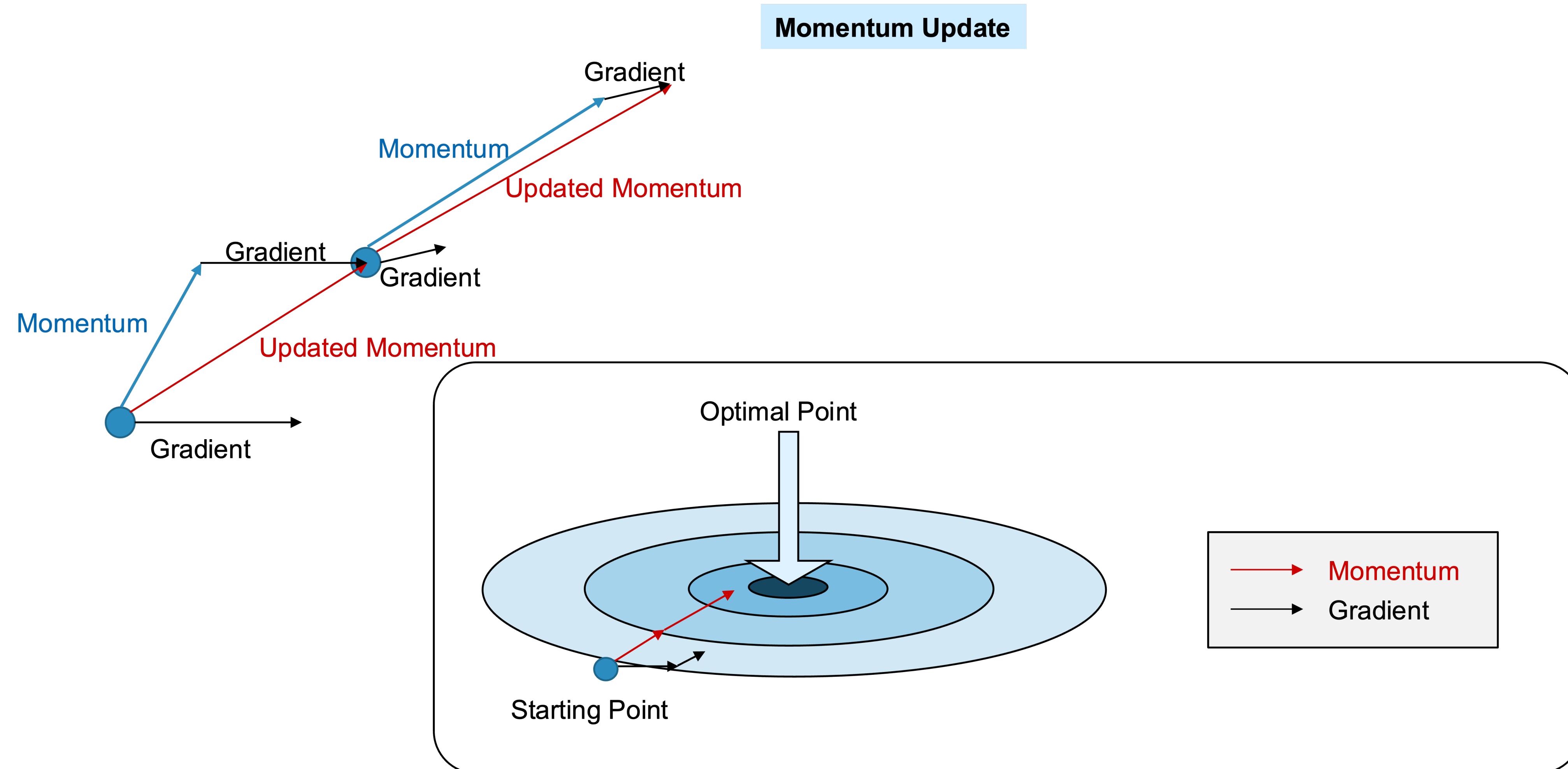
21



Introduction to the basic terminology and notations

SGD vs SGD with momentum

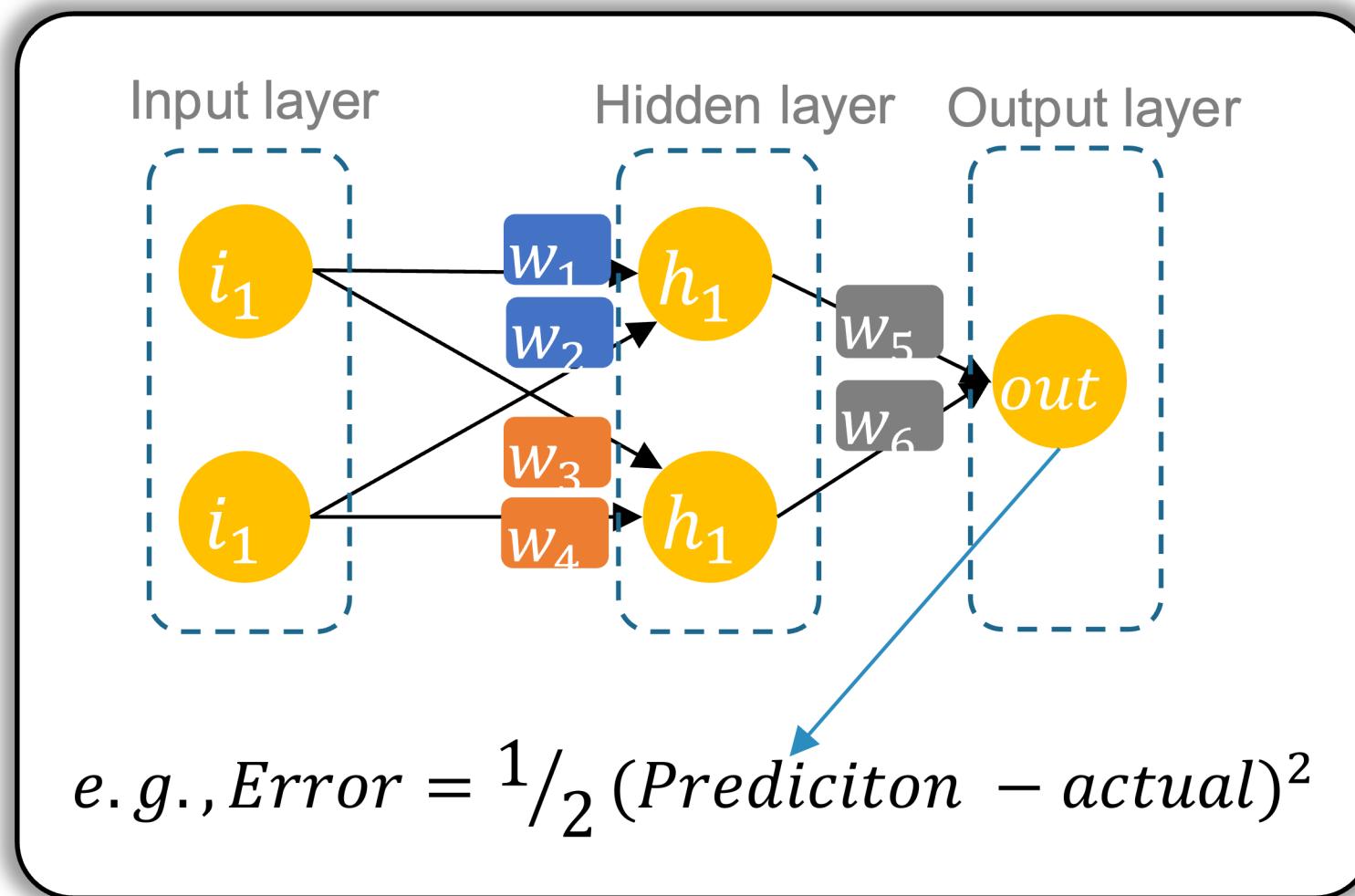
22



Introduction to the basic terminology and notations

23

SGD without momentum vs SGD with momentum



Derivative of error with respect to each weight

Old weight

New weight

Learning rate

$$W_x = W_x - \gamma \cdot \left(\frac{\partial Error}{\partial W_x} \right)$$

(a) SGD without momentum

Old momentum

Derivative of error with respect to each weight

New momentum

Learning rate

$$v_x = v_x + \gamma \cdot \left(\frac{\partial Error}{\partial W_x} \right)$$

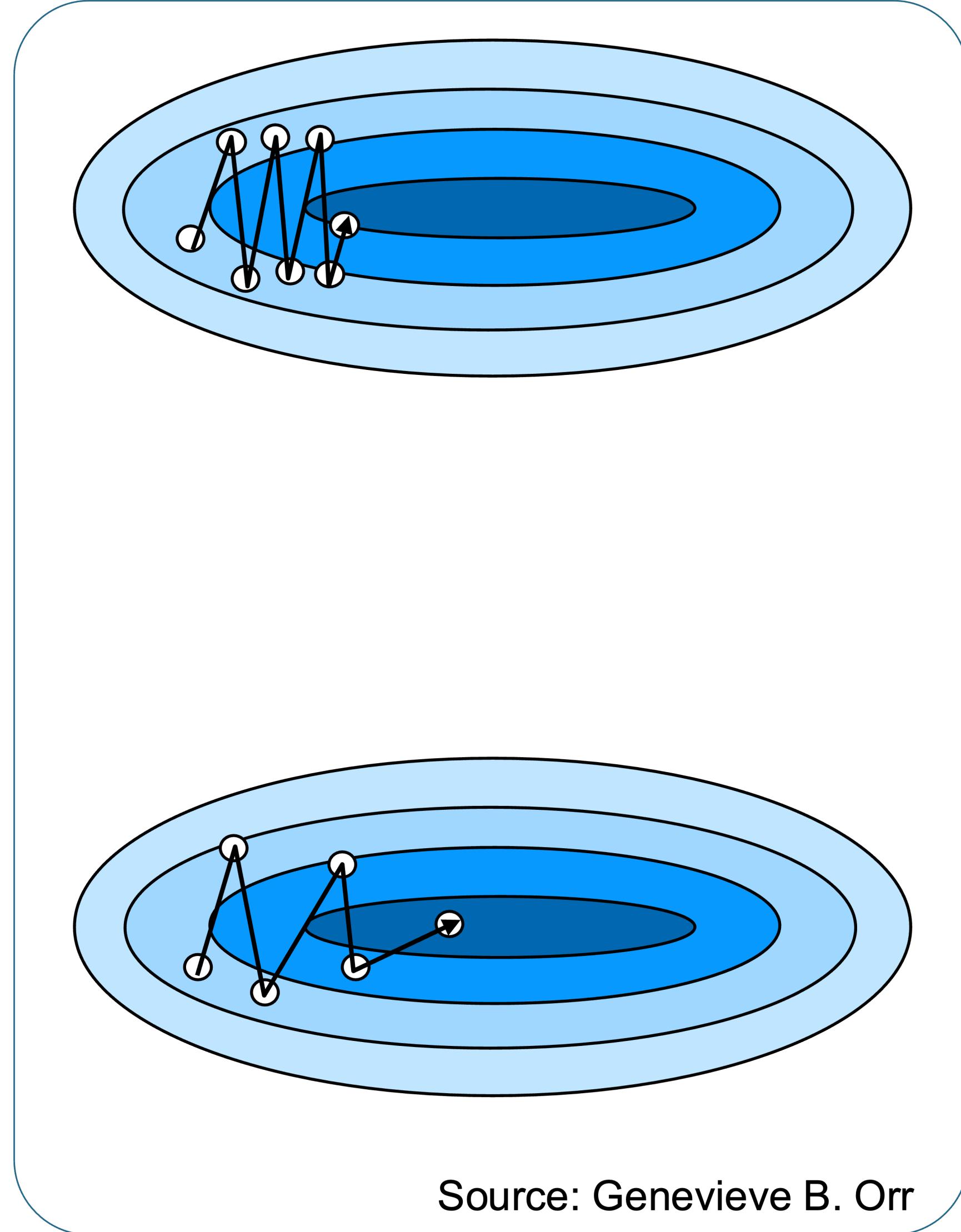
Old weight

New weight

Updated momentum

$$W_x = W_x - v_x$$

(b) SGD with momentum





Q&A