
CHAPTER 1

TRAINING SIMPLE MACHINE LEARNING ALGORITHMS FOR CLASSIFICATION

So, we will see basic algorithms for classification and some other useful things.

1.1 INTRODUCTION

In this section we will make use of the first two algorithmically described machine learning algorithms for classification: the **perceptron** and **adaptive linear neurons**.

Frist, we'll start by implementing a perceptron step by step in Python and training it to classify different flower species in the iris dataset.

Observation 1.1.1

Using this algorithm and discussing the basics of optimization using adaptive linear neurons layse the groundwork for using more sophisticated classifiers.

The goals of this chapter are the following:

- Building an understanding of machine learning algorithms
- Using pandas, NumPy, and Matplotlib to read in, process, and visualize data
- Implementing linear classifiers for 2-class problems in Python

1.2 ARTIFICIAL NEURONS: EARLY STAGES OF MACHINE LEARNING

Warren McCulloch and Walter Pitts published the first concept of a simplified brain cell, the so-called McCulloch-Pitts (MCP) neuron, in 1943 (A Logical Calculus of the Ideas Immanent in Nervous Activity by W. S. McCulloch and W. Pitts, Bulletin of Mathematical Biophysics, 5(4): 115-133, 1943).

Observation 1.2.1 (Biological Neurons)

Biological Neurons are interconnected nerve cells in the brain that are involved in the processing and transmitting of chemical and electrical signals.

Idea 1.2.1

They described a nerve cell as a **simple logic gate with binary outputs**.

Multiple signals arrive at the dendrites, they are then integrated into the cell body and, if the accumulated signals exceeds a certain threshold, an output signal is generated that will be passed on by the axon.

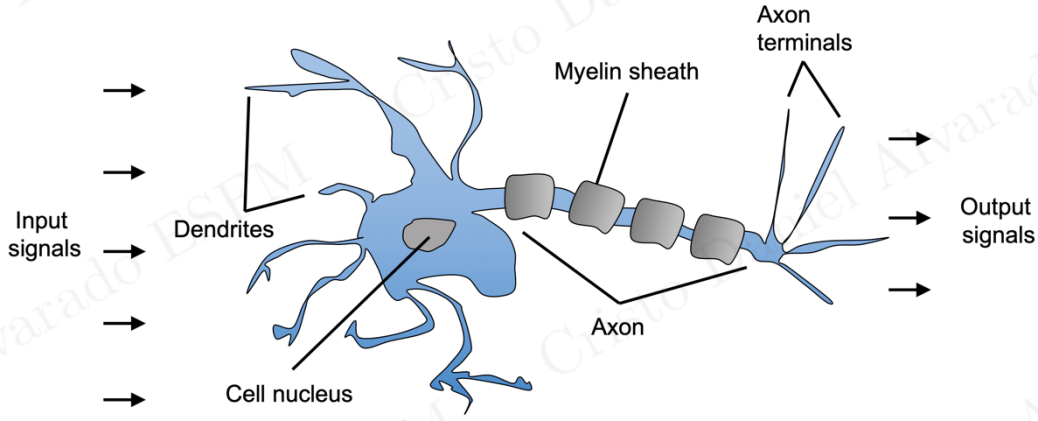


Figure 1.1: Biological Neuron.

Observation 1.2.2 (Perceptron Learning Rule)

Only a few years later, Frank Rosenblatt published the first concept of the **perceptron learning rule** based on the MCP neuron model (*The Perceptron: A Perceiving and Recognizing Automaton* by F. Rosenblatt, Cornell Aeronautical Laboratory, 1957). With his perceptron rule, Rosenblatt *proposed an algorithm that would automatically learn the optimal weight coefficients that would then be multiplied with the input features in order to make the decision of whether a neuron fires (transmits a signal) or not.*

In the context of supervised learning and classification, *such an algorithm could then be used to predict whether a new data point belongs to one class or the other.*

1.3 FORMAL DEFINITION OF AN ARTIFICIAL NEURON

We can put the idea of **artificial neurons** into the context of a binary classification task with two classes: 0 and 1.

Definition 1.3.1 (Decision Function)

A **decision function** is a function $\sigma : X \rightarrow \{0, 1\}$ that takes a linear combination of a certain input values (called x) and a corresponding weight vector (called w), where z is the net input:

$$z = w_1x_1 + w_2x_2 + \cdots + w_mx_m = \sum_{i=1}^m w_ix_i$$

so, $\sigma(z) \in \{0, 1\}$.

If the net input of a particular example is greater than a defined threshold (lets say $\vartheta \in \mathbb{R}$, so $z \geq \vartheta$), we predict class 1, and class 0 otherwise. In the perceptron algorithm, the decision function, $\sigma(z)$ is a variant of a **unit step function**:

$$\sigma(z) = \begin{cases} 1 & \text{if } z \geq \vartheta \\ 0 & \text{otherwise} \end{cases}$$

Observation 1.3.1

We can modify the setup with a couple of steps. So, we make the **bias unit** $b = -\vartheta$, and we make:

$$z' = w_1x_1 + w_2x_2 + \cdots + w_mx_m + b = \vec{w}^T \vec{x} + b$$

And now replacing every input z by z' we obtain that:

$$\sigma(z') = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Observation 1.3.2

The transpose A^T of matrix $A \in \mathcal{M}_{m \times n}$ is defined as:

$$A^T = \begin{bmatrix} a_1^{(1)} & a_2^{(1)} & \cdots & a_n^{(1)} \\ a_1^{(2)} & a_2^{(2)} & \cdots & a_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ a_1^{(m)} & a_2^{(m)} & \cdots & a_n^{(m)} \end{bmatrix}$$

where:

$$A = \begin{bmatrix} a_1^{(1)} & a_2^{(1)} & \cdots & a_m^{(1)} \\ a_1^{(2)} & a_2^{(2)} & \cdots & a_m^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ a_1^{(n)} & a_2^{(n)} & \cdots & a_m^{(n)} \end{bmatrix}$$

Observation 1.3.3 (Notation Problem)

This notation is horrible and differs from the one Mathematicians use in all his textbooks, please, don't use this notation.

1.4 PERCEPTRON LEARNING RULE

The whole idea behind the MCP neuron and Rosenblatt's thresholded perceptron model is to *use a reductionist approach to mimic how a single neuron in the brain works: it either fires or it doesn't*. Thus, Rosenblatt's classic perceptron rule is fairly simple, and the perceptron algorithm can be summarized by the following steps:

1. Initialize the weights and bias unit to 0 or small random numbers
2. For each training example, $x^{(i)}$.
3. Compute the output value, $\hat{y}^{(i)}$.
4. Update the weight and bias unit.

Observation 1.4.1 (Output Value)

Here, the **output value** is the **class label** predicted by the unit step function defined earlier.

The simultaneous update of each weight w_j in the weight vector \vec{w} can be formally written as:

$$w_j \equiv w_j + \Delta w_j$$

and,

$$b \equiv b + \Delta b$$

Observation 1.4.2 (Use of the Symbol \equiv)

Here, the symbol \equiv states that the value of the variable to the left is updated to that of the variable on the right.

Where, the updated values (or **deltas**) are computed as follows:

$$\Delta w_j = \eta(y^{(i)} - \hat{y}^{(i)})x_j^{(i)}$$

and,

$$\Delta b = \eta(y^{(i)} - \hat{y}^{(i)})$$

Observation 1.4.3 (Notation Abuse)

This thing is absolutely ridiculous in terms of notation, the correct notation should be something like $\Delta w_j^{(i)}$, but since this variable is not going to be used anymore, we simply forgot about it, this is because in the process this variable is going to be updated every single time and we don't care about the different values depending on the i .

Here, η is a learning rate (usually between 0 and 1). $y^{(i)}$ is the **true class label** and $\hat{y}^{(i)}$ is the **predicted class label**.

Observation 1.4.4

The bias unit and all weights in the weight vector are updated simultaneously, which means that the predicted label $\hat{y}^{(i)}$ is not recomputed before the bias unit and all the weights are updated via the respective update values, Δw_j and Δb .

Example 1.4.1

Concretely, for a two-dimensional dataset, we would write the update as:

$$\Delta w_1 = \eta(y^{(i)} - \text{output}^{(i)})x_1^{(i)}$$

$$\Delta w_2 = \eta(y^{(i)} - \text{output}^{(i)})x_2^{(i)}$$

$$\Delta b = \eta(y^{(i)} - \text{output}^{(i)})$$

1. Before we implement the perceptron rule in Python, let's go through a simple thought experiment to illustrate how beautifully simple this learning rule really is. In the two scenarios where the perceptron predicts the class label correctly, the bias unit and weights remain unchanged, since the update values are 0:

2. **When prediction is correct:**

- If $y^{(i)} = 0$ and $\hat{y}^{(i)} = 0$:

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

$$\Delta b = \eta(0 - 0) = 0$$

- If $y^{(i)} = 1$ and $\hat{y}^{(i)} = 1$:

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

$$\Delta b = \eta(1 - 1) = 0$$

3. When prediction is wrong:

- If $y^{(i)} = 0$ and $\hat{y}^{(i)} = 1$:

$$\begin{aligned}\Delta w_j &= \eta(0 - 1)x_j^{(i)} = -\eta x_j^{(i)} \\ \Delta b &= \eta(0 - 1) = -\eta\end{aligned}$$

- If $y^{(i)} = 1$ and $\hat{y}^{(i)} = 0$:

$$\begin{aligned}\Delta w_j &= \eta(1 - 0)x_j^{(i)} = \eta x_j^{(i)} \\ \Delta b &= \eta(1 - 0) = \eta\end{aligned}$$

4. To get a better understanding of the feature value as a multiplicative factor $x_j^{(i)}$, consider another example where $y^{(i)} = 1$, $\hat{y}^{(i)} = 0$, $\eta = 1$. Assume that $x_j^{(i)} = 1.5$ and this example is misclassified as class 0. In this case, we would increase the corresponding weight so that the net input $z = x_j^{(i)}w_j + b$ would be more positive the next time we encounter this example and thus more likely to be above the threshold of the unit step function to classify the example as class 1:

$$\begin{aligned}\Delta w_j &= (1 - 0) \times 1.5 = 1.5 \\ \Delta b &= (1 - 0) = 1\end{aligned}$$

5. The weight update Δw_j is proportional to the value of $x_j^{(i)}$. For instance, if we have another example $x_j^{(i)} = 2$ that is incorrectly classified as class 0, we will push the decision boundary by an even larger extent to classify this example correctly the next time:

$$\begin{aligned}\Delta w_j &= (1 - 0) \times 2 = 2 \\ \Delta b &= (1 - 0) = 1\end{aligned}$$

6. The convergence of the perceptron is only guaranteed if the two classes are linearly separable, which means that the two classes can be perfectly separated by a linear decision boundary. If the two classes cannot be separated by a linear decision boundary, we can set a maximum number of passes over the training dataset (epochs) and/or a threshold for the number of tolerated misclassifications—the perceptron would never stop updating the weights otherwise.
7. The perceptron receives the inputs of an example (x) and combines them with the bias unit (b) and weights (w) to compute the net input. The net input is then passed on to the threshold function, which generates a binary output of 0 or 1—the predicted class label of the example. During the learning phase, this output is used to calculate the error of the prediction and update the weights and bias unit.

1.5 IMPLEMENTING A PERCEPTRON LEARNING ALGORITHM IN PYTHON

Now, we will do an implementation of the Perceptron Rule in Python.

1.6 AND OBJECT-ORIENTED PERCEPTRON API

We will take an object-oriented approach to defining the perceptron interface as a Python class, which will allow us to initialize new Perceptron objects that can learn from data via a `fit` method and make predictions via a separate `predict` method.

Observation 1.6.1

We append an underscore `_` to attributes that are not created upon the initialization of the object, but we do this by calling the object's other methods, for example, `self._w_`.

The code is the following:

```
1 import numpy as np
2
3 class Perceptron:
4     """Perceptron classifier.
5
6     Parameters
7     -----
8     eta : float
9         Learning rate (between 0.0 and 1.0)
10    n_iter : int
11        Passes over the training dataset.
12    random_state : int
13        Random number generator seed for random weight
14        initialization.
15
16    Attributes
17    -----
18    w_ : 1d-array
19        Weights after fitting.
20    b_ : Scalar
21        Bias unit after fitting.
22    errors_ : list
23        Number of misclassifications (updates) in each epoch.
24
25    """
26    def __init__(self, eta=0.01, n_iter=50, random_state=1):
27        self.eta = eta
28        self.n_iter = n_iter
29        self.random_state = random_state
30
31    def fit(self, X, y):
32        """Fit training data.
33
34        Parameters
35        -----
36        X : {array-like}, shape = [n_examples, n_features]
37            Training vectors, where n_examples is the number of
38            examples and n_features is the number of features.
39        y : array-like, shape = [n_examples]
40            Target values.
41
42        Returns
43        -----
44        self : object
45
46        """
```



```

47     rgen = np.random.RandomState(self.random_state)
48     self.w_ = rgen.normal(loc=0.0, scale=0.01, size=X.shape
49         [1])
49     self.b_ = np.float_(0.)
50     self.errors_ = []
51
52     for _ in range(self.n_iter):
53         errors = 0
54         for xi, target in zip(X, y):
55             update = self.eta * (target - self.predict(xi))
56             self.w_ += update * xi
57             self.b_ += update
58             errors += int(update != 0.0)
59         self.errors_.append(errors)
60     return self
61
62     def net_input(self, X):
63         """Calculate net input"""
64         return np.dot(X, self.w_) + self.b_
65
66     def predict(self, X):
67         """Return class label after unit step"""
68         return np.where(self.net_input(X) >= 0.0, 1, 0)

```

Code 1.1: Python Perceptron Rule Implementation.

First, we will list the parameters and attributes of the code, in order to understand it better:

- **eta**: Is a float, which represents the learning rate (a number between 0.0 and 1.0).
- **n_iter**: Is an integer, which is the number of passes over the training dataset.
- **random_state**: It is an integer, which is a **random number generator seed** for random weight initialization.

Definition 1.6.1 (Random Number Generator (RNG))

A **random number generator (RNG) seed** is an initial value used by a computer algorithm to generate a sequence of pseudo-random numbers. For random weight initialization in a neural network, this seed sets the starting point for the sequence of numbers used to initialize the model's weights and biases.

The attributes are the following:

- **w_**: It is a 1-dimensional array, which represents the weights after fitting.
- **b_**: It is a float, which represents the bias unit after fitting.
- **errors_**: It is a list, which is the number of misclassifications (or updates) in each **epoch**.

Definition 1.6.2 (Epoch)

In machine learning, an **epoch** is one complete pass through the entire training dataset. During a single epoch, the learning algorithm processes every training example once to update the model's internal parameters. Training a model over multiple epochs allows it to progressively learn and refine its understanding of the data, with a well-tuned number of epochs balancing performance

between underfitting and overfitting.

In this code, we have three methods in the class `Perceptron`:

- `__init__`: It initializes the parameters of the class, in this case, the `eta` is set to 0.01, `n_iter=50` and `random_state=1`.
- `fit`: It is the code that fits the data in a dataset, given an array `X` and another array `y` of target values.
- `net_input`: It calculates the net input of a dataset `X`.
- `predict`: It calculates the class label after a unit step of a dataset `X`.

Observation 1.6.2 (Use of the Perceptron Implementation)

Using this perceptron implementation, we can initialize new `Perceptron` objects with a given learning rate, `eta`, and the number of epochs, `n_iter` (or iterations over the training dataset).

Via the `fit` method, we initialize the bias `self.b_` to an initial value 0 and the weights in `self.w_` to a vector in \mathbb{R}^m . Here, m stands for the number of dimensions (features) in the dataset.

Observation 1.6.3 (Decision of Initialization of Weight Vectors)

The **initial weight** vector contains small random numbers drawn from a normal distribution with a standard deviation of 0.01 via `rgen.normal(loc=0.0, scale=0.01, size=1 + X.shape[1])`, where `rgen` is a NumPy random number generator seeded with a user-specified random seed so that results can be reproduced when desired.

Technically, we could initialize the weights to zero (this is done in the original perceptron algorithm). However, if we did that, the **learning rate** η (`eta`) would have no effect on the decision boundary. If all the weights are initialized to zero, the learning rate parameter affects only the scale of the weight vector, not the direction.

Observation 1.6.4 (Use of Normal Distribution)

Our decision to draw the random numbers from a normal distribution—for example, instead of from a uniform distribution—and to use a standard deviation of 0.01 was arbitrary; remember, we are just interested in small random values to avoid the properties of all-zero vectors, as discussed earlier.

Exercise 1.6.1

As an optional exercise, you can change `self.w_ = rgen.normal(loc=0.0, scale=0.01, size=X.shape[1])` to `self.w_ = np.zeros(X.shape[1])` and run the perceptron training code presented in the next section with different values for `eta`. You will observe that the decision boundary does not change.

After the weights have been initialized, the `fit` method loops over all individual examples in the training dataset and updates the weights according to the perceptron learning rule discussed in the previous section.

Observation 1.6.5 (Prediction)

The **class labels are predicted by the predict method**, which is called in the *fit method* during training to get the class label for the weight update; but *predict* can also be used to predict the class labels of new data after the model has been fitted. Furthermore, we collect the number of misclassifications during each epoch in the `self.errors_` list so that we can later analyze how well the perceptron performed during training. The `np.dot` function used in the `net_input` method calculates the vector dot product, $wTx + bwTx + b$.

Idea 1.6.1

Instead of using **NumPy** to calculate the vector dot product between two arrays, `a` and `b`, via `a.dot(b)` or `np.dot(a, b)`, we could also perform the calculation in pure Python via `sum([i * j for i, j in zip(a, b)])`. However, the advantage of using **NumPy** over classic Python for loop structures is that **its arithmetic operations are vectorized**. Vectorization means that an elemental arithmetic operation is automatically applied to all elements in an array. By formulating our arithmetic operations as a sequence of instructions on an array, rather than performing a set of operations for each element at a time, we can make better use of our modern central processing unit (CPU) architectures with single instruction, multiple data (SIMD) support. Furthermore, **NumPy** uses highly optimized linear algebra libraries, such as **Basic Linear Algebra Subprograms (BLAS)** and **Linear Algebra Package (LAPACK)**, that have been written in **C** or **Fortran**. Lastly, **NumPy** also allows us to write our code in a more compact and intuitive way using the basics of linear algebra, such as vector and matrix dot products.

1.7 TRAINING OF THE PERCEPTRON ALGORITHM ON THE IRIS DATASET

Observation 1.7.1 (Restriction of Dimensions)

To test our perceptron implementation, we will restrict the following analyses and examples in the remainder of this section to two feature variables (dimensions). Although the perceptron rule is not restricted to two dimensions, considering only two features—sepal length and petal length—allows us to visualize the decision regions of the trained model in a scatterplot.

Idea 1.7.1 (Restriction to Only Two Classes)

We will also only consider two flower classes, *setosa* and *versicolor*, from the Iris dataset for practical reasons (remember, the **perceptron is a binary classifier**). However, the perceptron algorithm can be extended to multi-class classification using the one-versus-all (OvA) technique.

1.7.1 OVA METHOD FOR MULTI-CLASS CLASSIFICATION

We can extend any binary classifier to a multi-class problems using the method called **OvA** or **one-versus-rest (OvR)**.

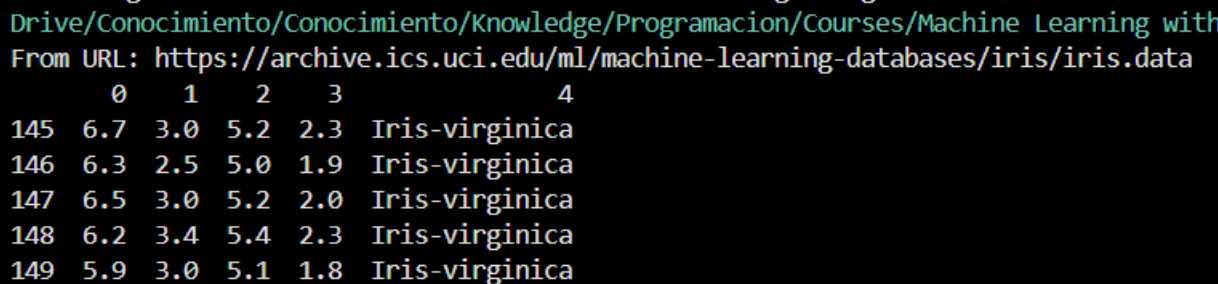
Definition 1.7.1 (One-versus-rest (OvR))

The method **OvR** allows to extend any binary classifier to a multiple classifier. To classify a new, unlabeled instance, we use our n classifiers (where n is the number of class labels) and assign the class label with the highest confidence to the instance. For the perceptron, we choose the class label associated with the largest absolute net input value.

First, we use the `pandas` library to load the Iris dataset directly from the **UCI Machine Learning Repository** into a `DataFrame` object and *print the last five lines via the tail method to confirm that the data loaded correctly*:

```
1 import os
2 import pandas as pd
3 s = 'https://archive.ics.uci.edu/ml/'\
4     'machine-learning-databases/iris/iris.data'
5 print('From URL:', s)
6 #From URL: https://archive.ics.uci.edu/ml/machine-learning-
7   databases/iris/iris.data
8 df = pd.read_csv(s,
9                 header=None,
10                encoding='utf-8')
11 print(df.tail())
```

Code 1.2: Read Dataset from UCI.



```
Drive/Conocimiento/Conocimiento/Knowledge/Programacion/Courses/Machine Learning with
From URL: https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data
   0    1    2    3    4
145 6.7  3.0  5.2  2.3 Iris-virginica
146 6.3  2.5  5.0  1.9 Iris-virginica
147 6.5  3.0  5.2  2.0 Iris-virginica
148 6.2  3.4  5.4  2.3 Iris-virginica
149 5.9  3.0  5.1  1.8 Iris-virginica
```

Figure 1.2: Tail of the Iris Dataset.

Next, we extract the first 100 class labels that correspond to the 50 `Iris-setosa` and 50 `Iris-versicolor` flowers and convert the class labels into the two integer class labels 1 (`versicolor`) and 0 (`setosa`), assigning them to a vector `y`. Similarly, we extract the first feature column (`sepal length`) and the third feature column (`petal length`) of those 100 training examples and assign them to a feature matrix `X`, which we then visualize via a two-dimensional scatterplot:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 # select setosa and versicolor
4 y = df.iloc[0:100, 4].values
5 y = np.where(y == 'Iris-setosa', 0, 1)
6 # extract sepal length and petal length
7 X = df.iloc[0:100, [0, 2]].values
8 # plot data
9 plt.scatter(X[:50, 0], X[:50, 1],
10           color='red', marker='o', label='Setosa')
11 plt.scatter(X[50:100, 0], X[50:100, 1],
12           color='blue', marker='s', label='Versicolor')
13 plt.xlabel('Sepal length [cm]')
14 plt.ylabel('Petal length [cm]')
15 plt.legend(loc='upper left')
16 plt.show()
```

Code 1.3: Extraction and Visualization of Iris Dataset.

Observation 1.7.2 (Use of the Code)

Add this code to the remaining code.

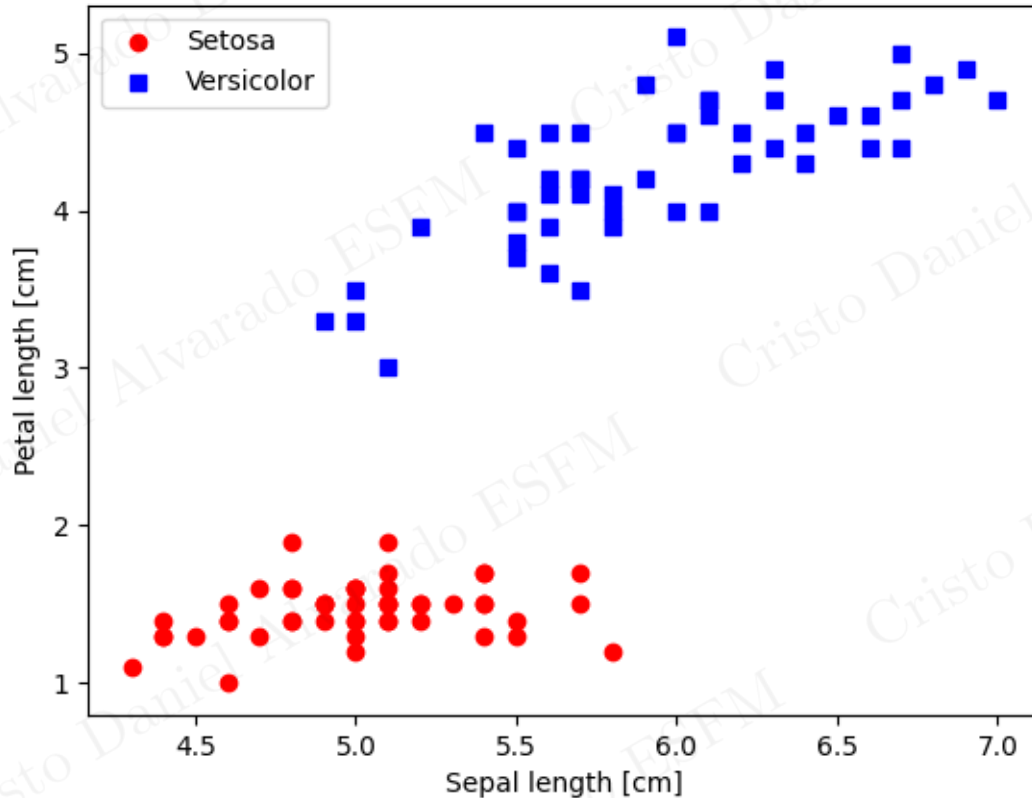


Figure 1.3: Plot of the Setosa and Versicolor on the Iris Dataset.

This plot shows the distribution of flower examples in the Iris dataset, along with the two feature axes: petal and sepal length.

Observation 1.7.3

In this subspace, a linear decision boundary should be sufficient to separate setosa from versicolor flowers.

Thus, a linear classifier such as the perceptron should be able to classify the flowers in the Iris dataset perfectly.

1.7.2 TRAINING THE ALGORITHM

We will also plot the misclassification error for each epoch to check whether the algorithm converged and found a decision boundary that separates the two Iris flower classes:

```
1 import Perceptron
2 ppn = Perceptron(eta=0.1, n_iter=10)
3 ppn.fit(X, y)
4 plt.plot(range(1, len(ppn.errors_) + 1),
5          ppn.errors_, marker='o')
6 plt.xlabel('Epochs')
7 plt.ylabel('Number of updates')
```

```
8 plt.show()
```

Code 1.4: Training the Perceptron Algorithm with the Iris Dataset

Observation 1.7.4

In this scenario, since I put Perceptron in a different Python file, I have to import it in order to be able to use it.

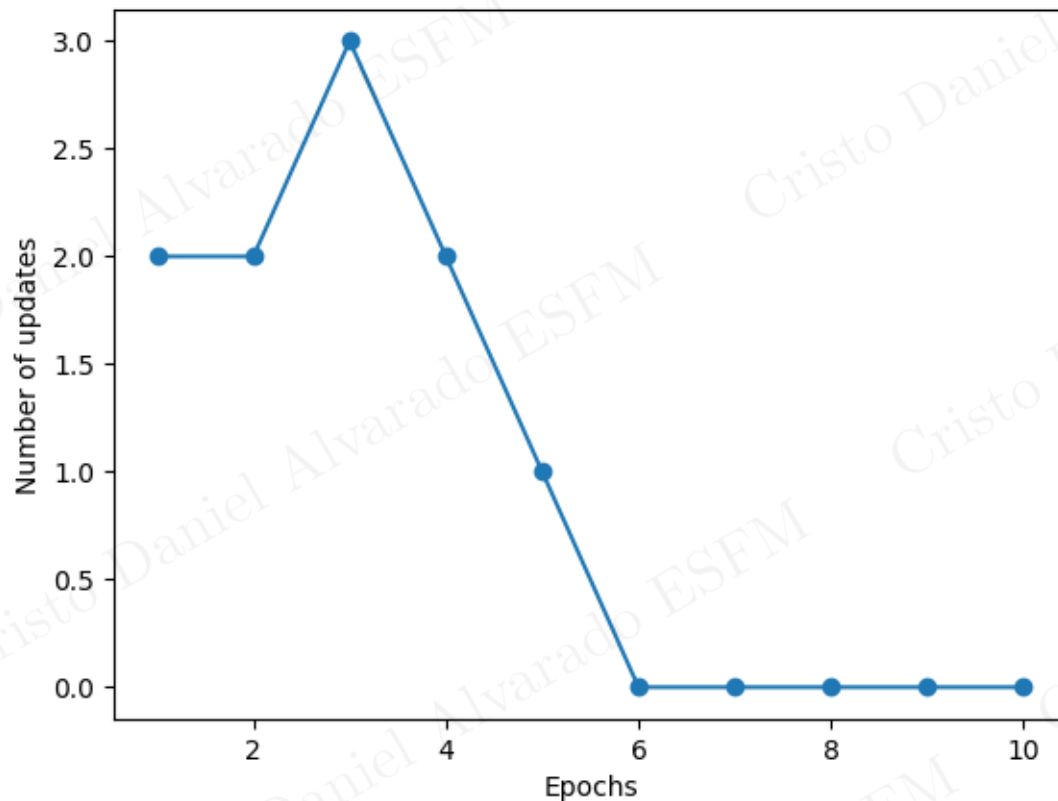


Figure 1.4: Epoch Plot and Number of Updates on the Iris Dataset.

Observation 1.7.5 (Convergence of the Perceptron)

As we can see, our **perceptron** converged after the sixth epoch and should now be able to classify the training examples perfectly.

Let's *implement a small convenience function to visualize the decision boundaries for two-dimensional datasets*:

```
1 from matplotlib.colors import ListedColormap
2
3 def plot_decision_regions(X, y, classifier, resolution=0.02):
4     # setup marker generator and color map
5     markers = ('o', 's', '^', 'v', '<')
6     colors = ('red', 'blue', 'lightgreen', 'gray', 'cyan')
7     cmap = ListedColormap(colors[:len(np.unique(y))])
8
```

```

9      # plot the decision surface
10     x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
11     x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
12     xx1, xx2 = np.meshgrid(np.arange(x1_min, x1_max, resolution)
13                             ,
14                             np.arange(x2_min, x2_max, resolution)
15                             )
16     lab = classifier.predict(np.array([xx1.ravel(), xx2.ravel()
17                                       ]).T)
18     lab = lab.reshape(xx1.shape)
19     plt.contourf(xx1, xx2, lab, alpha=0.3, cmap=cmap)
20     plt.xlim(xx1.min(), xx1.max())
21     plt.ylim(xx2.min(), xx2.max())
22
23     # plot class examples
24     for idx, cl in enumerate(np.unique(y)):
25         plt.scatter(x=X[y == cl, 0],
26                     y=X[y == cl, 1],
27                     alpha=0.8,
28                     c=colors[idx],
29                     marker=markers[idx],
30                     label=f'Class {cl}',
31                     edgecolor='black')

```

Code 1.5: Plotting Regions on the Iris Dataset

First, we define a number of *colors* and *markers* and create a colormap from the list of colors via `ListedColormap`. Then, we determine the minimum and maximum values for the two features and use those feature vectors to create a pair of grid arrays, `xx1` and `xx2`, via NumPy's `meshgrid` function. Since we trained our perceptron classifier on two feature dimensions, we flatten the grid arrays and create a matrix with the same number of columns as the Iris training subset so that we can use the `predict` method to predict the class labels `lab` of the corresponding grid points.

After reshaping the predicted class labels into a grid with the same dimensions as `xx1` and `xx2`, we draw a contour plot via Matplotlib's `contourf` function, which maps the different decision regions to different colors for each predicted class in the grid array:

```

1 plot_decision_regions(X, y, classifier=ppn)
2 plt.xlabel('Sepal length [cm]')
3 plt.ylabel('Petal length [cm]')
4 plt.legend(loc='upper left')
5 plt.show()

```

Code 1.6: Plotting the Perceptron Info on the Iris Dataset

As we can see in the plot, the perceptron learned a decision boundary that can classify all flower examples in the Iris training subset perfectly.

Observation 1.7.6 (Convergence of the Perceptron)

Although the perceptron classified the two Iris flower classes perfectly, convergence is one of its biggest challenges. Rosenblatt proved mathematically that the **perceptron learning rule converges if the two classes can be separated by a linear hyperplane**. However, if the classes cannot be separated perfectly by such a linear decision boundary, the weights will never stop updating unless we set a maximum number of epochs. Interested readers can find a summary

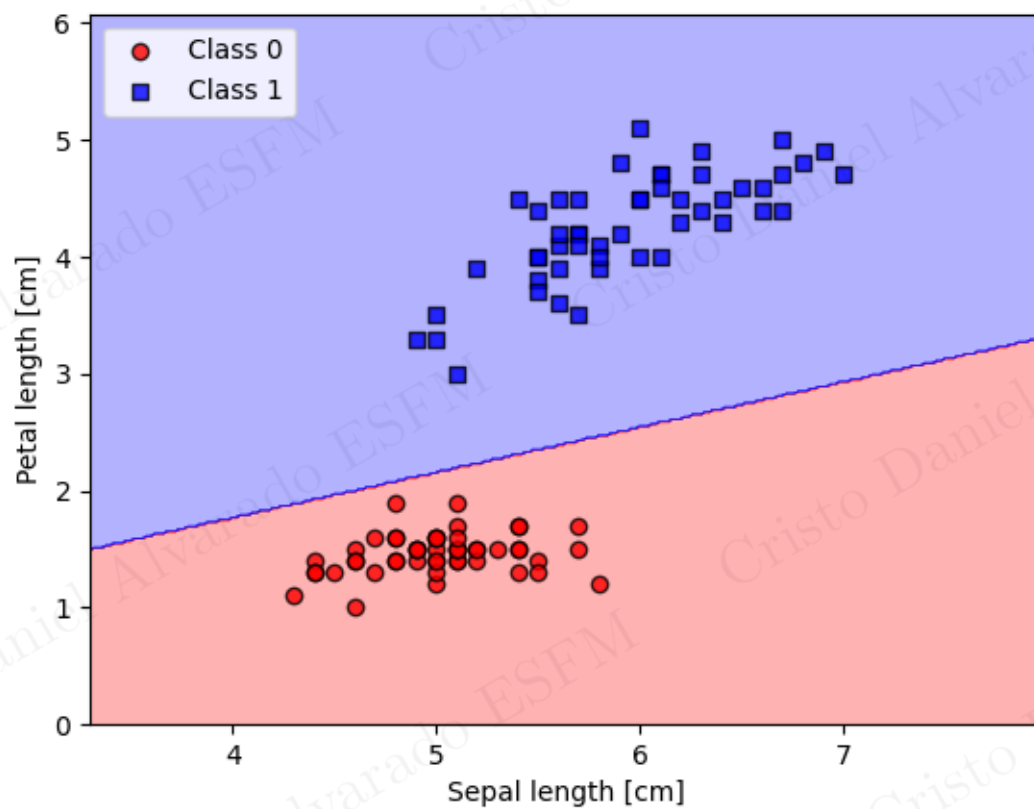


Figure 1.5: Perceptron Visual Classification on the Iris Dataset.

of the proof in the lecture notes at [Perceptron Slides](#).