# ESS2222

# Lecture 2 -  Feasibility of Learning

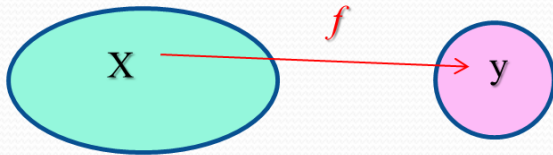*Hosein Shahnas*

*University of Toronto, Department of Earth Sciences,*

Machine learning: Learning from data
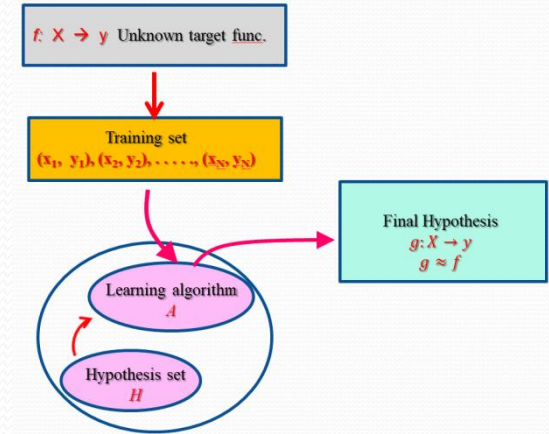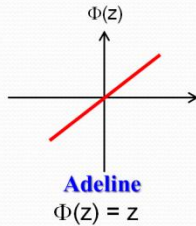
Criteria: Data ✓ , pattern ✓ , no formula ✓



$f$: X → y  Unknown target func.

Training set
$(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$

Final Hypothesis
$g: X \to y$
$g \approx f$

Learning algorithm
$A$

Hypothesis set
$H$

X  $f$  y

Learning model $\begin{cases} \text{The hypothesis set} \\ H = \{h\} \quad g \in H \\ \\ \text{The learning algorithm} \\ A \end{cases}$

$f$: X → y
$g$: X → y
$h \longrightarrow g \approx f$

$\Phi(z)$

1

□1

**PLA**
$\Phi(z) = \text{sign}(z)$

$\Phi(z)$

**Adeline**
$\Phi(z) = z$

*Perceptron:*

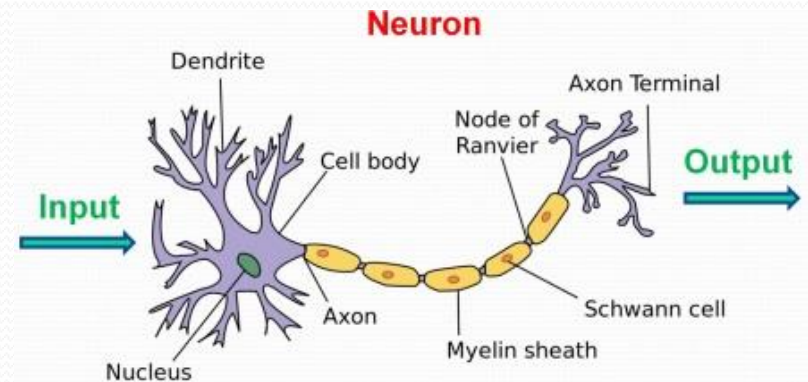$\Phi(z) = \text{sign}(z) = \begin{cases} 1, & z > 0 \\ -1, & otherwise \end{cases}$ $\hat{y} = \Phi(z)$

*Adeline:*

$\Phi(z) = z$ $\hat{y} = \begin{cases} 1, & \Phi(z) > 0 \\ -1, & otherwise \end{cases}$

$J(w) = \frac{1}{2} \Sigma_i \left( y^i - \phi(z^i) \right)^2$

1) **Supervised learning**
2) **Unsupervised learning**
3) **Reinforcement learning**

A) **Classification Problem**
B) **Regression Problem**

# Feasibility of Learning - Outline

- ❑ **Probabilistic Aspects of Learning**
- ❑ **Hoeffding's Inequality**
- ❑ **Generalization of Hoeffding's Inequality**
- ❑ **Permutation & Scaling**
- ❑ **Stochastic gradient decent**

# Is Learning Feasible?

Can we learn from a **finite** data set (samples) and generalize it (trough the mapping function) to the outside world?

The learned function (**g**)works on the **sample set**. How is the function **outside**?

The answer is the main subject of this lecture.

# A Probabilistic Situation
## An Experiment

Consider a bin with green and red marbles:
Pick N marbles independently (one by one).

Assume the fraction of the red marbles in the bin is $\mu$ and the size of bin is infinite.

Then:
P(picking a red marble) = $\mu$
P(picking a green marble) = $1 - \mu$

$\mu$ = unknown (for us) and will remain unknown.

**Bin**

**Sample (N)**

$\nu$ = fraction of red marbles

$\mu$ = Probability of red marbles

# A Probabilistic Situation
# An Experiment

How $\nu$ is related to $\mu$?
Can we say anything about $\mu$
having $\nu$?

Bin

Sample (N)

$\nu$ = fraction
of red marbles

$\mu$ = Probability
of red marbles

# A Probabilistic Situation
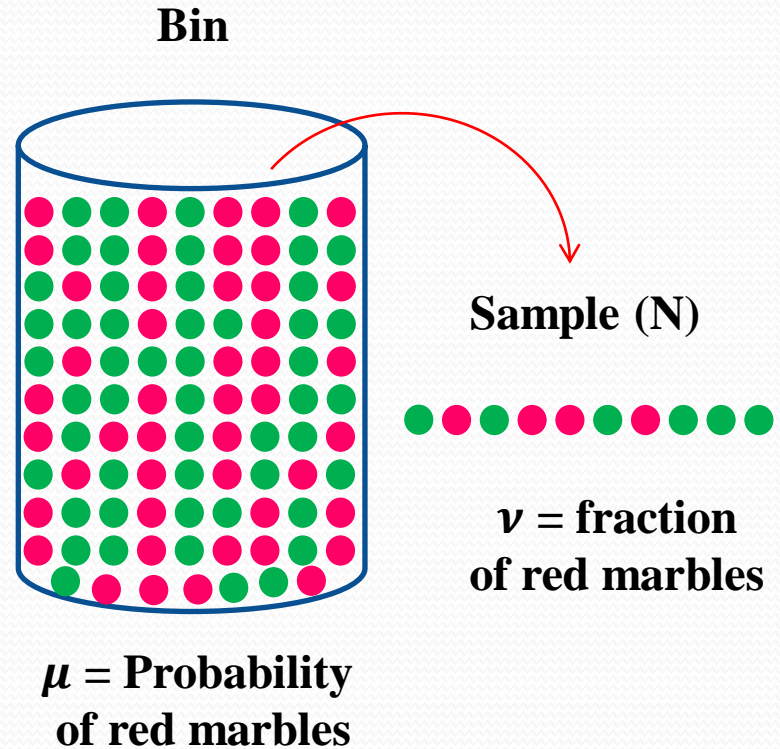# An Experiment

**How $\nu$ is related to $\mu$?
Can we say anything about $\mu$
having $\nu$?**

**No**

**Because sample can be mostly green while
bin is mostly red.** ●●●●●●●●●●●

**Bin**

**Sample (N)**

$\nu$ = fraction
of red marbles

$\mu$ = Probability
of red marbles

**How $\nu$ is related to $\mu$?**
**Can we say anything about $\mu$ having $\nu$?**

**No**
**Because sample can be mostly green while bin is mostly red.**

**And yes**
**Because if the sample is large enough, sample frequency $\nu$ is likely close to bin probability $\mu$.**

**Distinction between two answers: Possible versus probable**

**Bin**

**Sample (N)**

$\nu$ = fraction
of red marbles

$\mu$ = Probability
of red marbles

**What does $\nu$ say about $\mu$?**

**For a big sample (large N) $\nu$ is probably close to $\mu$ (within $\epsilon$).**
**Mathematically:**

$$\mathbf{P}[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$    **Hoeffding's inequality**

**bad event**

**Bound does not depend on $\mu$.**

# A Probabilistic Situation
## What does $\nu$ say about $\mu$?

**What does $\nu$ say about $\mu$?**

For a big sample **(large N)** $\nu$ is probably close to $\mu$ **(within $\epsilon$).**
Mathematically:

$$P[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$  **Hoeffding's inequality**

**bad event**

As the number of samples increases the probability of **bad event** decreases:    $P \longrightarrow 0$

However,  as $\epsilon \longrightarrow 0$,    $P \longrightarrow 2$

# A Probabilistic Situation
## What does $\nu$ say about $\mu$?

$$\mathbf{P}[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N} \qquad \forall\ \epsilon, \mathrm{N} \qquad \text{Hoeffding's inequality}$$

- ❑ So the statement   "$\mu = \nu$" is **probably approximately** correct (PAC).

- ❑ Despite P **depends on** $\mu$, the bound $(2e^{-2\epsilon^2 N})$ **does not**, which is good because $\mu$ is unknown.

- ❑ Hoeffding's inequality dictates that in order to have **lower tolerance** ($\epsilon$) we need **large number of samples** (N).

- ❑ Note that the inequality says that: $\nu \approx \mu$, because $\nu$ is affected by $\mu$ ($\mu$: the cuase, $\nu$ the effect), however, we infer $\mu \approx \nu$ due to the symmetry in the Hoeffding's inequality.

# What are unknowns?

**Bin**: The unknown quantity is $\mu$
**Learning**: The unknown is  f: x $\rightarrow$ y

**Each marble is a point in X space:  x $\in$ *X***

**Try a single (particular) hypothesis *h (an approx. to f)*:**

**Hypothesis predicts correctly**    $h(x) = f(x)$
**Hypothesis predicts wrong**    $h(x) \neq f(x)$

$\mu$

**X**

# What are unknowns?

Note that $h(x) = f(x)$ does not necessarily mean $h = f$

**Bin**: The unknown quantity is $\mu$
**Learning**: The unknown is f: x → y

**Each marble is a point in X space:** $x \in X$

**Try a single (particular) hypothesis h *(an approx. to f)*:**

**Hypothesis predicts correctly** $\quad h(x) = f(x)$
**Hypothesis predicts wrong** $\quad h(x) \neq f(x)$

$h(x) = f(x)$
$h(x) \neq f(x)$

$\mu$

*f(x)*

**X** → **y**

*h(x)*

**X** → **ŷ**

**X**

## More hypothesis:

❑ Learning is not trying a single hypothesis $h$. Trying a single hypothesis is just a **verification**.

❑ Before generalizing $\nu$ to $\mu$, we have to search for the set of hypothesis and find the best hypothesis.



Wrong prediction
Correct prediction

# Generalization

## More hypothesis:

❑ Learning is not trying a single hypothesis $h$. Trying a single hypothesis is just a **verification**.

❑ Before generalizing $\nu$ to $\mu$, we have to search for the set of hypothesis and find the best hypothesis.



$h_1$ $\qquad\qquad$ $h_2$ $\qquad\qquad\qquad\qquad$ $h_k$

$\mu_1$ $\qquad\qquad$ $\mu_2$ **Generalize to bin** $\qquad\qquad$ $\mu_k$

. . . . . . . . .

**good**

$\nu_1$ $\qquad\qquad\qquad$ $\nu_2$ $\qquad\qquad\qquad\qquad$ $\nu_k$

$g = h_2$

● **Wrong prediction**
● **Correct prediction**

15

# Generalization

- **Red marble**
- **Green marble**

$\mu$: **unknown (exact value)**

# Generalization

- **Red marble**
- **Green marble**

  **$\mu$: unknown (exact value)**

- **Wrong prediction**
- **Correct prediction**

  **Try to minimize the number of wrong predictions by search**

  **$\mu \approx \nu$ with the bounded probability P**

# Errors

$\nu$ and $\mu$ depend on $h$. We introduce the error rates corresponding to $\nu$ and $\mu$.

$\nu$ (in-sample) $\longrightarrow$ $E_{in}(h)$

$\mu$(out-sample) $\longrightarrow$ $E_{out}(h)$

**Then the Hoeffding's inequality becomes:**

$$P[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$$\longrightarrow P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$E_{out}(h)$



$E_{in}(h)$

# Errors

$\nu$ and $\mu$ depend on $h$. We introduce the error rates corresponding to $\nu$ and $\mu$.

$\nu$ (in-sample) $\longrightarrow$ $E_{in}(h)$
$\mu$(out-sample) $\longrightarrow$ $E_{out}(h)$

**Then the Hoeffding's inequality becomes:**

$E_{out}(h)$

$$P[|\nu - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

$\longrightarrow$ $P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$

The probability that in-sample performance deviates from out-sample performance by more than $\epsilon$, is less than $2e^{-2\epsilon^2 N}$.

$E_{in}(h) = \frac{1}{N} \sum_{1}^{N} [\![ h(x_n) \neq f(x_n) ]\!]$    for classification prob.
$[\![ h(x_n) \neq f(x_n) ]\!] = 1$   if $h(x) \neq f(x)$,    $= 0$ otherwise
$E_{out}(h) = P [\![ h(x) \neq f(x) ]\!]$

$E_{in}(h)$

$h_1$ $h_2$ $h_k$

$\boldsymbol{\mu_1}$ $\boldsymbol{\mu_2}$ $\boldsymbol{\mu_k}$

. . . . . . . . .

$\boldsymbol{\nu_1}$ $\boldsymbol{\nu_2}$ $\boldsymbol{\nu_k}$

**Wrong prediction**
**Correct prediction**

# Errors

$h_1$

$E_{out}(h_1)$

$E_{in}(h_1)$

$h_2$

$E_{out}(h_2)$

$E_{in}(h_2)$

$h_k$

$E_{out}(h_k)$

$E_{in}(h_k)$

● Wrong prediction
● Correct prediction

21

Hoeffding's inequality **doesn't apply** to multiple bin.

**Consider coin analogy**

**Prob. 1:** **If you toss a fair coin 10 times, what is the probability that you get 10 heads?**

**Sol.:** $P = \dfrac{\# \ of \ events}{total \ \# \ of \ events} = \dfrac{1}{2^{10}}$

$P(k) = \dbinom{n}{k} p^k q^{n-k},$ $\qquad C^n_{\ k} = \dbinom{n}{k} = \dfrac{n!}{(n-k)!k!} = \dfrac{P^n_{\ k}}{k!},$

$n = 10,$ # of flips $\qquad where \quad P^n_{\ k} = \dfrac{n!}{(n-k)!}$ *(k-permutations of n)*

$k = 10,$ # of heads

$p = 1/2.$ Prob. of success

$q = 1 - p = 1/2,$ Prob. of failure

$P(10) = \dfrac{1}{2^{10}} \approx 0.001 \ \rightarrow \ P(10) \approx 0.1\%$

| A | B | C |
|---|---|---|
| A | C | B |
| B | A | C |
| B | C | A |
| C | A | B |
| C | B | A |

3!

| A | B |   |   |
|---|---|---|---|
|   | A | B |   |
|   |   | A | B |
| A |   |   | B |
| A |   | B |   |
|   | A |   | B |
| B | A |   |   |
|   | B | A |   |
|   |   | B | A |
| B |   |   | A |
| B |   | A |   |
|   | B |   | A |

$P^4_2$

# A Problem!

Hoeffding's inequality **doesn't apply** to multiple bin.

**Consider coin analogy**

**Prob. 2:** If you toss **1000** fair coins **10** times each, what is the probability that some coin gets **10** heads?

**Sol.:** Two steps:

$$P_1 = \binom{10}{0} \left(\frac{1}{2}\right)^{10} \left(1 - \frac{1}{2}\right)^0 = \frac{1}{2^{10}}, \qquad P_2 = \binom{1000}{1} (P_1)^1 (1 - P_1)^{999},$$

$$P_2 = \binom{1000}{1} \left(\frac{1}{2^{10}}\right)^1 \left(1 - \frac{1}{2^{10}}\right)^{999} \approx 0.368, \qquad \text{exactly 1 out of 1000 coins}$$
shows up 10 heads on 10 tosses

$$P = 1 - \binom{1000}{0} \left(\frac{1}{2^{10}}\right)^0 \left(1 - \frac{1}{2^{10}}\right)^{1000} \approx 0.624, \qquad \text{at least 1 out of 1000 coins}$$
shows up 10 heads on 10 tosses

Prob. of getting no 10-heads

Now suppose $\mu = 0.5$ (i.e. half of marbles green and half of them red). The same probability we have in flipping a fair coin.



This is not reflecting the reality (the real probability),

# A Simple Solution

Getting 10 heads in tossing a fair coin 10 times does not reflect the reality. But if we try too hard, something bad will happen somewhere.

In mathematical language: Hoeffding's inequality applies for a single experiment (not multiple).

*P(g)* for the final hypothesis is less than any *P(h)* and therefore less than the union of them:

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq \quad P[ \quad |E_{in}(h_1) - E_{out}(h_1)| > \epsilon$$

or
$$|E_{in}(h_2) - E_{out}(h_2)| > \epsilon$$

**g: best of h** $\quad \cdots\cdots\cdots$

or
$$|E_{in}(h_k) - E_{out}(h_k)| > \epsilon \quad ]$$

$$\leq \sum_{m=1}^{k} P[|E_{in}(h_m) - E_{out}(h_m)| > \epsilon]$$

**Union bound**

This explains that why we got high probability (~0.63) for that bad event.

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq \sum_{m=1}^{k} P[|E_{in}(h_m) - E_{out}(h_m)| > \epsilon]$$

$$\leq \sum_{m=1}^{k} 2e^{-2\epsilon^2 N}$$

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

The factor M at the right had side of the Hoeffding's inequality increases the probability bound which is not good, however, as we will see later, the inequality can be improved.

# Learning Diagram

# Scaling for Optimal Performance

**Feature scaling:**

1) Min-max normalization: $\quad x' = \dfrac{x - x_{min}}{x_{max} - x_{min}}$

2) Mean normalization: $\quad x' = \dfrac{x - \mu}{x_{max} - x_{min}}$

3) Normalization to a range ($r_{min}$, $r_{max}$): $x' = \dfrac{x - x_{min}}{x_{max} - x_{min}}(r_{max} - r_{min}) + r_{min}$

4) Standardization: gives data the property of a standard normal distribution.

$x' = \dfrac{x - \mu}{\sigma}, \quad \mu : mean, \quad \sigma : standard\ diviation$

The set of new data ($X'$) has zero mean and unit variance, but not bounded.

# Stochastic Gradient Decent Method

**1 - Gradient decent (GD) method ((Batch gradient descent):**

$w = w + \Delta w$

$\Delta w = -\eta\, \Delta J(w)$

$\Delta wj = -\eta\, \dfrac{\partial J}{\partial w_j} = \eta \sum_i (y^i - \phi(z^i))\, x^i_j$      Based on all samples

i: Sample index,    j: Feature index

**2 – Stochastic gradient decent (SGD) method
(an approx. for GD for large data):**

$\Delta wj = -\eta\, \dfrac{\partial J}{\partial w_j} = \eta\, (y^i - \phi(z^i))\, x^i_j$      Based on random samples

$\eta$ : Variable learning rate (decreasing with iteration)

In SGD **a)** the error is noisier than in GD (because $\Delta w's$ are based on single samples),
**b)** The convergence is faster than GD (more frequent updates), **c)** the local minima can be escaped faster which is good, **d)** to get accurate results the samples must be shuffled,

```python
'''
Simultaneous permutation of class labels and features
Hosein Shahnas
'''
import sys
import os.path
#from sklearn.preprocessing import scale
import os
import numpy as np
#===================================================================== import data
save_path = os.path.dirname(os.path.abspath(__file__))
print(save_path)

name_of_file = 'Features'                           # filename for features
completeName = os.path.join(save_path, name_of_file+".dat")    # complete filename (path included)
Feature_data = np.loadtxt(completeName)             # load the file as text

name_of_file = 'Classe_Labels'                      # filename for class labels
completeName = os.path.join(save_path, name_of_file+".dat")
Target_data = np.loadtxt(completeName)

print ('Feature_data.shape = ', Feature_data.shape)
print ('Target_data.shape = ', Target_data.shape)

np.unique(Target_data)
print ('np.unique(Target_data) = ', np.unique(Target_data))   #Returns the sorted unique elements of an array.


y_size = Target_data.size               # get the size of target array
X_size = Feature_data.size              # get the size of features array
Feature_size = int(X_size/y_size)       # get the number of features

print ('y_size = ', y_size)
print ('X_size = ', X_size)
print ('Feature_size = ', Feature_size)

#===================================================================== import data
```

# Permutation

```python
38
39#============================================================== shuffle data
40perm = np.random.permutation(Target_data.size)    # get the index numbers for random shuffle (permutation )
41print ('perm = ', perm)
42
43Feature_data = Feature_data[perm]                  # based on perm shuffle features
44Target_data = Target_data[perm]                    # based on perm shuffle targets
45#============================================================== shuffle data
46
47#============================================================== write data
48name_of_file = 'Perm_Classe_Labels_Features'
49completeName = os.path.join(save_path, name_of_file+".dat")
50file1 = open(completeName, "w")
51for i in range(0,Target_data.size):
52    file1.write("%5i %5i  " % (perm [i], Target_data [i]))
53    for j in range(0,Feature_size):
54        file1.write(" %20.12e " % (Feature_data [i,j]))
55    file1.write(" \n " )                            # go to the next line
56file1.close();
57
58name_of_file = 'Perm_Classe_Labels'
59completeName = os.path.join(save_path, name_of_file+".dat")
60file1 = open(completeName, "w")
61for i in range(0,Target_data.size):
62    file1.write("%5i \n" % (Target_data [i]))
63file1.close();
64
65name_of_file = 'Perm_Features'
66completeName = os.path.join(save_path, name_of_file+".dat")
67file1 = open(completeName, "w")
68for i in range(0,Target_data.size):
69    for j in range(0,Feature_size):
70        file1.write(" %20.12e " % ( Feature_data [i,j]))
71    file1.write(" \n " )
72file1.close();
73#============================================================== write data
74
75#sys.exit('Program stopped here')
76
```

31

# Scaling

```python
1 '''
2 Scaling of features
3 Hosein Shahnas
4 '''
5 import sys
6 import os.path
7 #from sklearn.preprocessing import scale
8 from sklearn import preprocessing
9 import os
10 import numpy as np
11 #========================================================================= import data
12
13 save_path = os.path.dirname(os.path.abspath(__file__))
14 print(save_path)
15
16 name_of_file = 'Perm_Features'
17 completeName = os.path.join(save_path, name_of_file+".dat")
18 Fearures = np.loadtxt(completeName)
19
20 name_of_file = 'Perm_Classe_Labels'
21 completeName = os.path.join(save_path, name_of_file+".dat")
22 Class_Labels = np.loadtxt(completeName)
23
24 print ('Fearures.shape = ', Fearures.shape)
25 print ('Class_Labels.shape = ', Class_Labels.shape)
26
27 np.unique(Class_Labels)
28 print ('np.unique(Class_Labels) = ', np.unique(Class_Labels))    #Returns the sorted unique elements of an array.
29 y_size = Class_Labels.size
30 X_size = Fearures.size
31 Feature_size = int(X_size/y_size)
32 #========================================================================= import data
33
```

# Scaling

```python
33
34 #=========================================================== scale
35 min_max_scaler = preprocessing.MinMaxScaler()
36 Fearures_s = min_max_scaler.fit_transform(Fearures)
37 #=========================================================== scale
38
39 #=========================================================== write data
40 name_of_file = 'Perm_Classe_Labels_Features_Scaled'
41 completeName = os.path.join(save_path, name_of_file+".dat")
42 file1 = open(completeName, "w")
43 for i in range(0,Class_Labels.size):
44     file1.write("%5i %5i  " % (i, Class_Labels [i]))
45     for j in range(0,Feature_size):
46         file1.write(" %20.12e " % (Fearures_s [i,j]))
47     file1.write(" \n " )
48 file1.close();
49
50
51 name_of_file = 'Perm_Features_Scaled'
52 completeName = os.path.join(save_path, name_of_file+".dat")
53 file1 = open(completeName, "w")
54 for i in range(0,Class_Labels.size):
55     for j in range(0,Feature_size):
56         file1.write(" %20.12e " % ( Fearures_s [i,j]))
57     file1.write(" \n " )
58 file1.close();
59 #=========================================================== write data
60
61 #sys.exit('Program stopped here')
62
```

# Perceptron, Adaline-GD, Adaline-SGD Algorithms

```
 1# Sebastian Raschka, 2015 (http://sebastianraschka.com)
 2# Python Machine Learning - Code Examples
 3#
 4# Chapter 2 - Training Machine Learning Algorithms for Classification
 5#
 6# S. Raschka. Python Machine Learning. Packt Publishing Ltd., 2015.
 7# GitHub Repo: https://github.com/rasbt/python-machine-learning-book
 8#
 9# License: MIT
10# https://github.com/rasbt/python-machine-learning-book/blob/master/LICENSE.txt
11
12''' Iris Classification problem - Modified version'''
13#===============================================================
14import sys
15import numpy as np
16import pandas as pd
17import matplotlib.pyplot as plt
18from matplotlib.colors import ListedColormap
19#===============================================================
```

## Perceptron, Adaline-GD, Adaline-SGD Algorithms

```python
20#=============================================================              Perceptron Algorithm
21class Perceptron(object):
22    """Perceptron classifier.
23
24    Parameters
25    ------------
26    eta : float
27        Learning rate (between 0.0 and 1.0)
28    n_iter : int
29        Passes over the training dataset.
30
31    Attributes
32    -----------
33    w_ : 1d-array
34        Weights after fitting.
35    errors_ : list
36        Number of misclassifications (updates) in each epoch.
37
38    """
39    def __init__(self, eta=0.01, n_iter=10):
40        self.eta = eta
41        self.n_iter = n_iter
42
43    def fit(self, X, y):
44        """Fit training data.
45
46        Parameters
47        ----------
48        X : {array-like}, shape = [n_samples, n_features]
49            Training vectors, where n_samples is the number of samples and
50            n_features is the number of features.
51        y : array-like, shape = [n_samples]
52            Target values.
53
54        Returns
55        -------
56        self : object
57
58        """
```

```python
59         self.w_ = np.zeros(1 + X.shape[1])                    # dimension of w array = number of features
60         self.errors_ = []
61
62         for _ in range(self.n_iter):
63             errors = 0
64             for xi, target in zip(X, y):
65                 update = self.eta * (target - self.predict(xi))
66                 self.w_[1:] += update * xi
67                 self.w_[0] += update
68                 errors += int(update != 0.0)
69             self.errors_.append(errors)
70         return self
71
72     def net_input(self, X):
73         """Calculate net input"""
74         return np.dot(X, self.w_[1:]) + self.w_[0]
75
76     def predict(self, X):
77         """Return class label after unit step"""
78         return np.where(self.net_input(X) >= 0.0, 1, -1)
79 #=============================================================  Perceptron Algorithm
80
81 #=============================================================  import iris data from web source
82 print(50 * '=')
83 print('Section: Training a perceptron model on the Iris dataset')
84 print(50 * '-')
85
86 df = pd.read_csv('https://archive.ics.uci.edu/ml/'
87                  'machine-learning-databases/iris/iris.data', header=None)
88 print(df.tail())
89
90 #=============================================================  import iris data from web source
91
```

```python
92 #================================================================ cunstruct arrays from raw data
93 print(50 * '=')
94 print('Plotting the Iris data')
95 print(50 * '-')
96
97 # select setosa and versicolor
98 y = df.iloc[0:100, 4].values
99 y = np.where(y == 'Iris-setosa', -1, 1)
100
101 # extract sepal length and petal length
102 X = df.iloc[0:100, [0, 2]].values
103 #================================================================ cunstruct arrays from raw data
104
105 #================================================================ scatter plot for Setosa and Versicolor iris
106 # plot data
107 plt.scatter(X[:50, 0], X[:50, 1],
108             color='red', marker='o', label='setosa')
109 plt.scatter(X[50:100, 0], X[50:100, 1],
110             color='blue', marker='x', label='versicolor')
111
112 plt.xlabel('sepal length [cm]')
113 plt.ylabel('petal length [cm]')
114 plt.legend(loc='upper left')
115
116 # plt.tight_layout()
117 # plt.savefig('./images/02_06.png', dpi=300)
118 plt.show()
119 #================================================================ scatter plot for Setosa and Versicolor iris
120
```

```python
121#============================================================= Apply Perceptron Algorithm
122print(50 * '=')
123print('Training the perceptron model')
124print(50 * '-')
125
126ppn = Perceptron(eta=0.1, n_iter=30)
127'''
128ppn = Perceptron()
129ppn = Perceptron(eta=0.1, n_iter=5)
130'''
131ppn.fit(X, y)
132#============================================================= Apply Perceptron Algorithm
133
134#============================================================= error plot for Perceptron
135print('len(ppn.errors_) = ', len(ppn.errors_))
136#plt.plot(range(1, len(ppn.errors_) + 1), ppn.errors_, marker='o')
137
138plt.plot(range(1,31), ppn.errors_, marker='o')
139
140plt.xlabel('Epochs')
141plt.ylabel('Number of misclassifications')
142
143# plt.tight_layout()
144# plt.savefig('./perceptron_1.png', dpi=300)
145plt.show()
146#============================================================= error plot for Perceptron
147
148#============================================================= make the grid for decision plot
149resolution=0.02
150x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
151x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
152xx1, xx2 = np.meshgrid(np.arange(x1_min, x1_max, resolution),
153                       np.arange(x2_min, x2_max, resolution))
154print('xx1.shape = ', xx1.shape)
155xx11 = xx1.ravel()
156print('xx11.shape = ', xx11.shape)
157#============================================================= make the grid for decision plot
158
```

38

```python
158
159#================================================= function for plotting decision region
160print(50 * '=')
161print('A function for plotting decision regions')
162print(50 * '-')
163
164def plot_decision_regions(X, y, classifier, resolution=0.02):
165
166    # setup marker generator and color map
167    markers = ('s', 'x', 'o', '^', 'v')
168    colors = ('red', 'blue', 'lightgreen', 'gray', 'cyan')
169    cmap = ListedColormap(colors[:len(np.unique(y))])
170
171    # plot the decision surface
172    x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
173    x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
174    xx1, xx2 = np.meshgrid(np.arange(x1_min, x1_max, resolution),
175                           np.arange(x2_min, x2_max, resolution))
176    print('xx1.shape = ', xx1.shape)
177    print('xx2.shape = ', xx2.shape)
178    '''
179    xx1[i,j]: the first feature at each grid point
180    xx2[i,j]: the second feature at each grid point
181    '''
182    features = np.array([xx1.ravel(), xx2.ravel()]).T
183    print('xx1.ravel().shape = ', xx1.ravel().shape)
184    print('xx2.ravel().shape = ', xx2.ravel().shape)
185    print('features.shape = ', features.shape)
186    '''
187    features has a shape of (n,m), where n is the number of smaples and m is the number of features
188    '''
```
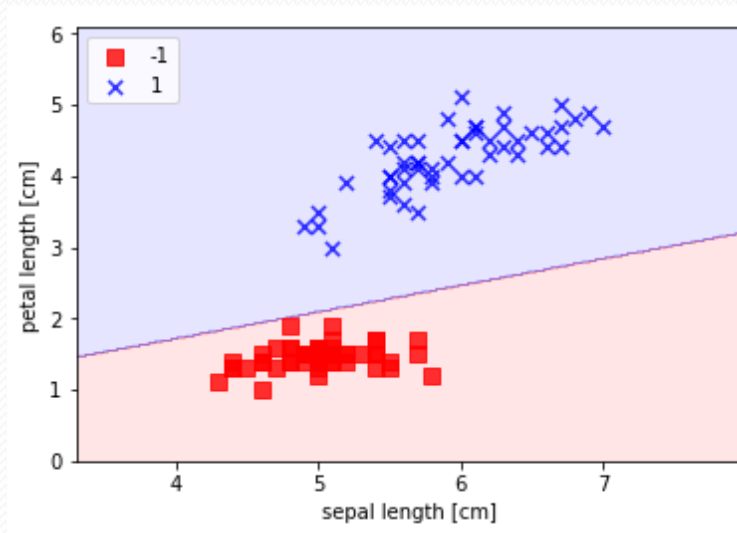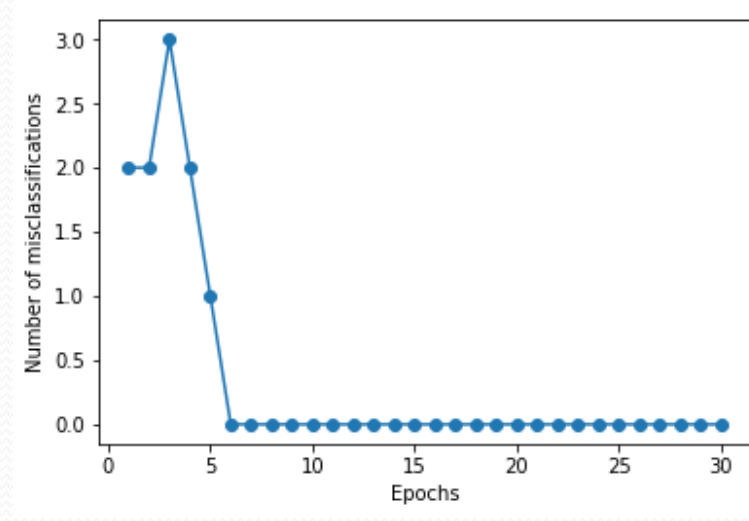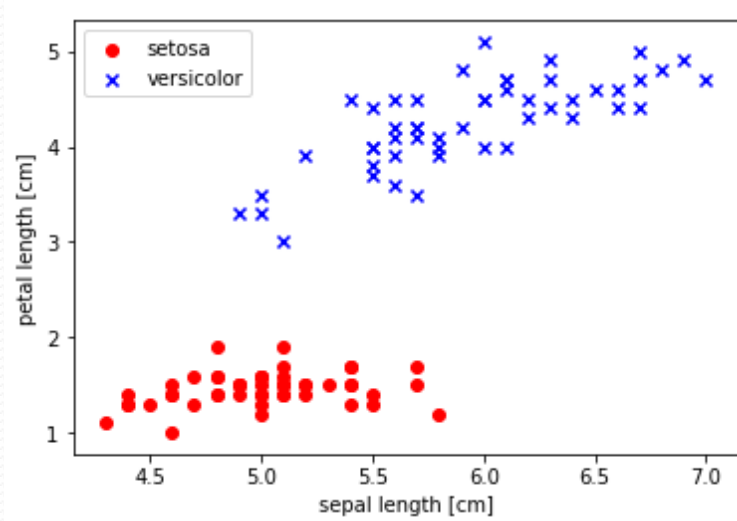
```
188     '''
189     Z = classifier.predict(features)
190
191     print('Z.shape1 = ', Z.shape)
192     Z = Z.reshape(xx1.shape)
193     print('Z.shape2 = ', Z.shape)
194
195     plt.contourf(xx1, xx2, Z, alpha=0.1, cmap=cmap)  # alpha: color level, cmap: color map
196     plt.xlim(xx1.min(), xx1.max())
197     plt.ylim(xx2.min(), xx2.max())
198
199     print('y.shape = ', y.shape)
200     print('np.unique(y) = ', np.unique(y))    # Returns the sorted unique elements of an array.
201     print('np.unique(y).shape = ', np.unique(y).shape)
202     #print('enumerate(np.unique(y)) = ', enumerate(-1,1)
203
204     # plot class samples                          # x=X[y == cl, 0]: x-componentx for which y = cl(-1 or 1)
205                                                    # y=X[y == cl, 1]: y-componentx for which y = cl (-1 or 1)
206     for idx, cl in enumerate(np.unique(y)):        # starting from idx = 0, cl takes the value of np.unique(y) members,
207         plt.scatter(x=X[y == cl, 0], y=X[y == cl, 1], # i.e., ixx = 1, cl = -1 and idx = 1, cl =1
208                     alpha=0.8, c=cmap(idx), s=50,      # s: size of the marker
209                     marker=markers[idx], label=cl)     # markers from the list
210 '''
211 x0=X[y == -1, 0]
212 print('x0 = ', x0)
213 print()
214 x0=X[y == 1, 0]
215 print('x0 = ', x0)
216 print()
217 x0=X[y == -1, 1]
218 print('x0 = ', x0)
219 print()
220 x0=X[y == 1, 1]
221 print('x0 = ', x0)
222 '''
223 #=========================================================== function for plotting decision region
224
```

# Perceptron, Adaline-GD, Adaline-SGD Algorithms

```
224
225#=========================================================== plotting decision region for Perceptron
226plot_decision_regions(X, y, classifier=ppn) # call for plot
227plt.xlabel('sepal length [cm]')
228plt.ylabel('petal length [cm]')
229plt.legend(loc='upper left')
230
231# plt.tight_layout()
232# plt.savefig('./perceptron_2.png', dpi=300)
233plt.show()
234print('===============================End of the Perceptron Algorithm')
235print()
236print()
237#=========================================================== plotting decision region for Perceptron
238
```

```
239 #================================================================ adaline gradient descent (GD) algorithm
240 print(50 * '=')
241 print('Implementing an adaptive linear neuron in Python (GD)')
242 print(50 * '-')
243
244
245 class AdalineGD(object):
246     """ADAptive LInear NEuron classifier.
247
248     Parameters
249     ------------
250     eta : float
251         Learning rate (between 0.0 and 1.0)
252     n_iter : int
253         Passes over the training dataset.
254
255     Attributes
256     -----------
257     w_ : 1d-array
258         Weights after fitting.
259     cost_ : list
260         Sum-of-squares cost function value in each epoch.
261
262     """
263     def __init__(self, eta=0.01, n_iter=50):
264         self.eta = eta
265         self.n_iter = n_iter
266
267     def fit(self, X, y):
```

```
268         """ Fit training data.
269
270         Parameters
271         ----------
272         X : {array-like}, shape = [n_samples, n_features]
273             Training vectors, where n_samples is the number of samples and
274             n_features is the number of features.
275         y : array-like, shape = [n_samples]
276             Target values.
277
278         Returns
279         -------
280         self : object
281
282         """
283         self.w_ = np.zeros(1 + X.shape[1])          # dimension of w array = number of features + 1
284         self.cost_ = []
285
286         for i in range(self.n_iter):
287             output = self.net_input(X)
288             errors = (y - output)
289             self.w_[1:] += self.eta * X.T.dot(errors)
290             self.w_[0] += self.eta * errors.sum()
291             cost = (errors**2).sum() / 2.0
292             self.cost_.append(cost)
293         return self
294
```
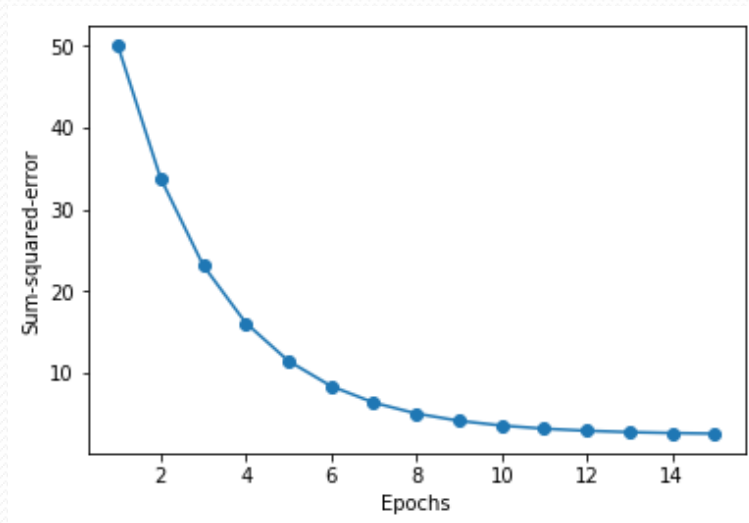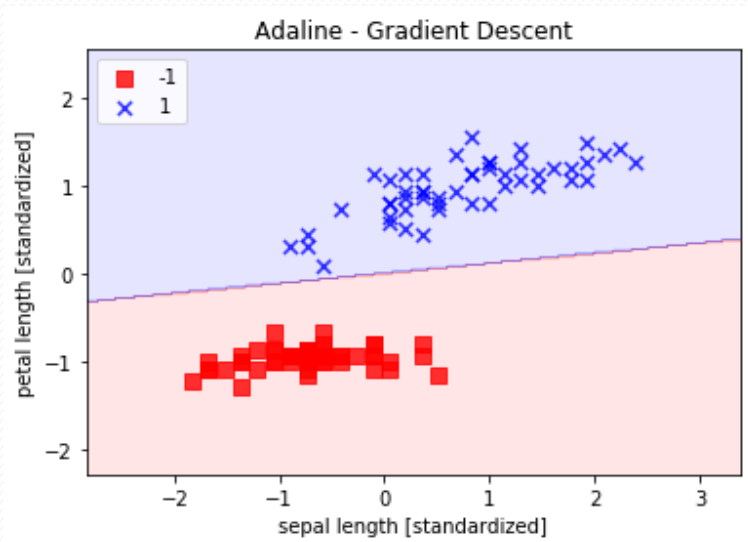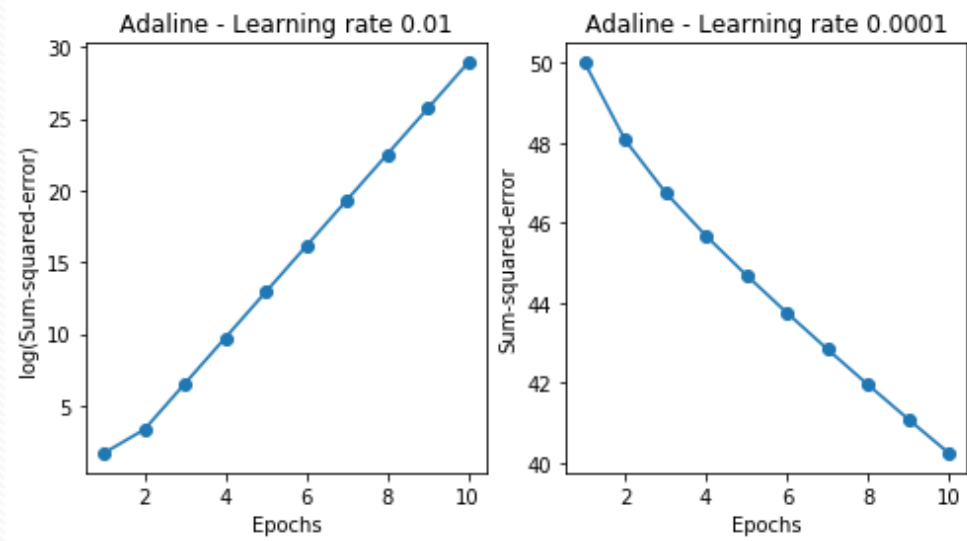
```python
295      def net_input(self, X):
296          """Calculate net input"""
297          return np.dot(X, self.w_[1:]) + self.w_[0]
298
299      def activation(self, X):
300          """Compute linear activation"""
301          return self.net_input(X)
302
303      def predict(self, X):
304          """Return class label after unit step"""
305          return np.where(self.activation(X) >= 0.0, 1, -1)
306
307 #===========================================================  adaline gradient descent (GD) algorithm
308
309 fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(8, 4))  # figure with sub-plots (one row, two columns)
310
311 #===========================================================  call Adaline-GD Learning algorithm with large eta
312 ada1 = AdalineGD(n_iter=10, eta=0.01).fit(X, y)
313 ax[0].plot(range(1, len(ada1.cost_) + 1), np.log10(ada1.cost_), marker='o')
314 ax[0].set_xlabel('Epochs')
315 ax[0].set_ylabel('log(Sum-squared-error)')
316 ax[0].set_title('Adaline - Learning rate 0.01')
317 #===========================================================  call Adaline-GD Learning algorithm with large eta
318
319 #===========================================================  call Adaline-GD Learning algorithm with small eta
320 ada2 = AdalineGD(n_iter=10, eta=0.0001).fit(X, y)
321 ax[1].plot(range(1, len(ada2.cost_) + 1), ada2.cost_, marker='o')
322 ax[1].set_xlabel('Epochs')
323 ax[1].set_ylabel('Sum-squared-error')
324 ax[1].set_title('Adaline - Learning rate 0.0001')
325 #===========================================================  call Adaline-GD Learning algorithm with small eta
326
327 #===========================================================  plot the errors for large and small eta(learning rate)
328 # plt.tight_layout()                                        Adaline-GD
329 # plt.savefig('./adaline_1.png', dpi=300)
330 plt.show()                          # show the figure
331 #===========================================================  plot the errors for large and small eta(learning rate)
332
```

```python
333#============================================================= scale the data using standardization
334print('standardize features')
335X_std = np.copy(X)        # copies x in X_std
336X_std[:, 0] = (X[:, 0] - X[:, 0].mean()) / X[:, 0].std()
337X_std[:, 1] = (X[:, 1] - X[:, 1].mean()) / X[:, 1].std()
338#============================================================= scale the data using standardization
339
340#============================================================= train using Adaline-GD
341ada = AdalineGD(n_iter=15, eta=0.01)
342ada.fit(X_std, y)
343#============================================================= train using Adaline-GD
344
345#============================================================= plot the decision regions and the errors for Adaline-GD
346plot_decision_regions(X_std, y, classifier=ada)
347plt.title('Adaline - Gradient Descent')
348plt.xlabel('sepal length [standardized]')
349plt.ylabel('petal length [standardized]')
350plt.legend(loc='upper left')
351# plt.tight_layout()
352# plt.savefig('./adaline_2.png', dpi=300)
353plt.show()
354
355plt.plot(range(1, len(ada.cost_) + 1), ada.cost_, marker='o')
356plt.xlabel('Epochs')
357plt.ylabel('Sum-squared-error')
358
359# plt.tight_layout()
360# plt.savefig('./adaline_3.png', dpi=300)
361plt.show()
362
363print('===============================End of the Adaline-GD Algorithm')
364print()
365print()
366#============================================================= plot the decision regions and the errors for Adaline-GD
```

```python
368 print(50 * '=')
369 print('Large scale machine learning and stochastic gradient descent (SGD)')
370 print(50 * '-')
371
372 class AdalineSGD(object):
373     """ADAptive LInear NEuron classifier.
374
375     Parameters
376     ------------
377     eta : float
378         Learning rate (between 0.0 and 1.0)
379     n_iter : int
380         Passes over the training dataset.
381
382     Attributes
383     -----------
384     w_ : 1d-array
385         Weights after fitting.
386     cost_ : list
387         Sum-of-squares cost function value averaged over all
388         training samples in each epoch.
389     shuffle : bool (default: True)
390         Shuffles training data every epoch if True to prevent cycles.
391     random_state : int (default: None)
392         Set random state for shuffling and initializing the weights.
393
394     """
395     def __init__(self, eta=0.01, n_iter=10, shuffle=True, random_state=None):
396         self.eta = eta
397         self.n_iter = n_iter
398         self.w_initialized = False
399         self.shuffle = shuffle
400
401         if random_state:
402             np.random.seed(random_state)
403
```

```python
404    def fit(self, X, y):
405        """ Fit training data.
406
407        Parameters
408        ----------
409        X : {array-like}, shape = [n_samples, n_features]
410            Training vectors, where n_samples is the number of samples and
411            n_features is the number of features.
412        y : array-like, shape = [n_samples]
413            Target values.
414
415        Returns
416        -------
417        self : object
418
419        """
420                                                # initialize the weight factors using _initialize_weights
421        self._initialize_weights(X.shape[1]) # dimension of w array = number of features (exclude w0)
422        self.cost_ = []
423        for i in range(self.n_iter):
424            if self.shuffle:                    # if shuffle = true then shuffle data by the defined function _shuffle
425                X, y = self._shuffle(X, y)
426            cost = []                           # initiate cost array with unknown diimension
427            for xi, target in zip(X, y):
428                cost.append(self._update_weights(xi, target)) # find the
429            avg_cost = sum(cost) / len(y)
430            self.cost_.append(avg_cost)
431        return self
432
```

```python
def partial_fit(self, X, y):
    """Fit training data without reinitializing the weights"""
    if not self.w_initialized:
        self._initialize_weights(X.shape[1])
    if y.ravel().shape[0] > 1:
        for xi, target in zip(X, y):
            self._update_weights(xi, target)
    else:
        self._update_weights(X, y)
    return self

def _shuffle(self, X, y):                    # (2)
    """Shuffle training data"""
    r = np.random.permutation(len(y))   # get the index numbers for random shuffle (permutation )
    return X[r], y[r]                    # shuffle X and y the same way
                                         # (features X and targets y must be shuffled exactly the way)
def _initialize_weights(self, m):        # (1)
    """Initialize weights to zeros"""
    self.w_ = np.zeros(1 + m)            # set w_ to zero (including w0)
    self.w_initialized = True            # after w's are initialized, set w_initialized = True

def _update_weights(self, xi, target):   # (3)
    """Apply Adaline learning rule to update the weights"""
    output = self.net_input(xi)
    error = (target - output)
    self.w_[1:] += self.eta * xi.dot(error)    # w_[1:]: elements of w starting from inxed 1 to the end
    self.w_[0] += self.eta * error             # w_[0]: the first elements  of w (index 0)
    cost = 0.5 * error**2
    return cost
```

```
462
463    def net_input(self, X):                                    # calculate w.x
464        """Calculate net input"""
465        return np.dot(X, self.w_[1:]) + self.w_[0]
466
467    def activation(self, X):                                   # calculate activation, since it is linear: z= w.x
468        """Compute linear activation"""
469        return self.net_input(X)
470
471    def predict(self, X):
472        """Return class label after unit step"""
473        return np.where(self.activation(X) >= 0.0, 1, -1)    # predict the class labe: y_hat = 1 if z>=0, y_hat = -1 otherwise
474
475 #============================================================  adaline stochastic gradient descent (SGD) algorithm
476
477 #============================================================ call Adaline-SGD learning algorithm with initial eata=0.01
478 ada = AdalineSGD(n_iter=15, eta=0.01, random_state=1)
479 ada.fit(X_std, y)
480
481 print(ada.predict(X_std))
482 #sys.exit('Program stopped here')
483 #============================================================ call Adaline-SGD learning algorithm with initial eata=0.01
```

```
484
485 #================================================================= plot the decision regions and the errors for Adaline-SGD
486 plot_decision_regions(X_std, y, classifier=ada)        # call for plotting
487 plt.title('Adaline - Stochastic Gradient Descent')
488 plt.xlabel('sepal length [standardized]')
489 plt.ylabel('petal length [standardized]')
490 plt.legend(loc='upper left')
491
492 # plt.tight_layout()
493 # plt.savefig('./adaline_4.png', dpi=300)
494 plt.show()
495
496 plt.plot(range(1, len(ada.cost_) + 1), ada.cost_, marker='o')      # error plot
497 plt.xlabel('Epochs')
498 plt.ylabel('Average Cost')
499
500 # plt.tight_layout()
501 # plt.savefig('./adaline_5.png', dpi=300)
502 plt.show()
503 print('================================End of the Adaline-SGD Algorithm')
504 print()
505 print()
506 #================================================================= plot the decision regions and the errors for Adaline-SGD
507
508 #ada = ada.partial_fit(X_std[0, :], y[0])
509
510 sys.exit('Program stopped here')
```