



ESS2222

**Machine Learning in Earth Science Problems
Lecture 1 - Learning from Data**

Hosein Shahnas

University of Toronto, Department of Earth Sciences,

Some Useful References

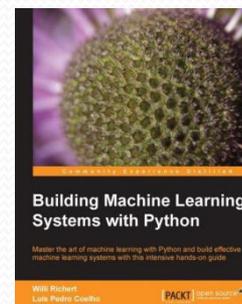
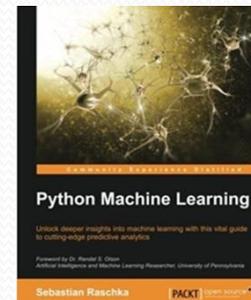
❑ **Learning from Data**
Y. S. Abu-Mostafa

❑ **Python Machine Learning**
S. Raschka, PACKT, Open Source

❑ **Building Machine Learning Systems with Python Mastering Machine**
W. Richert & L. P. Coelho, PACKT, Open Source

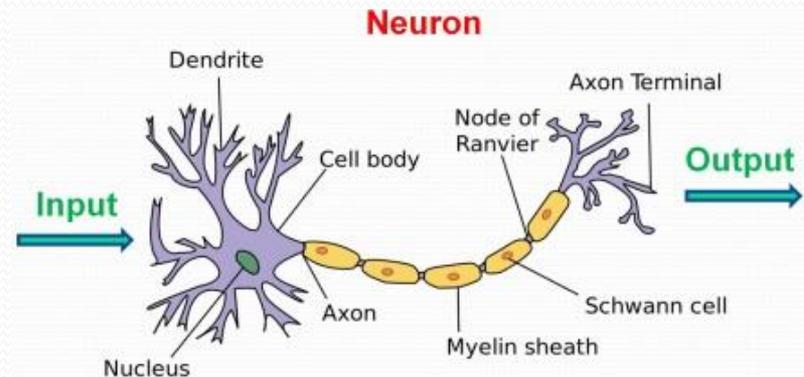
❑ **Learning with scikit-learn**
G. Hackeling, PACKT, Open Source

❑ **Look at the other relevant books published by PACKT**



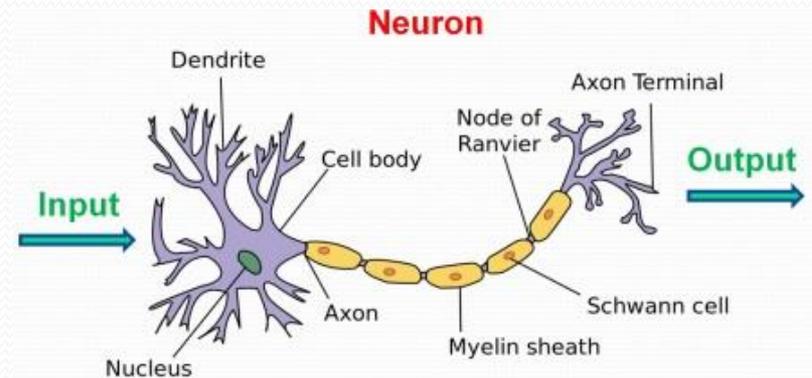
Outline of the Course

- ❑ Machine Learning and the Feasibility of Learning
- ❑ Linear and Nonlinear Models
- ❑ Testing, Training, Error and Noise
- ❑ Generalization Theory
- ❑ The VC Dimension
- ❑ Bias-Variance Trade-off
- ❑ Overfitting and Regularization
- ❑ Support Vector Machine
- ❑ Validation
- ❑ Kernel Methods
- ❑ Random Forest
- ❑ Neural Networks
- ❑ Convolutional Networks
- ❑ Regression Learning
- ❑ Multi Class Learning
- ❑ Machine Learning in Earth Sciences



Machine Learning

- ❑ Learning Components
- ❑ Illustrative Example
- ❑ PLA - A simple model
- ❑ Adaline Algorithm
- ❑ Types of learning

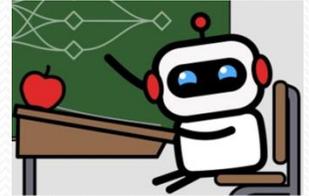


Learning from data

Machine Learning: Learning from data (exploring a target function)

Mathematical Aspects: Provides a conceptual framework

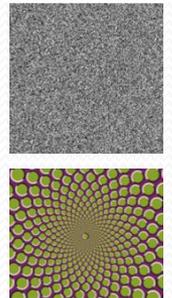
Practical Aspects: How to do it in real work



Components of learning – Criteria to be checked:

- 1- The problem cannot be elaborated mathematically ✓
- 2- There is data ✓
- 3- A pattern exists ✓

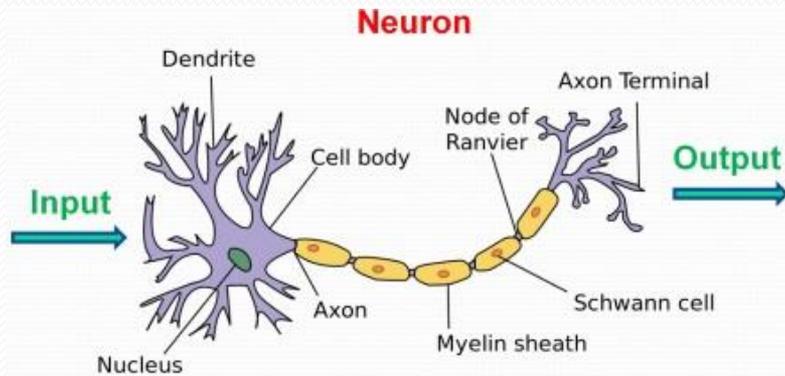
($F = m\gamma$?)



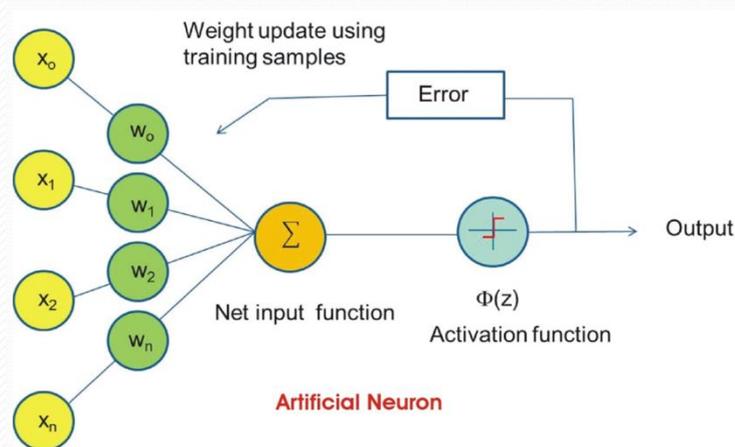
Perceptron Rule

First concept of a simplified brain cell ([McCulloch–Pitts \(MCP\) neuron, 1943](#))

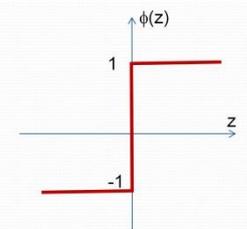
The first concept of the perceptron learning rule based on the MCP neuron model ([Frank Rosenblatt, 1957](#))



$$Z = X \cdot W = \sum_1^n W_i x_i = W_1 X_1 + W_2 X_2 + \dots + W_m X_m \quad \text{Net input}$$



$$\Phi(z) = \begin{cases} 1, & \text{if } z > \theta \\ -1 & \text{otherwise} \end{cases}$$



Credit Approval Example

Input:

Customer application =

[Age, Gender, Marital status, Credit limit, Past payment details, salary,
Current debt, Past debts, Employment status, Years in job,, Years in residence]



Output: good / bad (customer)

- 1- Is there a formula to solve this problem?
- 2- Do we have data?
- 3- Does a pattern exist?

Credit Approval Example

Input:

Customer application =

[Age, Gender, Marital status, Credit limit, Past payment details, salary,
Current debt, Past debts, Employment status, Years in job,, Years in residence]



Output: good / bad (customer)

- 1- Is there a formula to solve this problem? **No** ✓
- 2- Do we have data? **Yes** ✓
- 3- Does a pattern exist? **Yes** ✓

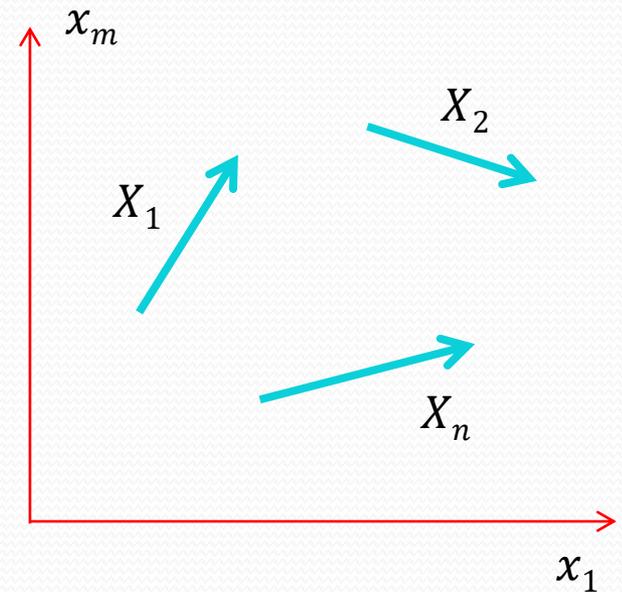
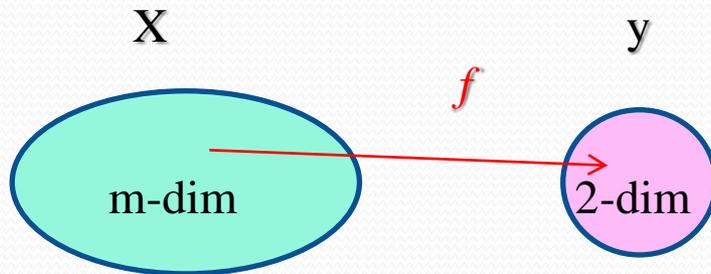
Credit Approval Example



Input: $X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ \dots \\ x_m \end{bmatrix}$

Output: $y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

We are looking for a **target** function: $f: X \rightarrow y$



Feature: An individual measurable property or characteristic of a phenomenon being observed,

Credit Approval Example



n: customers
m: features

Applicant information

	age	gender	salary	Yrs of residence	Yrs in job	Current debt	y	
X_1	X_{11}	X_{21}	X_{31}	X_{m1}	good	y_1
X_2	X_{12}	X_{22}	X_{32}	X_{m2}	bad	y_2
...										...
X_i										y_i
...										...
X_n	X_{1n}	X_{2n}	X_{3n}	X_{mn}	good	y_n

Formalization

Data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Known historical records

\mathbf{X}_i : previous customer's application records

y_i : customer's behaviour

We want: to Learn from these data or equivalently get a **hypothesis (h)**

Hypothesis: A function (h) to approximate the target function

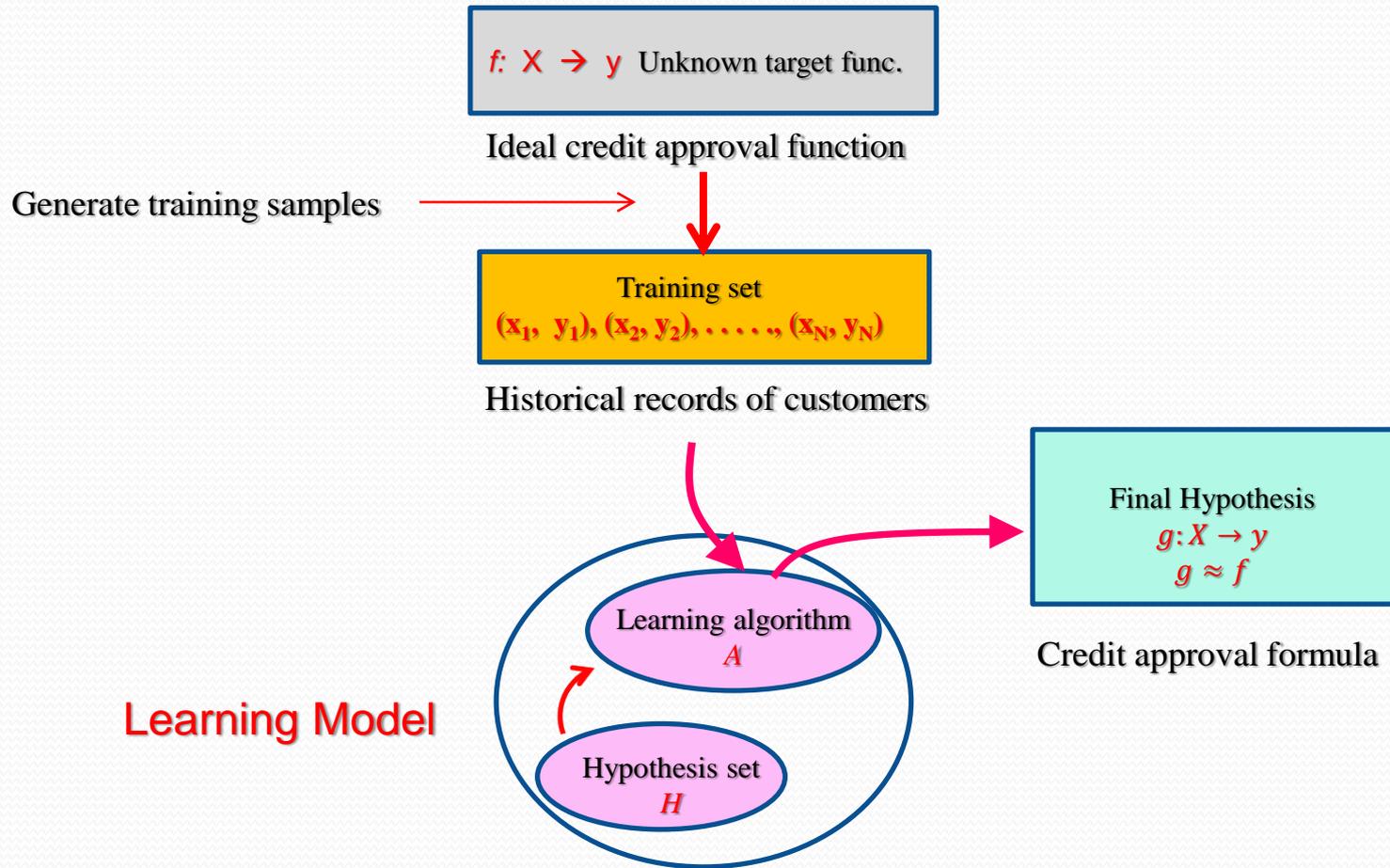
$$\begin{array}{l} g: X \rightarrow y \\ h \longrightarrow g \approx f \end{array}$$

f : unknown

g : known

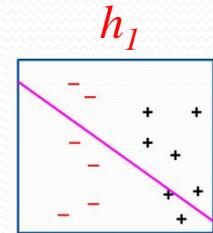
Learning: A process by which we start with h and make it as much as possible close to f ($g \approx f$).

Learning Diagram



Formalization

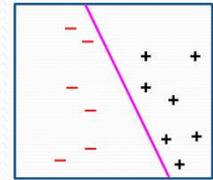
Learning model $\left\{ \begin{array}{l} \text{The hypothesis set} \\ H = \{h\} \quad g \in H \\ \text{The learning algorithm} \\ A \end{array} \right.$



...

...

$g = h_k$



g : Final hypothesis

...

...

e.g., $H = \{\text{percptron, SVM, DNN, CNN, etc.}\}$

h_n

e.g., $A = \{\text{perceptron algorithm, Back propogation, quadratic programming}\}$

Matrix Representation

$$X = \begin{bmatrix} x_{11} & x_{21} & x_{31} & \dots & x_{n1} \\ x_{12} & x_{22} & x_{32} & \dots & x_{n2} \\ & & \square & \dots & \\ x_{1i} & x_{2i} & x_{3i} & \dots & x_{ni} \\ & & \square & \dots & \\ & & \square & \dots & \\ & & \square & \dots & \\ x_{1m} & x_{2m} & x_{3m} & \dots & x_{nm} \end{bmatrix}$$

\updownarrow **m**

\updownarrow **n**

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix}$$

\leftarrow **n** \rightarrow

14

n: Number of samples
m: Number of features

$y_i = -1, +1$
or
 $y_i = 0, 1$
or
 $y_i = F, T$

A Simple Hypothesis: Perceptron

What does the perceptron do?

$$\sum_{i=1}^m w_i x_i > \text{threshold} \quad \text{Approve credit}$$

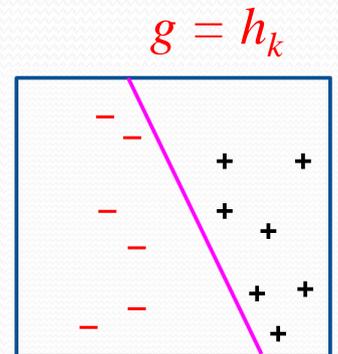
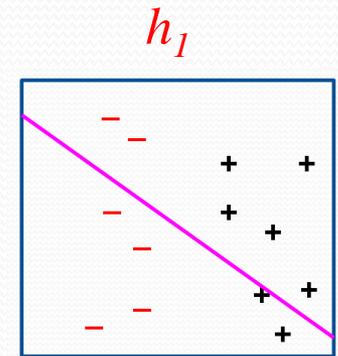
$$\sum_{i=1}^m w_i x_i < \text{threshold} \quad \text{Deny credit}$$

$$h(x) = \text{sign} \left\{ \underbrace{\left[\sum_{i=1}^m w_i x_i \right]}_{\text{Credit score}} - \text{threshold} \right\} \quad \text{set of hypothesis}$$

For convenience: *threshold* \longrightarrow $w_0 x_0$
 where $x_0 = 1$ Artificial feature (coordinate)

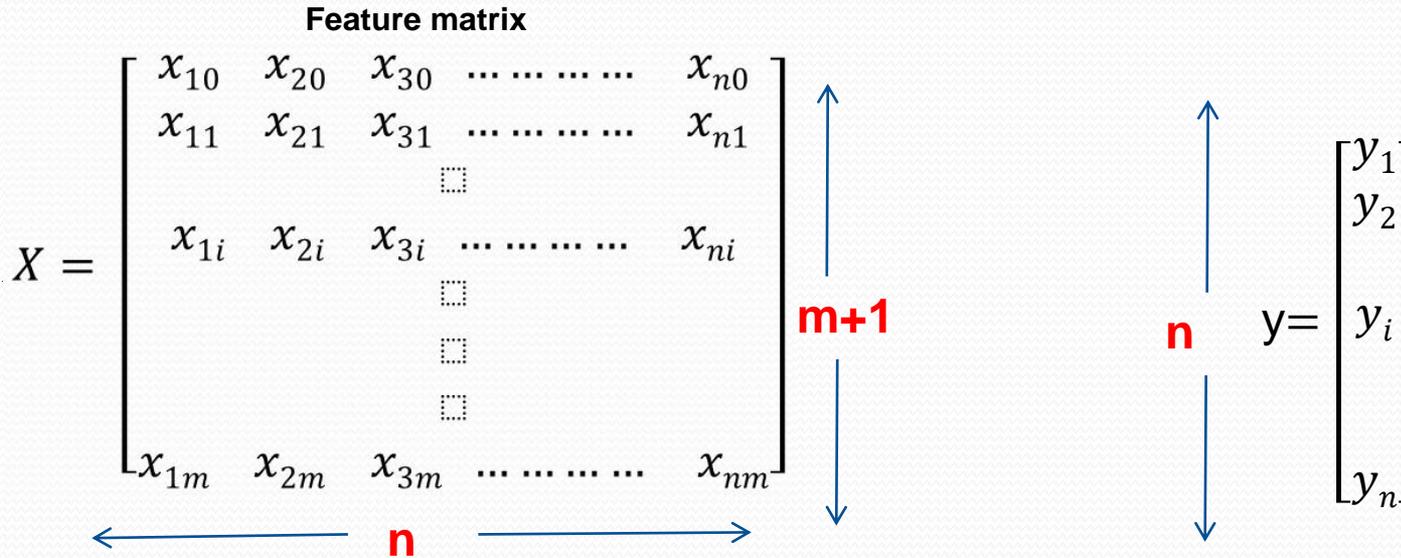
$$h(x) = \text{sign} \left(\sum_{i=0}^m w_i x_i \right)$$

Linearly separable



Final hypothesis

Matrix Representation



n: Number of samples
m+1: Number of features

$y_i = -1, +1$
 or
 $y_i = 0, 1$
 or
 $y_i = F, T$

In vector form

$$\hat{y} = h(x) = \text{sign}(w^T X), \quad y: \text{label}, \hat{y}: \text{calculated label}$$

Matrix Representation

$$W = \begin{bmatrix} W_0 \\ W_1 \\ \square \\ W_i \\ \square \\ \square \\ W_m \end{bmatrix}$$

$$X = \begin{bmatrix} x_{10} & x_{20} & x_{30} & \dots & \dots & \dots & x_{n0} \\ x_{11} & x_{21} & x_{31} & \dots & \dots & \dots & x_{n1} \\ & & & \square & & & \\ x_{1i} & x_{2i} & x_{3i} & \dots & \dots & \dots & x_{ni} \\ & & & \square & & & \\ & & & \square & & & \\ & & & \square & & & \\ x_{1m} & x_{2m} & x_{3m} & \dots & \dots & \dots & x_{nm} \end{bmatrix}$$

$$W^T X = [W_0, W_1, W_2, \dots, W_j, \dots, W_m] \begin{bmatrix} x_{10} & x_{20} & x_{30} & \dots & \dots & \dots & x_{n0} \\ x_{11} & x_{21} & x_{31} & \dots & \dots & \dots & x_{n1} \\ & & & \square & & & \\ x_{1i} & x_{2i} & x_{3i} & \dots & \dots & \dots & x_{ni} \\ & & & \square & & & \\ & & & \square & & & \\ & & & \square & & & \\ x_{1m} & x_{2m} & x_{3m} & \dots & \dots & \dots & x_{nm} \end{bmatrix} = \begin{bmatrix} Z_1 \\ Z_2 \\ \square \\ \square \\ Z_i \\ \square \\ \square \\ Z_n \end{bmatrix}$$

Perceptron Learning Algorithm (PLA)

Set of hypothesis: $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{X}) \in H$

Data: $(\mathbf{X}_1, y_1), (\mathbf{X}_2, y_2), \dots, (\mathbf{X}_m, y_m)$

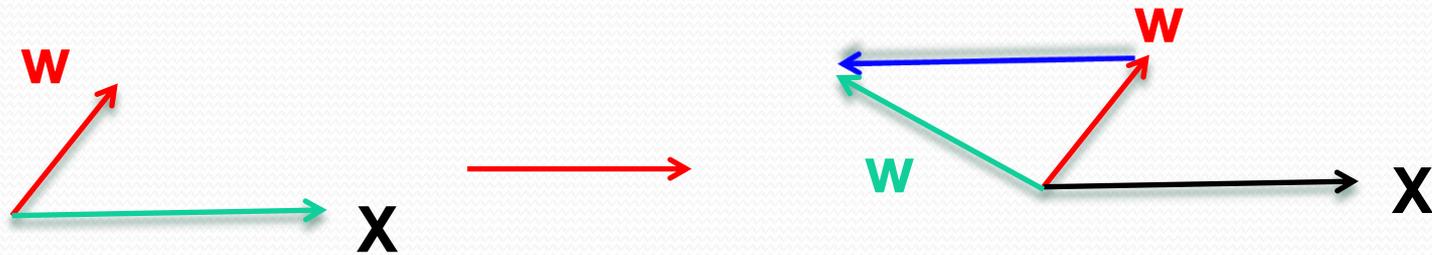
Current customers' records and their behaviour

Algorithm: Takes data and searches for misclassified items,

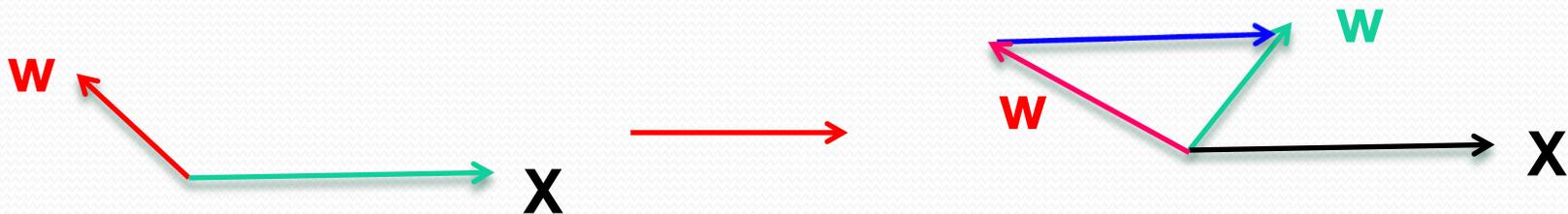
If $\text{sign}(\mathbf{w}^T \mathbf{X}_n) \neq y_n$ then update the weight vector $\mathbf{w} \rightarrow \mathbf{w} + \Delta \mathbf{w}$

Perceptron Learning Algorithm (PLA)

Suppose $\hat{y}_n = \Phi(\mathbf{w}^T \mathbf{X}_n) = +1$, but $y_n = -1$



Or $\hat{y}_n = \Phi(\mathbf{w}^T \mathbf{X}_n) = -1$, but $y_n = +1$



Perceptron Learning Algorithm (PLA)

PLA iteration (Rosenblatt's initial perceptron rule):

I - Initialize the weights w to 0 or small random numbers

II- For each training sample from $(X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)$ and apply PLA on it:

1 - Compute output value \hat{y}^i

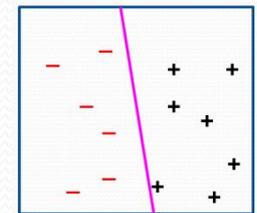
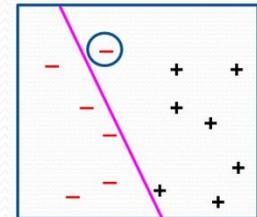
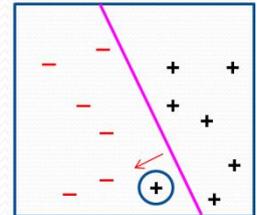
2 - Update the weights:

$$w_j \rightarrow w_j + \Delta w_j \quad (y^i = \pm 1, \hat{y}^i = \pm 1)$$

$$\Delta w_j = \eta (y_i - \hat{y}_i) X_j^i$$

$$0 < \eta \ll 1 \quad \text{Learning rate}$$

It can be proved that if the data are linearly separable, the iteration converges.



Matrix Representation

$$z = w^T X = [w_0, w_1, w_2, \dots, w_j, \dots, w_m]$$

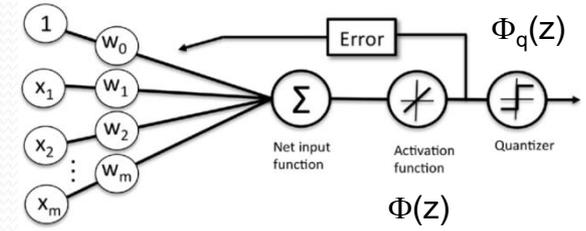
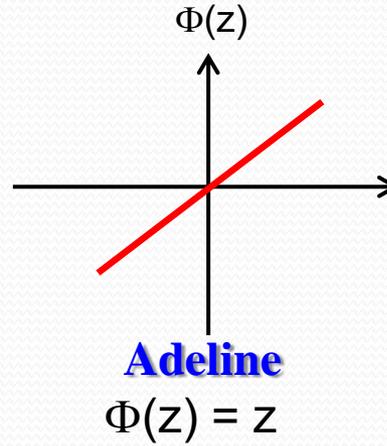
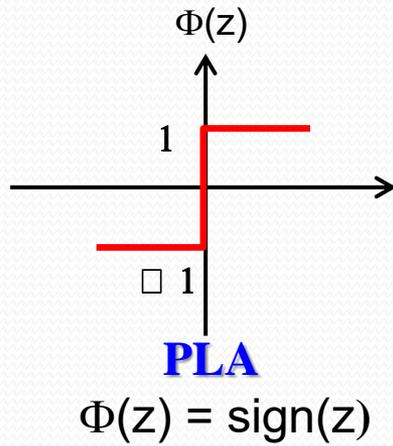
$$\begin{bmatrix} x_{10} & x_{20} & x_{30} & \dots & x_{n0} \\ x_{11} & x_{21} & x_{31} & \dots & x_{n1} \\ & & \square & & \\ x_{1i} & x_{2i} & x_{3i} & \dots & x_{ni} \\ & & \square & & \\ & & \square & & \\ & & \square & & \\ x_{1m} & x_{2m} & x_{3m} & \dots & x_{nm} \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ \square \\ \square \\ z_k \\ \square \\ \square \\ z_n \end{bmatrix}$$

Quantizer:

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

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_i \\ \hat{y}_n \end{bmatrix}$$

ADaptive Linear NEuron (Adaline) Algorithm



In Adaline algorithm (*Widrow-Hoff rule*) the weights are updated based on a linear activation function.

ADaptive LInear NEuron (Adaline) Algorithm

Minimizing cost functions with gradient descent ("batch" gradient descent):

1 - Define a cost function as the Sum of Squared Errors (SSE):

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

2- Use gradient decent optimization algorithm to find weights that minimise the cost function and therefore the **error.**

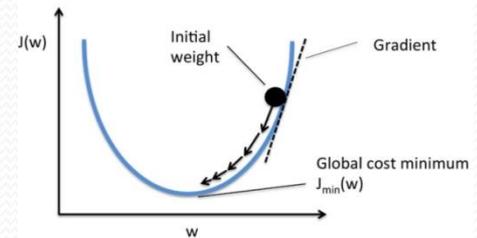
$$w = w + \Delta w$$

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

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

i: Sample index, j: Feature index

Updates: Based on **all samples (batch of samples)**. In Perceptron updates are incrementally after each sample.



Types of Learning

- 1) Supervised learning**
 - 2) Unsupervised learning**
 - 3) Reinforcement learning**
-
- A) Classification Problem**
 - B) Regression Problem**

Types of Learning

1) **Supervised learning:** **If the outputs of the data are explicitly given, we have supervised learning,**

$(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)$

y_i

$X_i \rightarrow y_i$

Previous customer

Credit behaviour

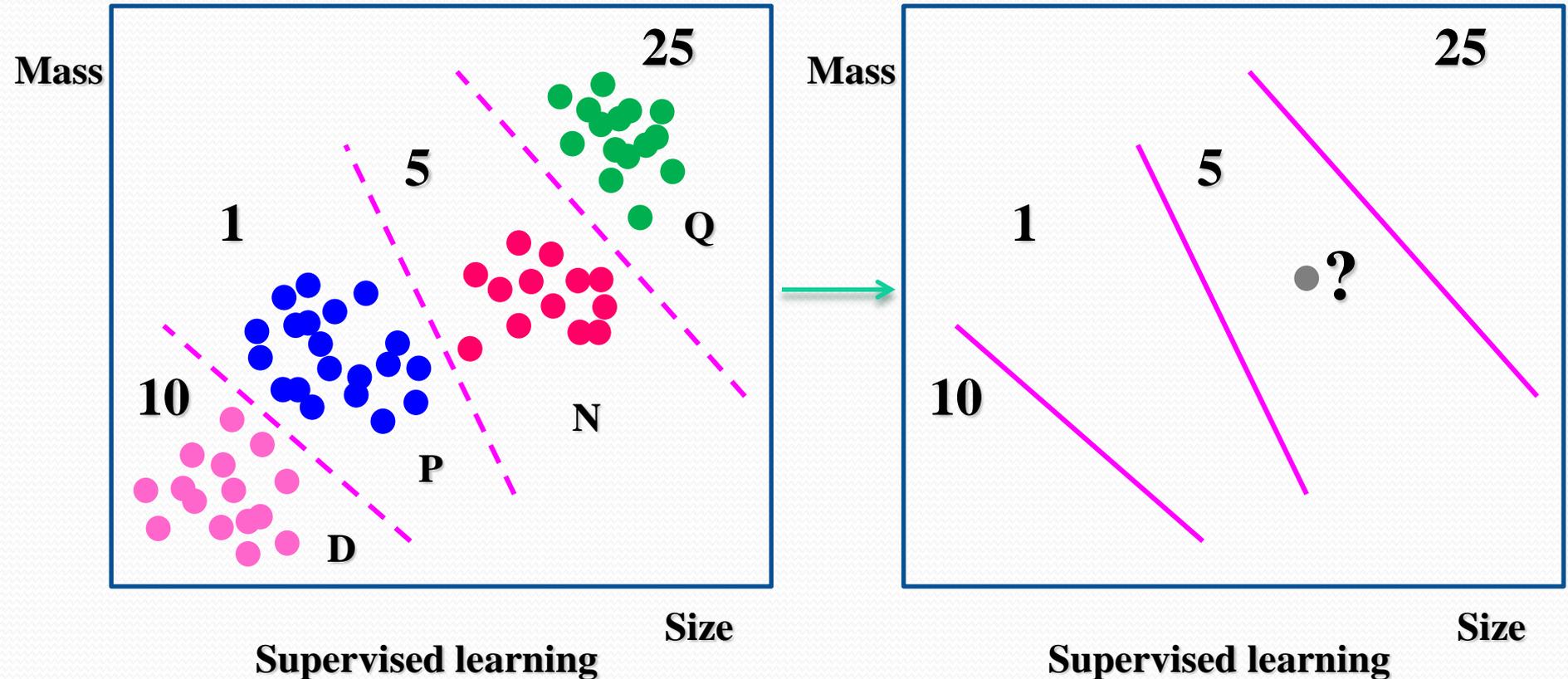
To classify the future customer's credit

Types of Learning

1) Supervised learning: The outputs **are** known. **Ex: Vending machine – Coin recognition**

$$(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)$$

Learned Model

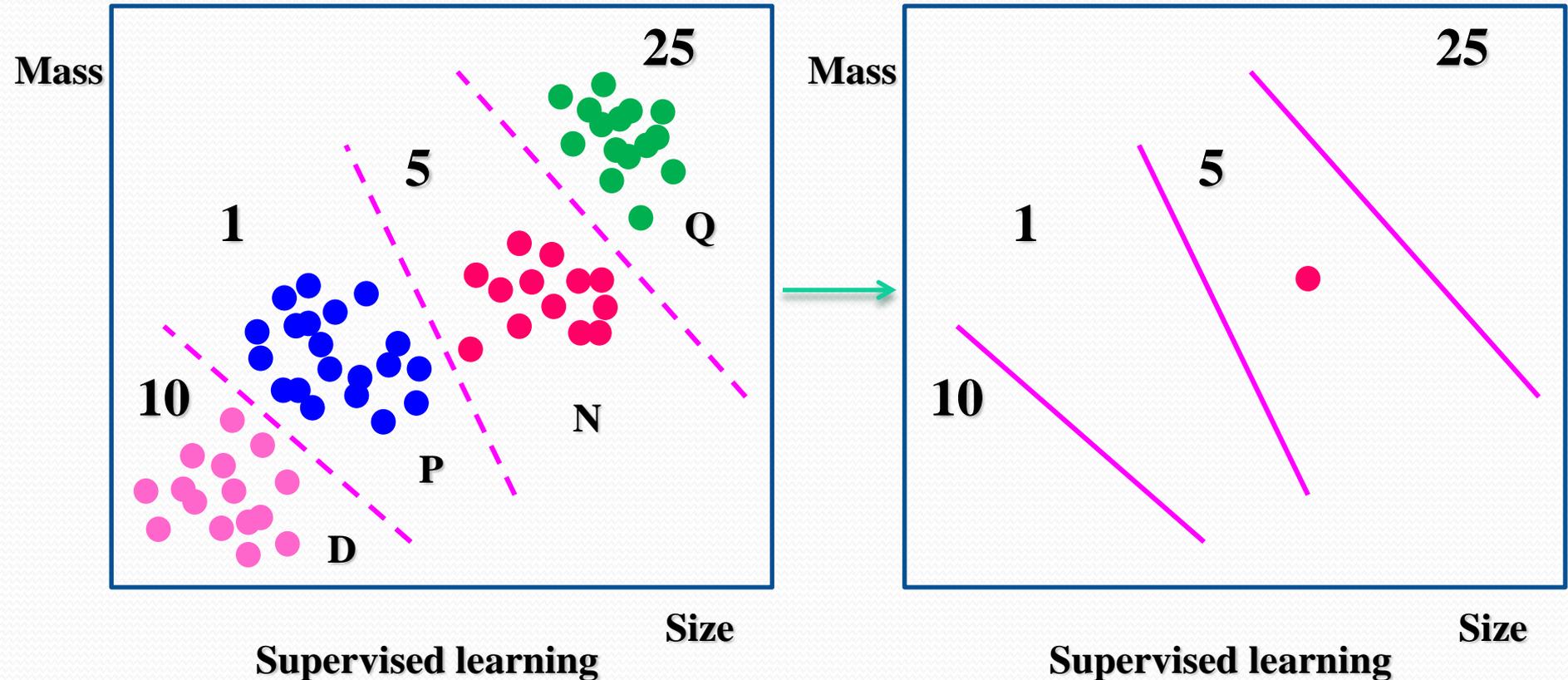


Types of Learning

1) Supervised learning: The outputs **are** known. **Ex: Vending machine – Coin recognition**

$$(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)$$

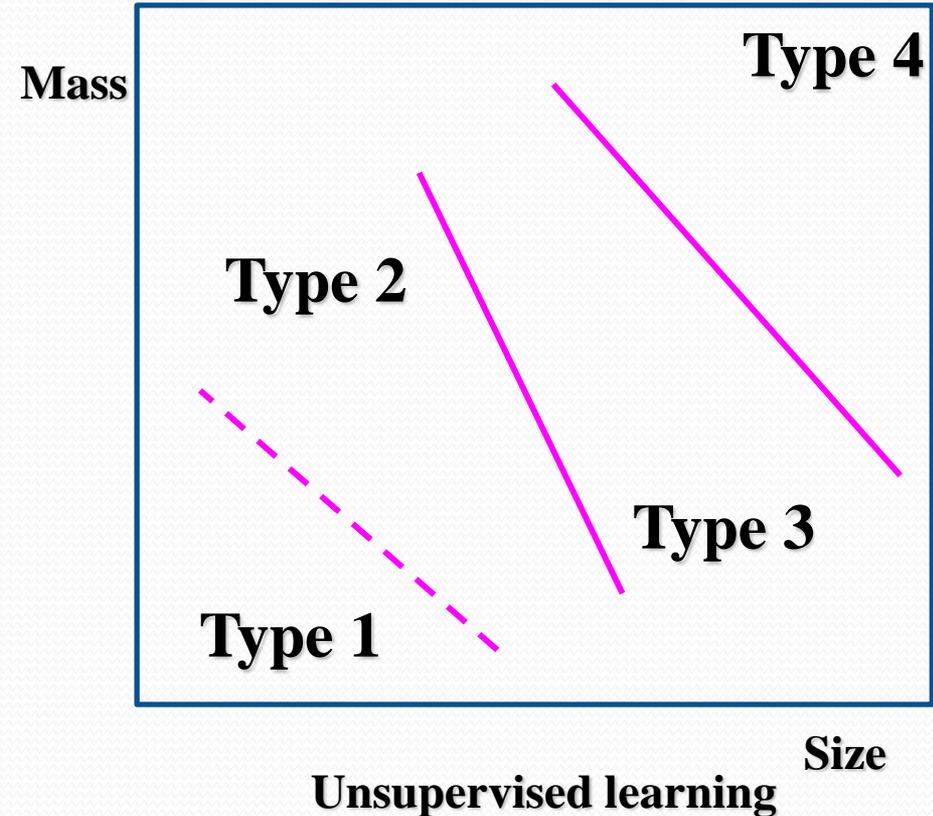
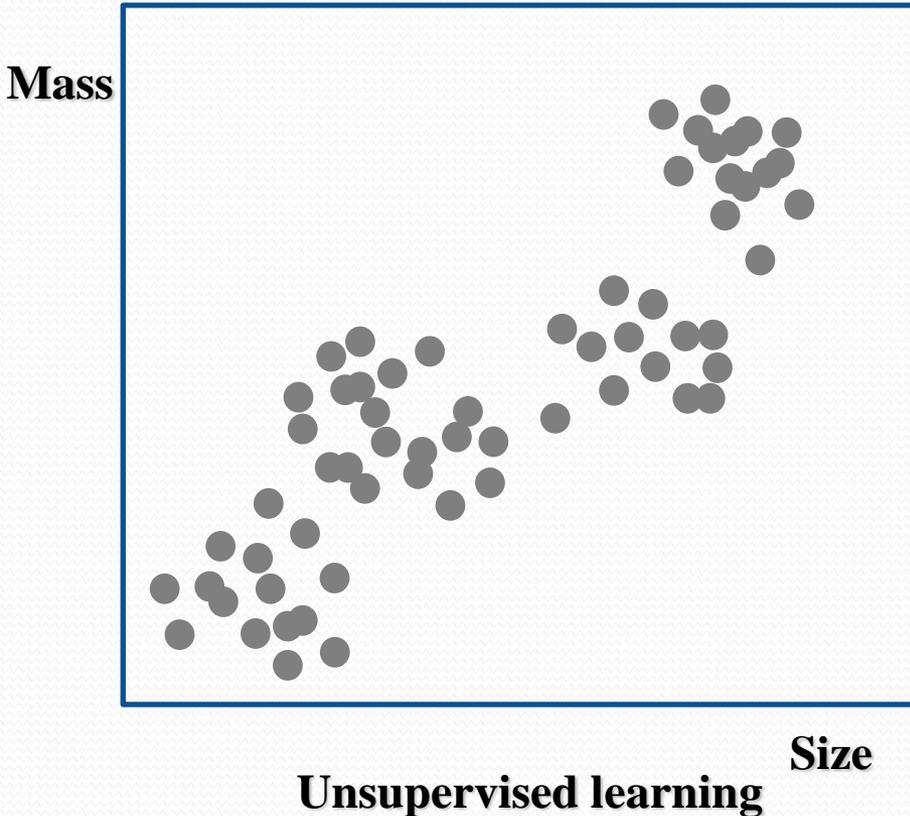
Learned Model



Types of Learning

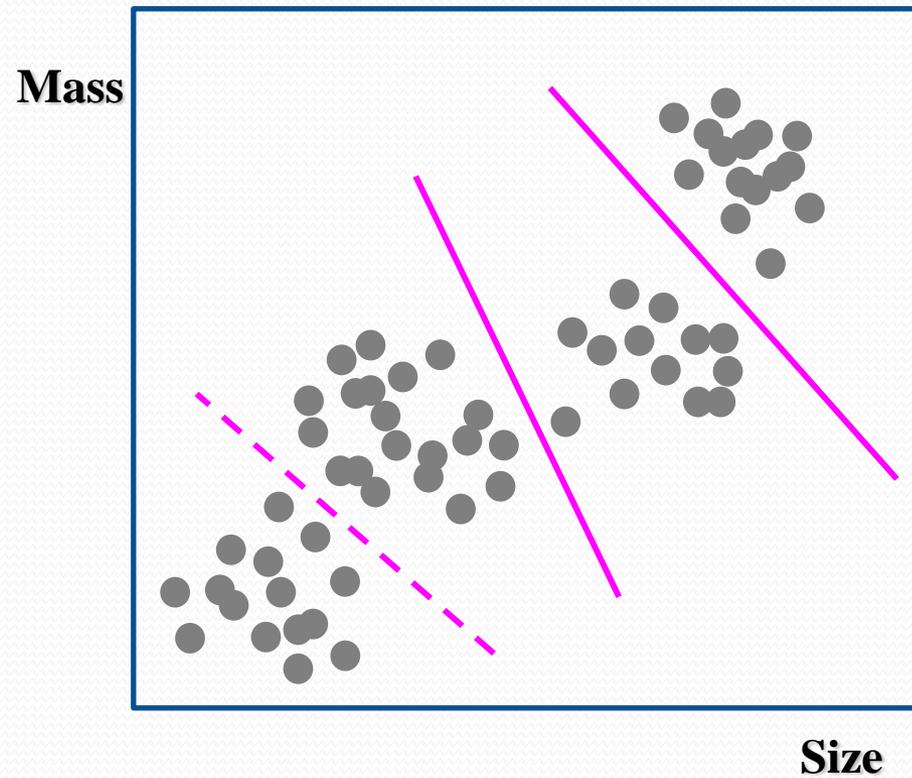
2) Unsupervised learning: The outputs **are not** known. Ex: Vending machine – Coin recognition

$$(X_1, \cancel{y_1}), (X_2, \cancel{y_2}), \dots, (X_n, \cancel{y_n})$$



Types of Learning

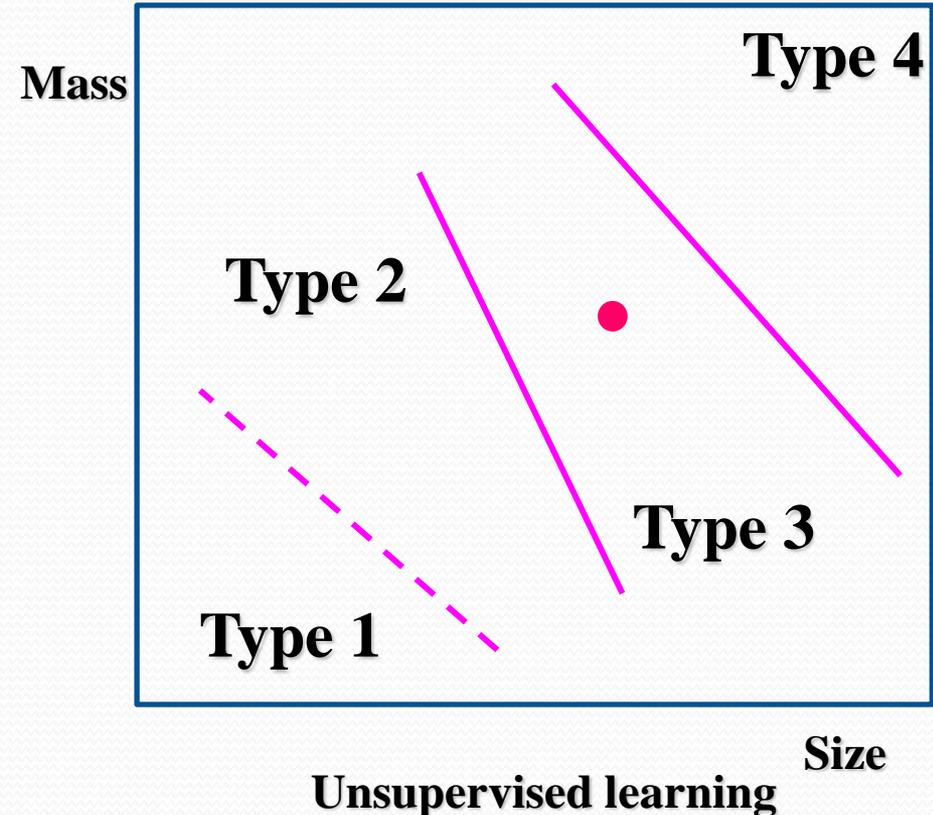
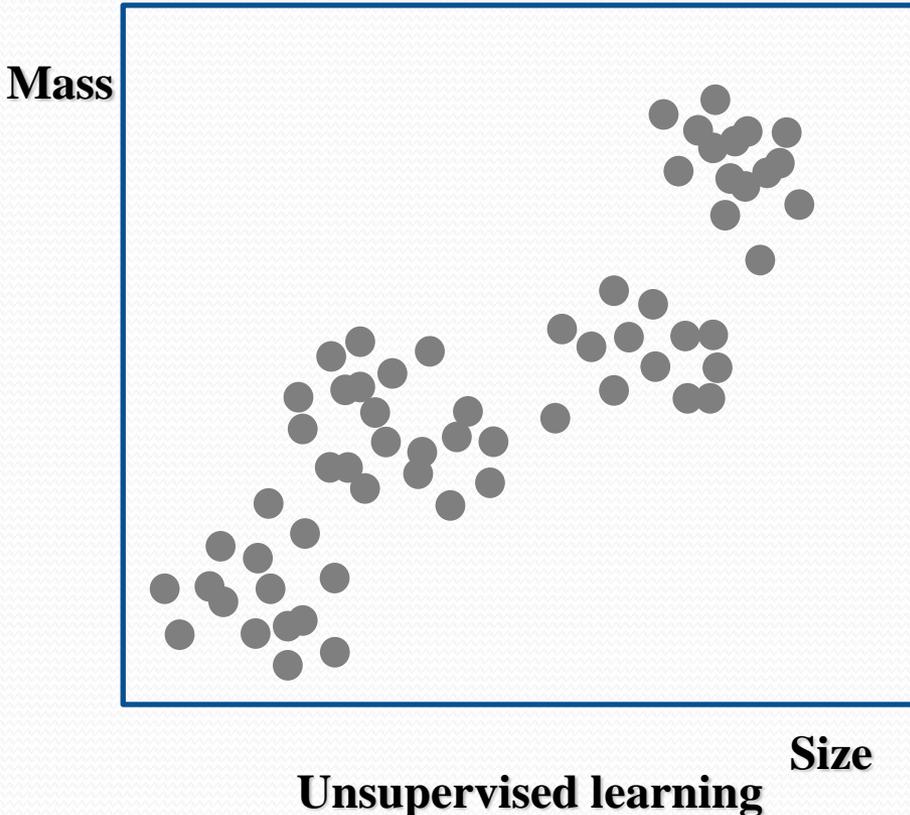
Higher Level Representation of Input Data



Types of Learning

2) Unsupervised learning: The outputs **are not** known. Ex: Vending machine – Coin recognition

$$(X_1, \cancel{y_1}), (X_2, \cancel{y_2}), \dots, (X_n, \cancel{y_n})$$



Types of Learning

3) Reinforcement learning

$$(X_i, y_i)$$

Input + some output (but **not** very clear)

Learning about the outputs by: **Positive** and **Negative** Rewards

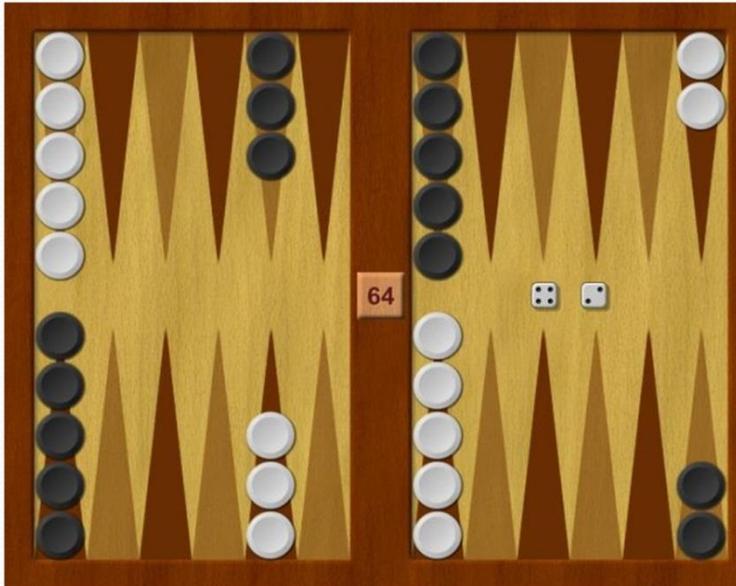


Types of Learning

3) Reinforcement learning

Input + some output (but **not** very clear)

$$(X_i, y_i)$$



Start with a crazy move
Win or lose
Propagate back the credit
Do this hundred billion times

Backgammon

Classification vs Regression

In **classification** machine learning the samples are classified in one of two (binary) or more (multi) classes and the outputs (y) form a discrete spectrum.

Example: Classifying the customers of a bank in **bad** and **good** credit customers.

In **regression** machine learning the outputs (y) are continuous.

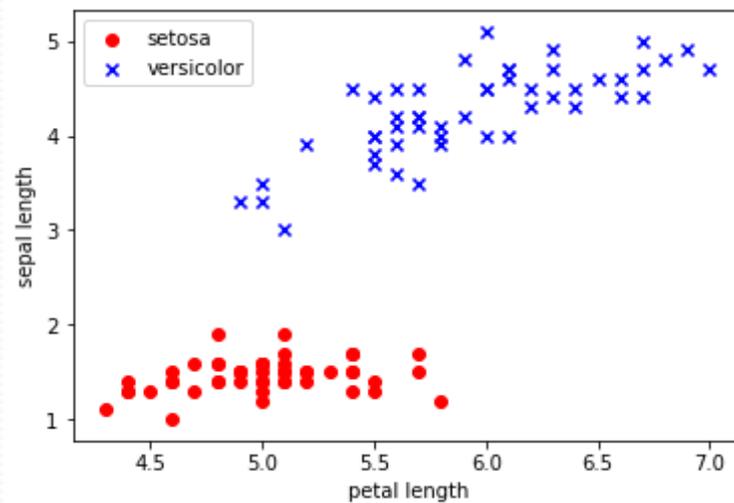
Example: The amount of the credit.

Perceptron Algorithm in Python

```
1'''
2Perceptron Algorithm - Python Machine Learning (Sebastian Raschka)
3Iris Classification problem - Modified version
4'''
5
6import sys # System-specific parameters and functions
7#sys.exit(1)
8'''
9sys – System-specific parameters and functions. This module provides access to some variables
10used or maintained by the interpreter and to functions that interact strongly with the interpreter.
11'''
12#===== import iris data from web source
13import pandas as pd
14'''
15pandas is a Python package providing fast, flexible, and expressive data structures designed to make
16working with “relational” or “labeled” data both easy and intuitive.
17'''
18df = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data', header=None)
19print('df.head() = ', df.head(7)) # print the first 7 lines
20print('df.tail() = ', df.tail(5)) # print the last 5 lines
21#print('df = ', df)
22
23#===== import iris data from web source
24
25#===== construct arrays from raw data
26import matplotlib.pyplot as plt # import lib. MATLAB like graphic
27import numpy as np # import numpy package
28'''
29NumPy is the fundamental package for scientific computing with Python.
30It contains among other things: a powerful N-dimensional array object.
31sophisticated (broadcasting) functions. tools for integrating C/C++ and
32Fortran code.
33'''
34
35y = df.iloc[0:100, 4].values # get data at the fifth (4) column for the first 100 lines
36#print('y = ', y)
37
38y = np.where(y == 'Iris-setosa', -1, 1) # set y = 1 for y = 'Iris-setosa' and y = -1 otherwise
39print('y = ', y)
40
41#X = df.iloc[0:100, [0, 1, 2, 3]].values
42X = df.iloc[0:100, [0, 2]].values # get data from the first (0) and third (2) columns for the first 100 lines
43#print('X = ', X)
44print('X.shape = ', X.shape, 'X.shape[0] = ', X.shape[0], 'X.shape[1] = ', X.shape[1])
45#sys.exit(1)
46#===== construct arrays from raw data
47
```


Perceptron Algorithm in Python

```
47
48 #===== scatter plot for Setosa and Versicolor iris
49
50 plt.scatter(X[:50, 0], X[:50, 1], color='red', marker='o', label='setosa')
51 plt.scatter(X[50:100, 0], X[50:100, 1], color='blue', marker='x', label='versicolor')
52 plt.xlabel('petal length')
53 plt.ylabel('sepal length')
54 plt.legend(loc='upper left')
55 plt.show()
56
57 #===== scatter plot for Setosa and Versicolor iris
```



Perceptron Algorithm in Python

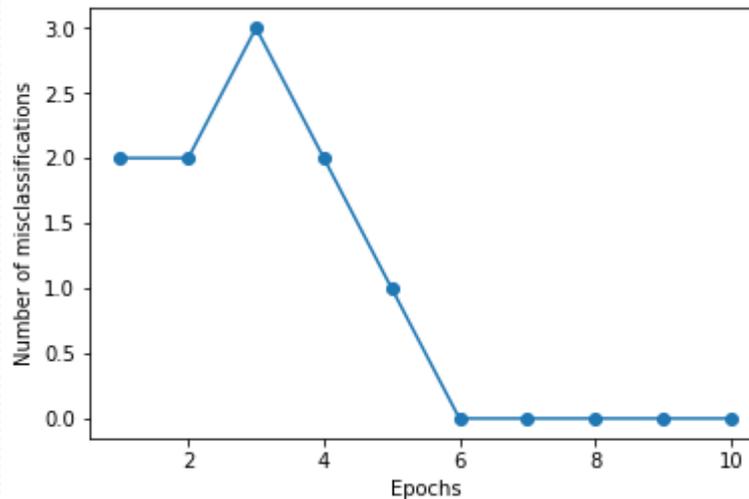
```
59 #----- Perceptron Algorithm
60 class Perceptron(object):
61     """Perceptron classifier.
62
63     Parameters
64     -----
65     eta : float
66         Learning rate (between 0.0 and 1.0)
67     n_iter : int
68         Passes over the training dataset.
69
70     Attributes
71     -----
72     w_ : 1d-array
73         Weights after fitting.
74     errors_ : list
75         Number of misclassifications in every epoch.
76
77     """
78     def __init__(self, eta=0.01, n_iter=10): #Learning parameters
79         self.eta = eta
80         self.n_iter = n_iter
81
```

Perceptron Algorithm in Python

```
82 def fitt(self, X, y):
83     """Fit training data.
84
85     Parameters
86     -----
87     X : {array-like}, shape = [n_samples, n_features]
88         Training vectors, where n_samples
89         is the number of samples and n_features is the number of features.
90     y : array-like, shape = [n_samples]
91         Target values.
92
93     Returns
94     -----
95     self : object
96
97     """
98     self.w_ = np.zeros(1 + X.shape[1]) # set w = 0, having dimension of d = 1 + number of features = 3
99     self.errors_ = [] # initialize error array
100     k = 0
101     for _ in range(self.n_iter):
102         k = k + 1
103         errors = 0
104         j = 0
105         for xi, target in zip(X, y): #iterable over (X,y) for any xi = (x1, x2) and target = y pairs
106             #for xi, target in (0,99):
107                 j = j + 1
108                 update = self.eta * (target - self.predict(xi)) # correction = eta * (y - y_pred)
109                 self.w_[1:] += update * xi # w = w + dw = w + correction * x for each sample
110                 self.w_[0] += update # w = w + dw = w + correction for each sample
111                 errors += int(update != 0.0) # error = error + 1 if correction !=0
112                 #print('k, j = ', k,j, xi, target, ' ', self.predict(xi), 'err = ', errors)
113
114             self.errors_.append(errors)
115         print('Epochs, errors = ', k, ', ', errors)
116     return self
117
118 def net_input(self, X):
119     """Calculate net input"""
120     return np.dot(X, self.w_[1:]) + self.w_[0] # calculate w.x
121
122 def predict(self, X):
123     """Return class label after unit step"""
124     return np.where(self.net_input(X) >= 0.0, 1, -1) # predict (depending on the sign of w.x)
125     '''where:Return elements depending on condition'''
126     '''returns 1 if the condition satisfies, and -1 otherwise'''
127
128 #===== Perceptron Algorithm
```

Perceptron Algorithm in Python

```
129
130#===== Apply Perceptron Algorithm
131ppn = Perceptron() # class object
132#ppn = Perceptron(eta=0.1, n_iter=10)
133ppn.fitt(X, y)
134
135plt.plot(range(1, len(ppn.errors_) + 1), ppn.errors_, marker='o')
136plt.xlabel('Epochs')
137plt.ylabel('Number of misclassifications')
138plt.show()
139#===== Apply Perceptron Algorithm
140#sys.exit('Program stopped here')
141
```



Applications and Packages for Practical Work



Open-source high-level programming language

python™

Python is an interpreted, object-oriented, high-level programming language. Since there is no compilation step, the edit-test-debug cycle is incredibly fast.

<https://www.python.org/>

<https://www.python.org/downloads/>

<https://docs.python.org/3.6/index.html>

(Documentation)

Applications and Packages for Practical Work

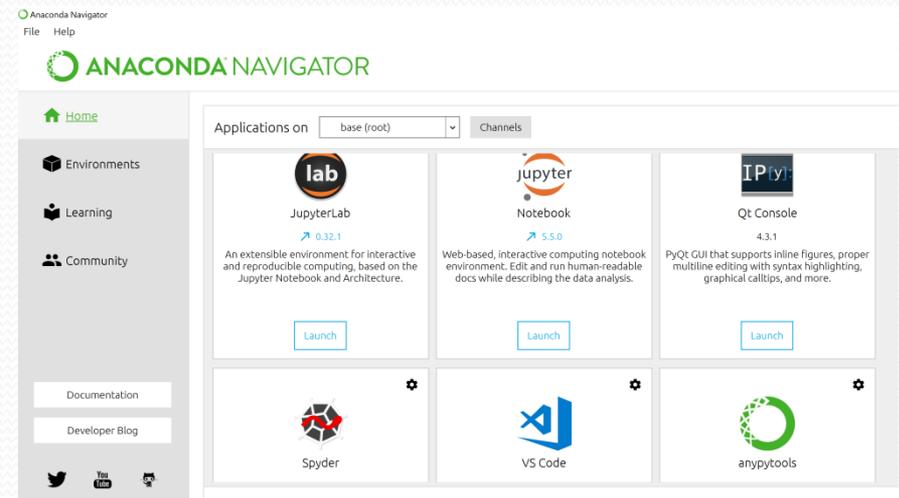


Open-source Python Data Science Platform

The open source Anaconda Distribution is the fastest and easiest way to do Python and R data science and machine learning on Linux, Windows, and Mac OS X. It's the industry standard for developing, testing, and training on a single machine.

<https://anaconda.org/>

<https://www.anaconda.com/download/>



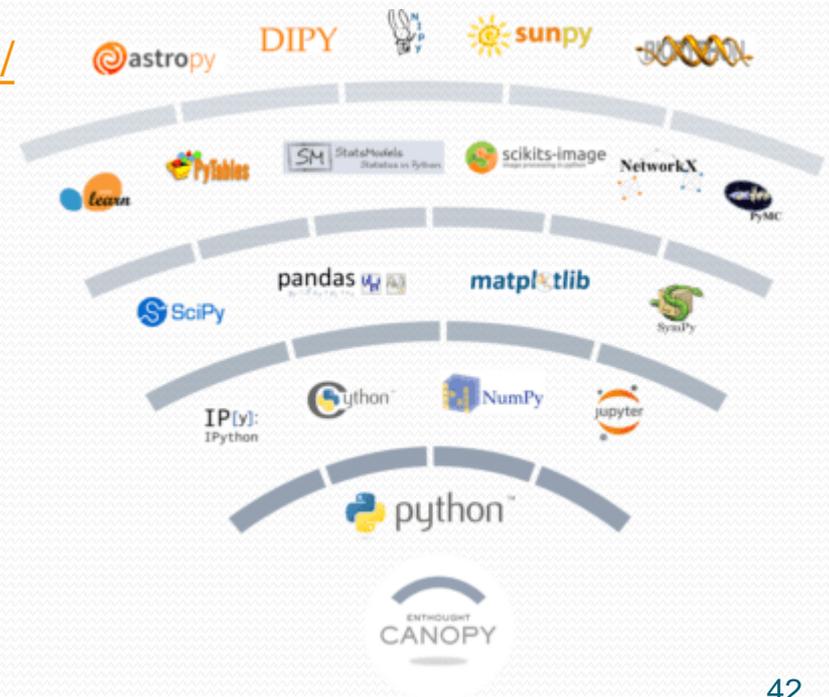
Applications and Packages for Practical Work



Open-source software library

Canopy is a tailor-made for the workflows of scientists and engineers, combining a streamlined integrated analysis environment over 450 proven scientific and analytic Python packages from the trusted Enthought Python Distribution. Canopy provides a complete, self-contained installer that gets you up and running with Python and a library of scientific and analytic tools – fast.

<https://www.enthought.com/product/canopy/>



Applications and Packages for Practical Work



Open-source software library

TensorFlow™ is an open source software library for high performance numerical computation. Its flexible architecture allows easy deployment of computation across a variety of platforms (CPUs, GPUs, TPUs), and from desktops to clusters of servers to mobile and edge devices. Originally developed by researchers and engineers from the Google Brain team within Google's AI organization, it comes with strong support for machine learning and deep learning and the flexible numerical computation core is used across many other scientific domains.

<https://www.tensorflow.org/>

<https://www.tensorflow.org/install/>

```
conda create --name tensorflow python=3.5
```

```
activate tensorflow
```

```
conda install jupyter
```

```
conda install scipy
```

```
pip install tensorflow-gpu
```

Applications and Packages for Practical Work

K Keras Open-source software library

The open source Anaconda Distribution is the fastest and easiest way to do Python and R data science and machine learning on Linux, Windows, and Mac OS X. It's the industry standard for developing, testing, and training on a single machine.

<https://anaconda.org/conda-forge/keras>

```
conda install -c conda-forge keras
```

Applications and Packages for Practical Work



Open-source web application

The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. Uses include: data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning, and much more.

<https://jupyter.org/>

<http://jupyter.org/install>