



機器學習基礎與演算法

陳弘軒 Hung-Hsuan Chen

<https://www.ncu.edu.tw/~hhchen>

中央大學資訊工程學系

「版權聲明頁」

本投影片已經獲得作者授權台灣人工智慧學校得以使用於教學用途，如需取得重製權以及公開傳輸權需要透過台灣人工智慧學校取得著作人同意；如果需要修改本投影片著作，則需要取得改作權；另外，如果有需要以光碟或紙本等實體的方式傳播，則需要取得人工智慧學校散佈權。

－台灣人工智慧學校

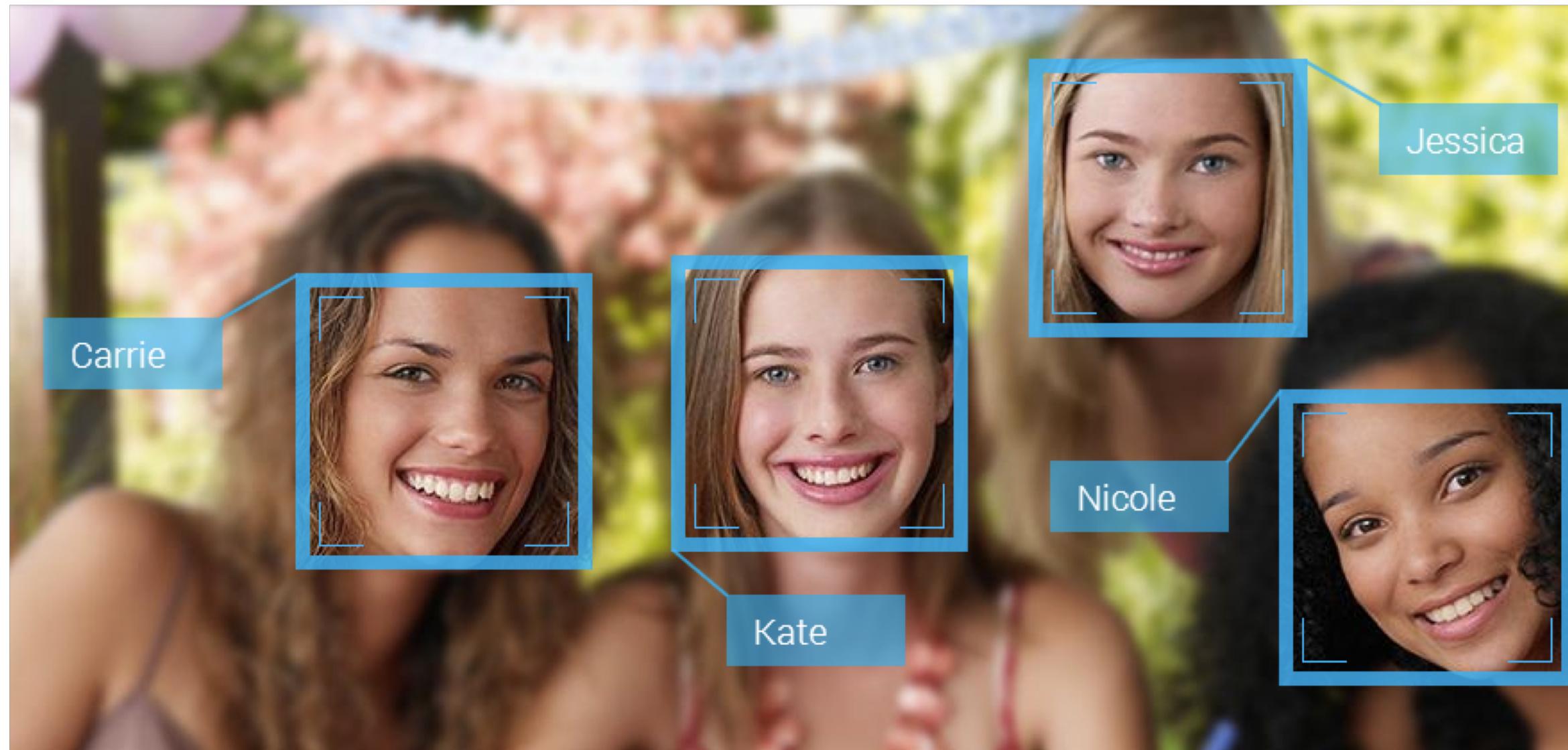
Outline

- Introduction to Machine Learning
- Regression
- Classification
- Support Vector Machine
- Decision trees and ensemble learning
- Dimension reduction

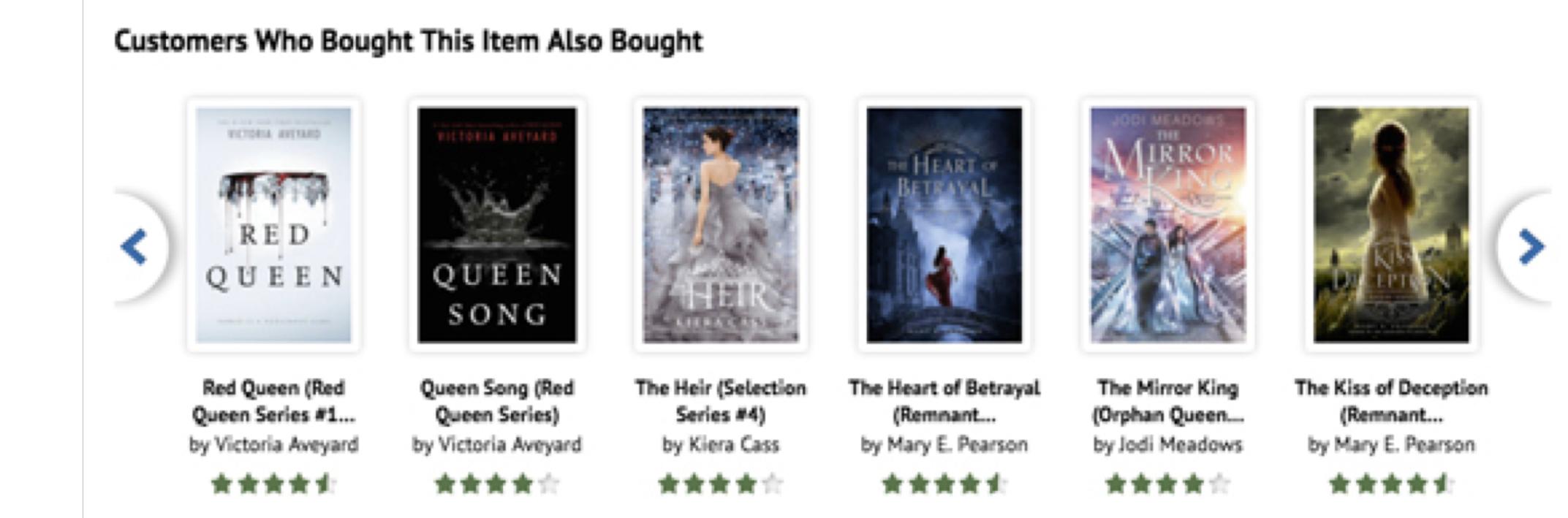


Introduction to Machine Learning

Examples of machine learning today



$7 \rightarrow 7$	$5 \rightarrow 5$
$8 \rightarrow 8$	$3 \rightarrow 3$
$2 \rightarrow 2$	$4 \rightarrow 4$



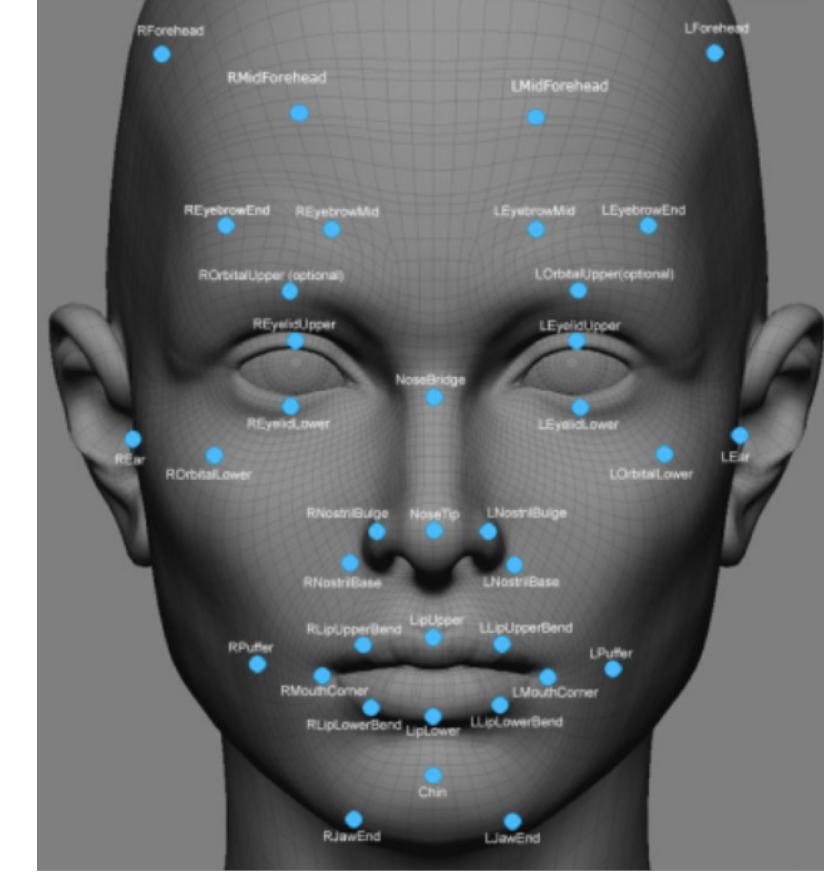
Traditional algorithm vs machine learning

- Traditional algorithm: requires “prior” knowledge of the target problem
 - E.g., sorting
- Some tasks are difficult to solve by prior knowledge
 - How to detect a spam mail? Answers may vary from person to person
- Machine learning: learning from examples
 - Data driven



How to determine the photos with faces?

- Traditional algorithm
 - Round shape, with two black circles (eyes), ...
 - Solve a problem based on your prior knowledge
- Data driven algorithm
 - Show many face/non-face photos to the machine, and let the machine identifies their differences
 - Let computer discover the patterns from the collected data

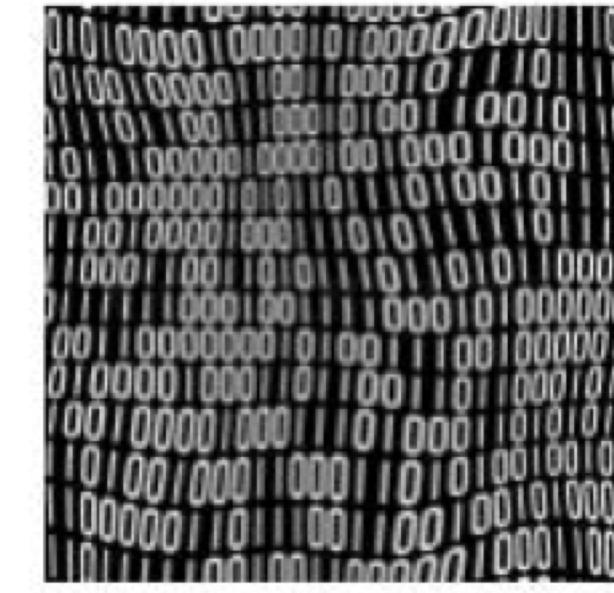


A naïve machine learning approach find faces

- From training data, find the neighboring photos of the new photo
- If most neighboring photos containing faces, the new photo probably contains faces
- If most neighboring photos have no faces, the new photo probably contains no faces



What people see



What computers see

$$\begin{bmatrix} 5 & 8 & 4 \\ 6 & 3 & 5 \\ 2 & 4 & 1 \end{bmatrix} - \begin{bmatrix} 3 & 4 & 2 \\ 2 & 0 & 3 \\ 1 & 3 & 0 \end{bmatrix}$$

Distance between two images



Types of Machine Learning

- Based on the input-output structure, ML can be categorized as (this is not a complete list):
 - Supervised Learning
 - Unsupervised Learning
 - Semi-supervised Learning
 - Reinforcement Learning
- We will mostly discuss supervised learning



Supervised Learning

- Given: a set of <input, output> pairs as the learning samples
- Goal: given an unseen input, predict the corresponding output
- Examples
 1. Input: X-ray photo of chests, output: whether it is cancerous
 2. Input: a sentence, output: whether a sentence is grammatical
 3. Input: some indicators of a company, output: whether it will make profit next year



Output types

- Categorical: ***classification problem***
- *Ordinal outputs: small, medium, large*
- *Non-ordinal outputs: blue, green, orange*
- Real values: ***regression problem***
- Other types are possible (e.g., output a list), but the above two are most widely studied

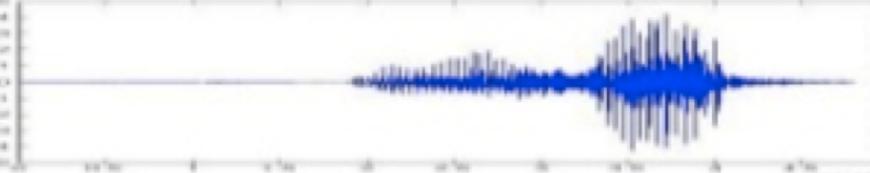


Different types of outputs

Source: Hung-Yi Lee:

https://www.slideshare.net/tw_dsconf/ss-62245351

- Speech Recognition

$f($  $) = \text{“How are you”}$

- Image Recognition

$f($  $) = \text{“Cat”}$

- Playing Go

$f($  $) = \text{“5-5”}$ (next move)

- Dialogue System

$f($ “Hi”
(what the user said) $) = \text{“Hello”}$
(system response)



Terminology

- Training data: a set of data used to discover potentially predictive relationships
- Test data: the data that has been specifically identified for use in tests
- Features (a.k.a. attributes, independent variables)
 - “Input” of a prediction task, usually denoted by x
- Target variable (a.k.a. outputs, dependent variables)
 - “Output” of a prediction task, usually denoted by y
 - In classification, target variables are also called classes



Example

<u>Weight</u>	<u>Wingspan</u>	<u>Webbed feet?</u>	<u>Back color</u>	<u>Species</u>
1000.1	125.0	No	Brown	Buteo jamaicensis
3000.7	200.0	No	Gray	Sagittarius serpentarius

Features

Target variable



Classification (1/3)

- It is a supervised learning task
- Goal: given a feature vector x , predicts which class in C may be associated with x .
- If $|C| = 2 \rightarrow$ Binary Classification
If $|C| > 2 \rightarrow$ Multi-class Classification



Classification (2/3)

- Training and predicting of a binary classification problem:

Training set (binary classification)

Feature Vector ($x_i \in R^d$)	Class
x_1	+1
x_2	-1
...	...
x_{n-1}	-1
x_n	+1

(1) Training

A new instance

Feature Vector ($x_{\text{new}} \in R^d$)	Class
x_{new}	?

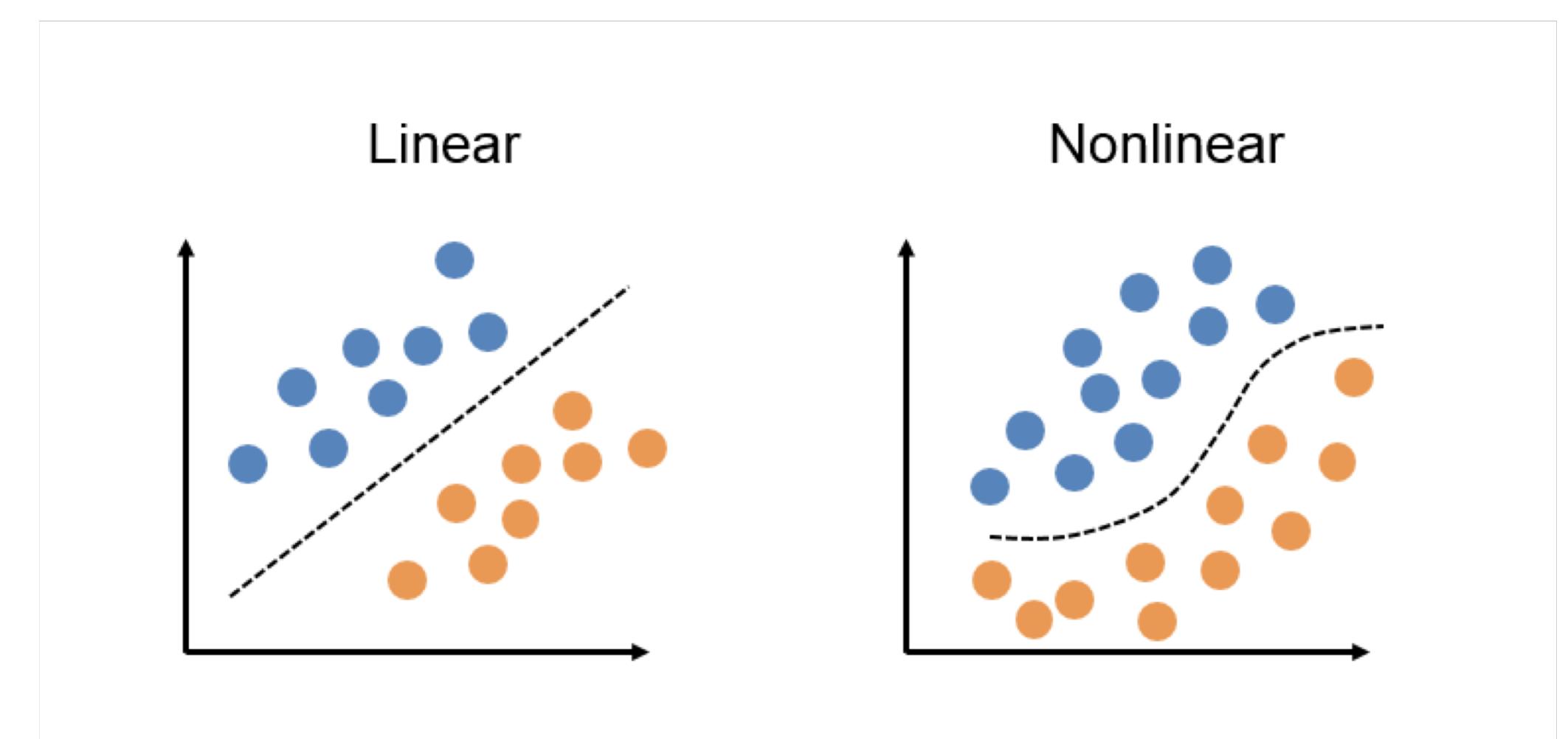
(2) Predicting

Classifier $f(x)$



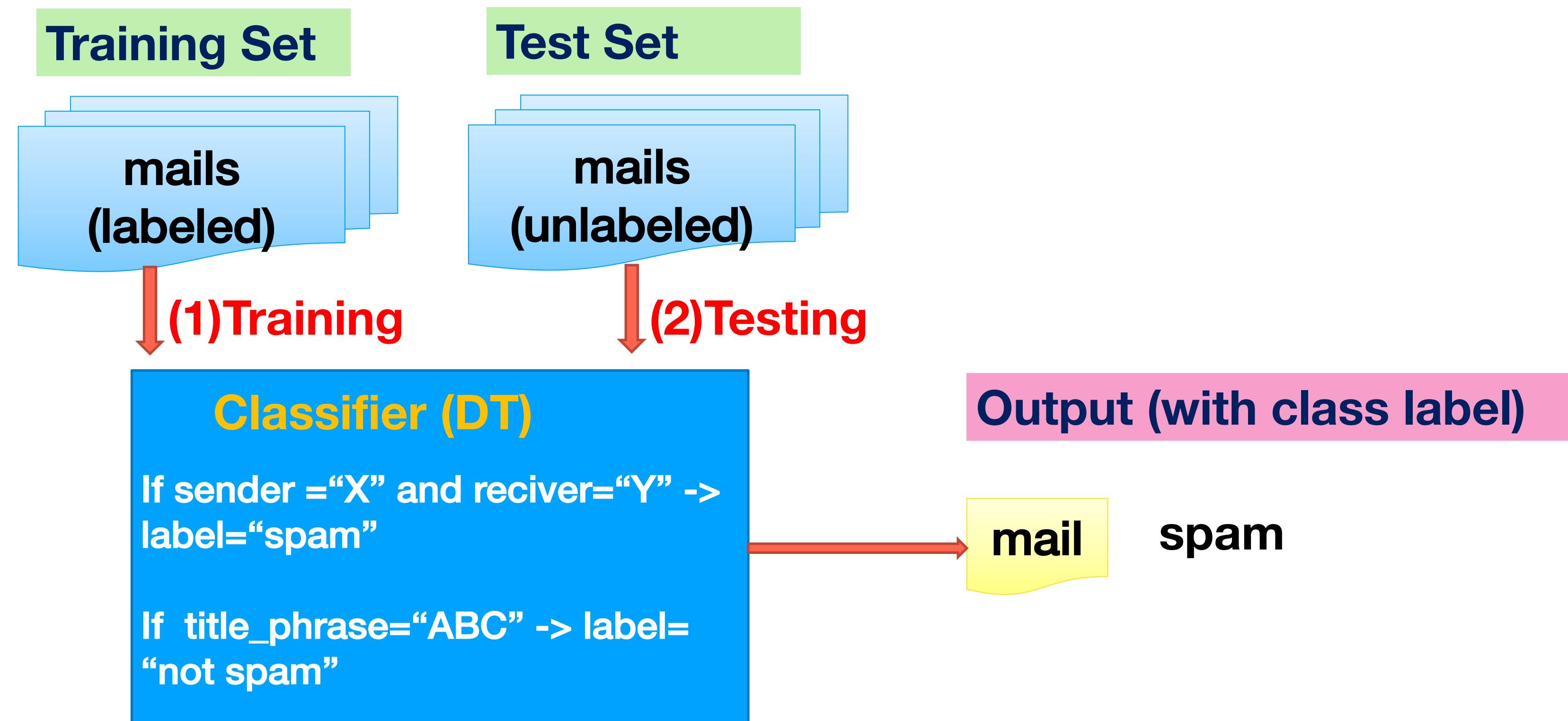
Classification (3/3)

- A classifier can be either **linear** or **non-linear**
- The geometric view of a linear and non-linear classifier
- Famous classification models:
 - k-nearest neighbor (kNN)
 - Decision Tree (DT)
 - Logistic regression
 - Support Vector Machine (SVM)



Real example: E-mail spam check

- Blocking the junk email and passing the normal email



Regression (1/2)

- A supervised learning task that, given a feature vector x , predicts the target value $y \in \mathbb{R}$.
- Training and predicting of a regression problem:

Training set (regression)

Feature Vector ($x_i \in R^d$)	y_i
x_1	+0.26
x_2	-3.94
...	...
x_{n-1}	-1.78
x_n	+5.31

A new instance

Feature Vector ($x_{\text{new}} \in R^d$)	y_i
x_{new}	?

(1) Training

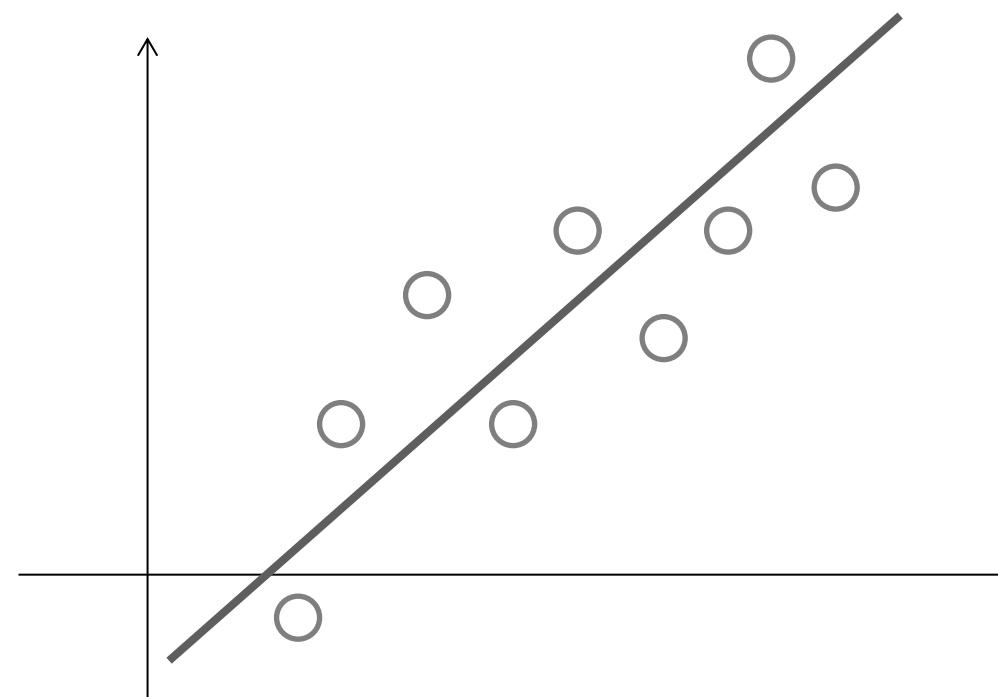
Regressor $f(x)$

(2) Predicting



Regression (2/2)

- The geometric view of a linear regression function

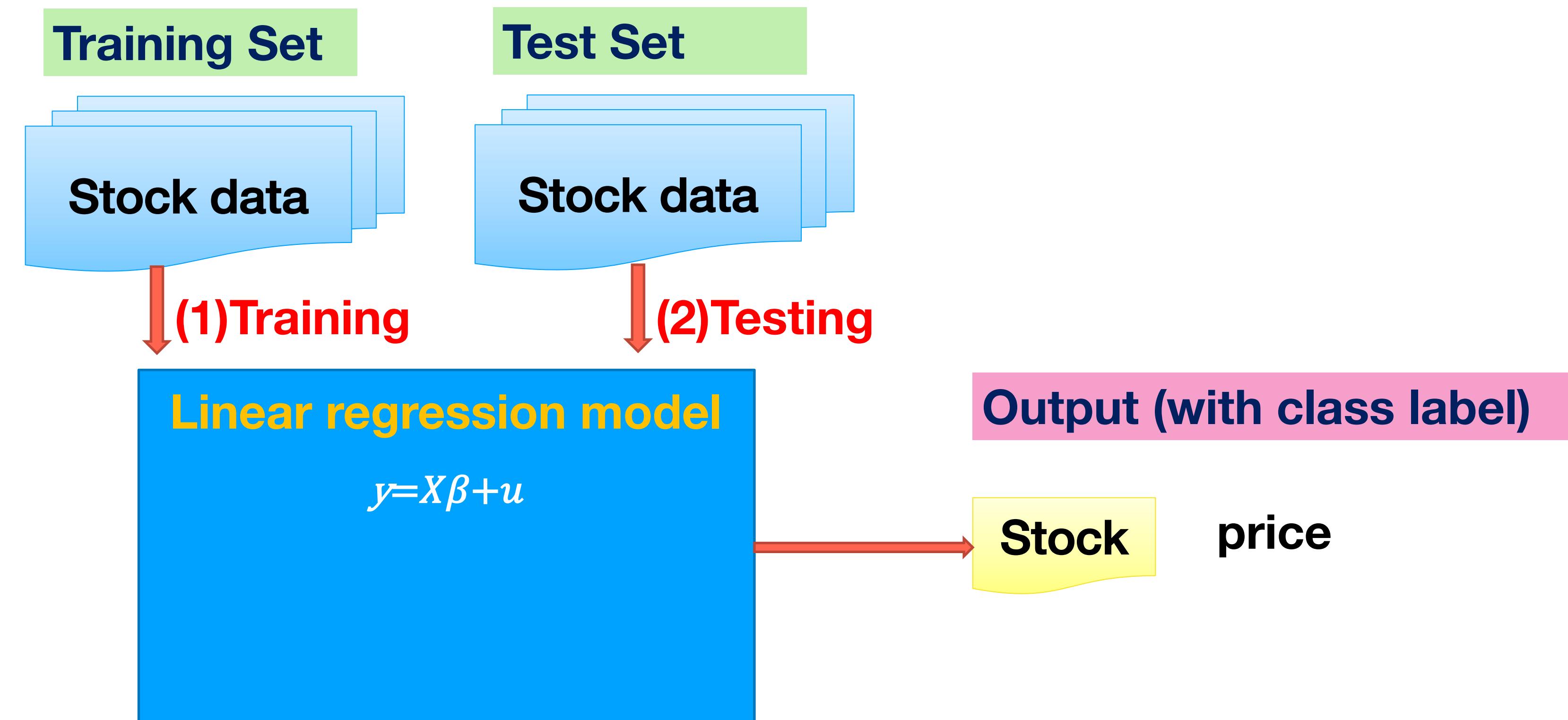


- Some types of regression: linear regression, support vector regression, ...

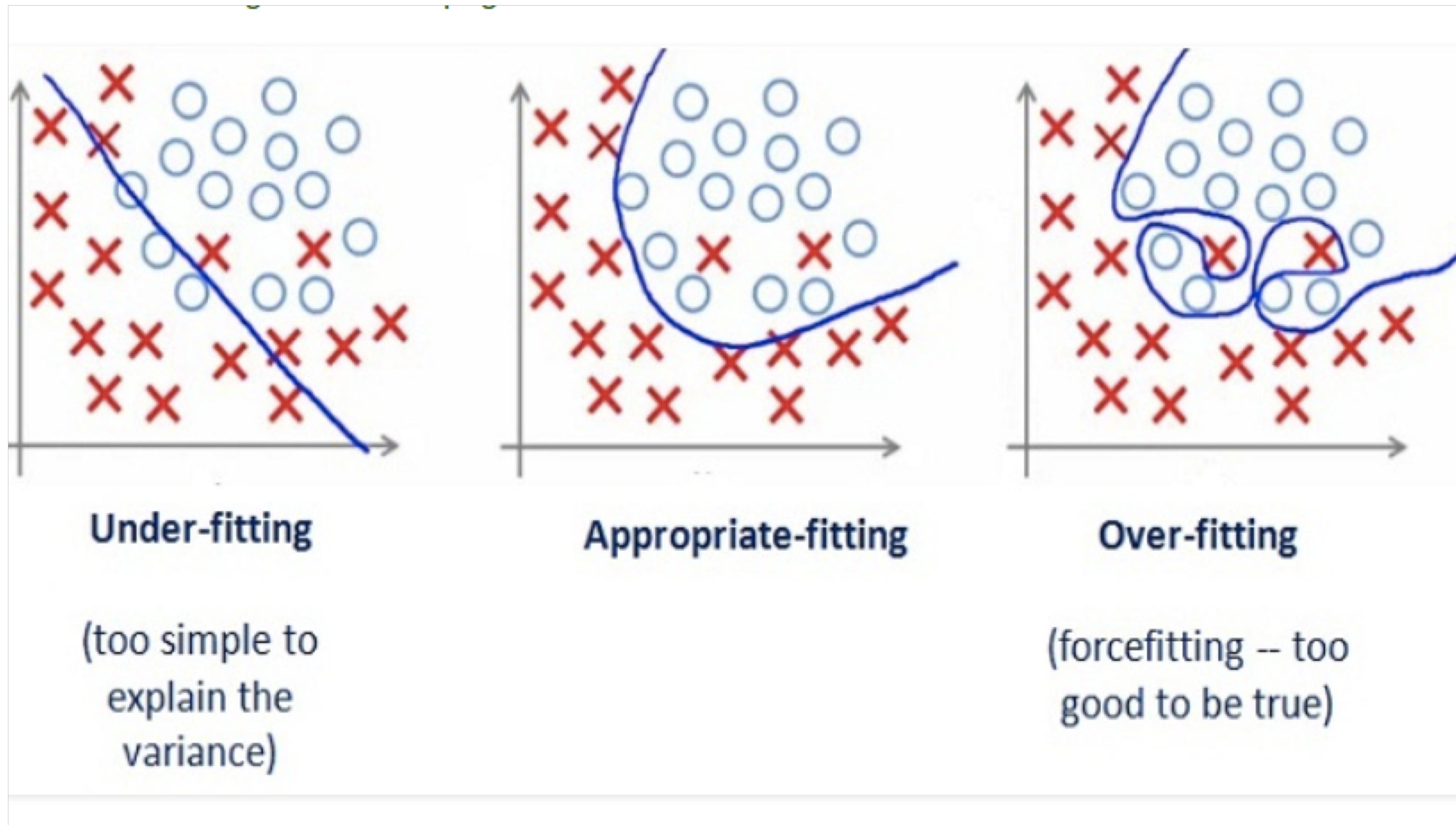


Real Example: Stock price prediction

- Predicting the price of a stock



Overfitting vs underfitting



Quiz

- Explain the difference between a supervised and an unsupervised algorithm
- Explain the difference between classification and regression
- Classify the normal mails and the junk mail based on the labeled datasets
 - Supervised or unsupervised?
 - Predicting the price of stock
 - Classification or regression?

