

FIRST STEP TO PRACTICAL MACHINE LEARNING

KNOWLEDGE SHARING FOR CPE/SKE STUDENTS

SIRAKORN LAMYAI

STUDENT, KASETSART U.

OCTOBER 30, 2018

BEFORE WE START...

Make sure these are installed on your computer.

This page is a guide for installing on Windows

- Python 3.6: Download and install at <https://www.python.org>
- NumPy, Scipy, Matplotlib, Scikit-learn, MLxtend:
Run `pip install numpy scipy matplotlib sklearn mlxtend`

1 Introduction to Machine Learning

- What is Machine Learning?

 - Traditional programming approach

 - Machine learning approach

2 Machine Learning Problems

- Supervised learning

- Unsupervised learning

- Reinforcement learning

3 Model

4 Machine Learning Process

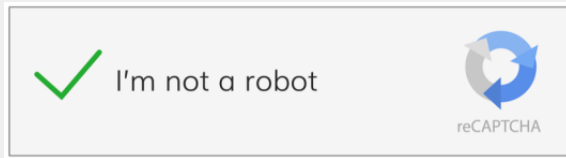
5 Problems for Machine Learning

- Handwriting recognition

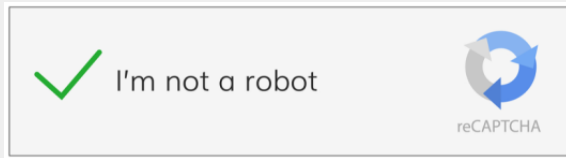
INTRODUCTION TO MACHINE LEARNING

WHAT IS MACHINE LEARNING?

WHAT IS MACHINE LEARNING?

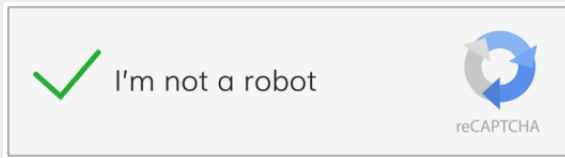


WHAT IS MACHINE LEARNING?



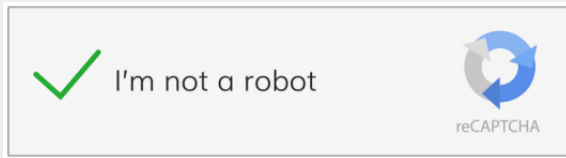
- This is Recaptcha.

WHAT IS MACHINE LEARNING?



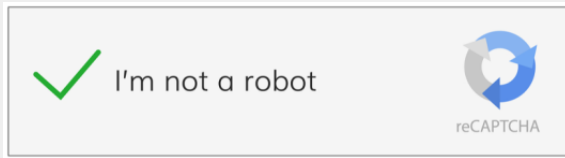
- This is Recaptcha.
 - ▶ Recaptcha helps stop millions of spam a day.

WHAT IS MACHINE LEARNING?



- This is Recaptcha.
 - ▶ Recaptcha helps stop millions of spam a day.
 - ▶ In some old days, we have to type Captcha texts to distinguish ourself from bots.

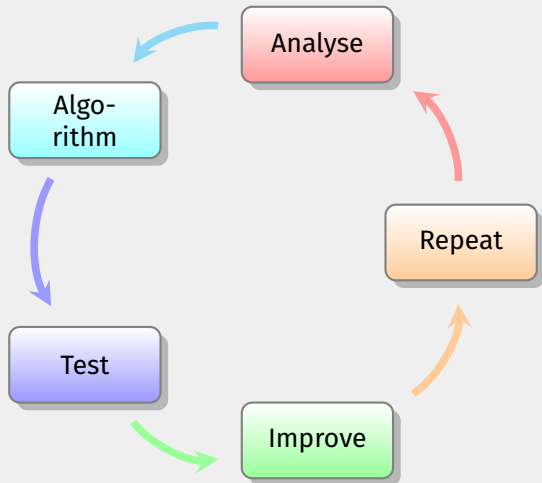
WHAT IS MACHINE LEARNING?



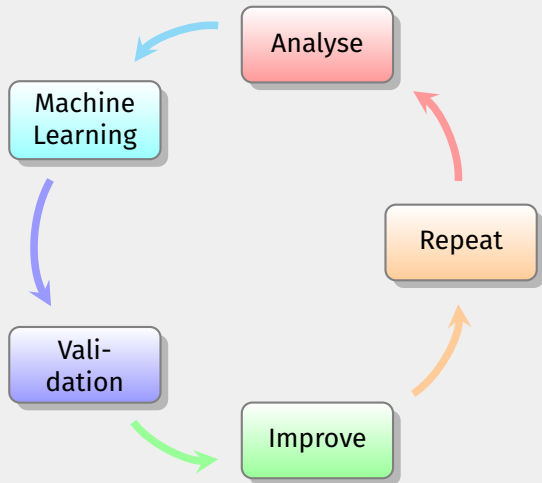
■ This is Recaptcha.

- ▶ Recaptcha helps stop millions of spam a day.
- ▶ In some old days, we have to type Captcha texts to distinguish ourself from bots.
- ▶ How is it possible that with a single click, an automated system can distinguish bots from humans?

TRADITIONAL PROGRAMMING APPROACH



MACHINE LEARNING APPROACH



Machine Learning

Machine Learning
= Data + Data analysis algorithm

Machine Learning
= Data + Data analysis algorithm
= Adapt to change

MACHINE LEARNING PROBLEMS

TYPES OF MACHINE LEARNING PROBLEMS

TYPES OF MACHINE LEARNING PROBLEMS

1. Supervised learning

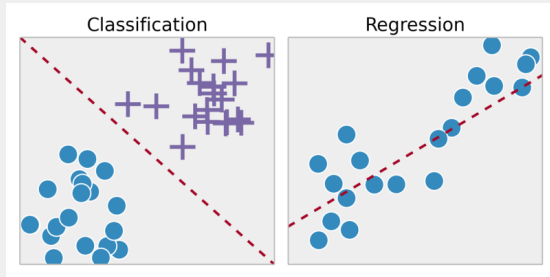
TYPES OF MACHINE LEARNING PROBLEMS

1. Supervised learning
2. Unsupervised learning

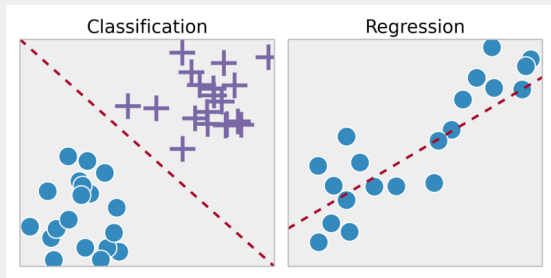
TYPES OF MACHINE LEARNING PROBLEMS

1. Supervised learning
2. Unsupervised learning
3. Reinforcement learning

SUPERVISED LEARNING

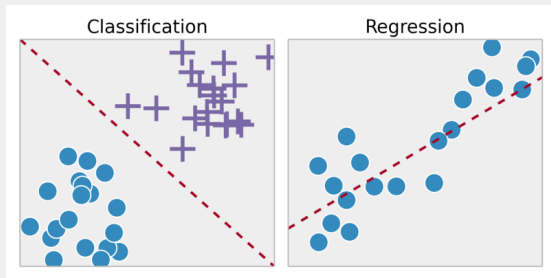


SUPERVISED LEARNING



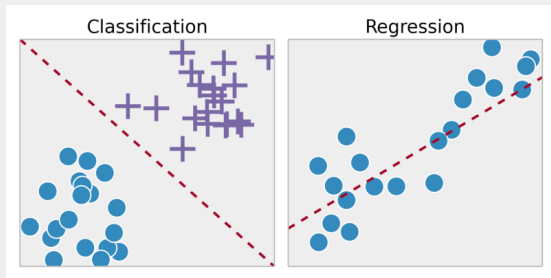
- Given a **training set** for the data, find a **model** to **generalise** well to **unseen** data.

SUPERVISED LEARNING



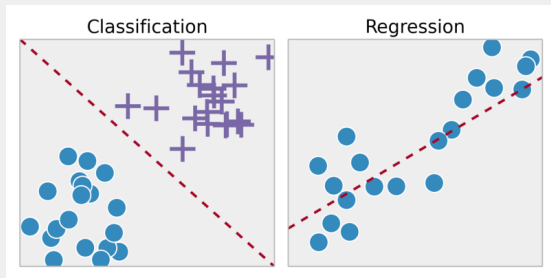
- Given a **training set** for the data, find a **model** to **generalise** well to **unseen** data.
- Two main supervised learning problems

SUPERVISED LEARNING



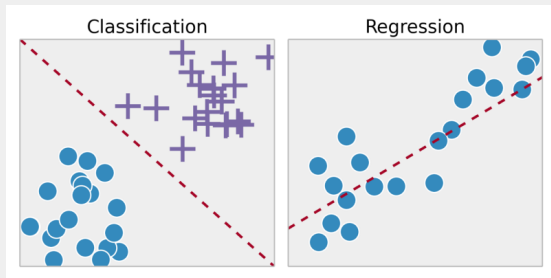
- Given a **training set** for the data, find a **model** to **generalise** well to **unseen** data.
- Two main supervised learning problems
 - ▶ Classification: On the discrete data

SUPERVISED LEARNING



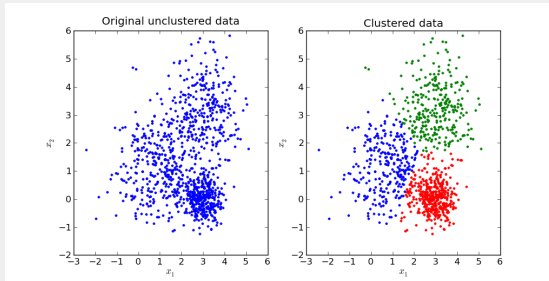
- Given a **training set** for the data, find a **model** to **generalise** well to **unseen** data.
- Two main supervised learning problems
 - ▶ Classification: On the discrete data
 - ▶ Regression: On the continuous data

SUPERVISED LEARNING

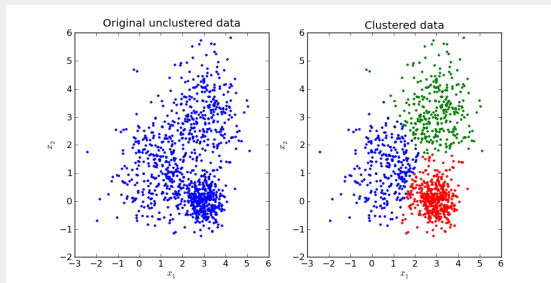


- Given a **training set** for the data, find a **model** to **generalise** well to **unseen** data.
- Two main supervised learning problems
 - ▶ Classification: On the discrete data
 - ▶ Regression: On the continuous data
- Example problems: Spam E-mail detection, Facial recognition

UNSUPERVISED LEARNING

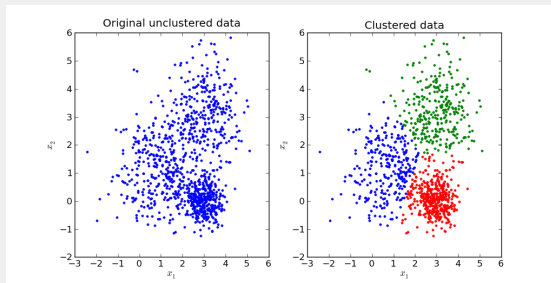


UNSUPERVISED LEARNING



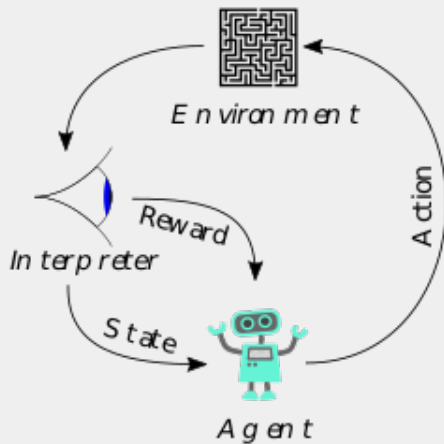
- Discover **hidden** structure in **non-labelled** data.

UNSUPERVISED LEARNING



- Discover **hidden** structure in **non-labelled** data.
- Example: Clustering, Generative models

REINFORCEMENT LEARNING



MODEL

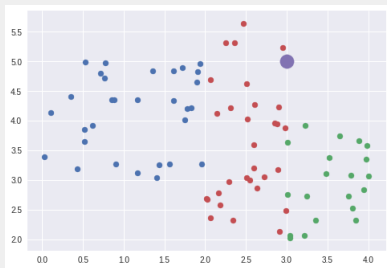
- A result of the combination between...

- A result of the combination between...
 - ▶ a **method** to recognise the data, and

- A result of the combination between...
 - ▶ a **method** to recognise the data, and
 - ▶ **sample datas** for such the method

MODEL

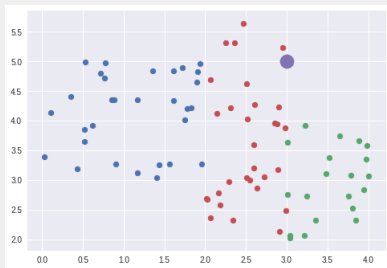
- A result of the combination between...
 - ▶ a **method** to recognise the data, and
 - ▶ **sample datas** for such the method



Data

MODEL

- A result of the combination between...
 - ▶ a **method** to recognise the data, and
 - ▶ **sample datas** for such the method



Data

Determine which group should the purple dot be in (red/green/blue) by **checking the colour of its nearest dot.**

Method

BEGINNING WITH OUR FIRST MODEL

- We're going to write our **first own** machine learning algorithm called ***k*-Nearest Neighbour** (*k*-NN)

- We're going to write our **first own** machine learning algorithm called ***k*-Nearest Neighbour** (*k*-NN)
 - ▶ *k*-NN is known to be very simple, with its concept as

- We're going to write our **first own** machine learning algorithm called ***k*-Nearest Neighbour** (*k*-NN)
 - ▶ *k*-NN is known to be very simple, with its concept as

k-NN algorithm

To classify label of a data point, get *k* nearest data points to the data point, and select the major label among those data points.

Coding time!

MACHINE LEARNING PROCESS

MACHINE LEARNING PROCESS

- Train
- Test

(There'll be more of this, trust me.)

CHOOSING THE PARAMETER FOR k -NN ALGORITHM

What is the bad way to choose k ?

- What if we choose $k = \#$ of all points?
 - ▶ What will happen if our dataset's got 3 labels of A, B, C with 10, 20, and 30 data points of each?
 - ▶ Answer: Our model will always answer the labels with the highest data point count.
- What if we choose $k = 1$?
 - ▶ Let's try!

Coding time!

TRAINING AND TESTING SET

- We separate our dataset into 2 parts: the **training set** and **testing set**

- We separate our dataset into 2 parts: the **training set** and **testing set**
 - ▶ Most of the time, the testing set will be around 10-25% of the entire dataset

- We separate our dataset into 2 parts: the **training set** and **testing set**
 - ▶ Most of the time, the testing set will be around 10-25% of the entire dataset
 - ▶ What will happen if we train on the testing set?

- We separate our dataset into 2 parts: the **training set** and **testing set**
 - ▶ Most of the time, the testing set will be around 10-25% of the entire dataset
 - ▶ What will happen if we train on the testing set?
 - ▶ What will happen if we test on the training set?

- We separate our dataset into 2 parts: the **training set** and **testing set**
 - ▶ Most of the time, the testing set will be around 10-25% of the entire dataset
 - ▶ What will happen if we train on the testing set?
 - ▶ What will happen if we test on the training set?
 - **Cheating!** Like letting the model *remembers* the answer instead of **generalising** the data pattern.

- We separate our dataset into 2 parts: the **training set** and **testing set**
 - ▶ Most of the time, the testing set will be around 10-25% of the entire dataset
 - ▶ What will happen if we train on the testing set?
 - ▶ What will happen if we test on the training set?
 - **Cheating!** Like letting the model *remembers* the answer instead of **generalising** the data pattern.
 - ▶ In other words, **don't test and train model on the same set of data.**

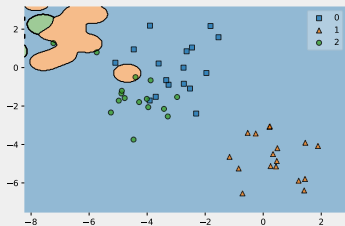
CHOOSING THE BEST k

- **Train** with the training set, to let our model know how will the data looks like.
- **Test** with the testing set, to see on how our model performs.

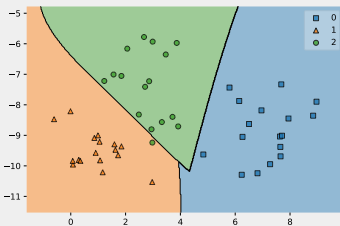
Warning! This is a simplified Machine Learning model training process, there are more to concerns!

OVERFITTING AND UNDERFITTING

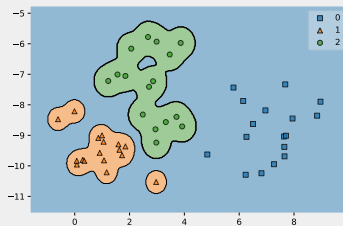
Which decision region is good?



Underfit: The model fails to recognise data pattern



Good fit: The model recognises data pattern generally



Overfit: The model **remembers** data pattern instead of generalising.

Good model must **generalise**