

Machine learning reading group, talk 1 (Introduction)

what is the reading group about?

A brief introduction to the core concepts in machine learning, assuming no prior knowledge of statistics or programming, but a strong (master/phd) background in mathematics.

Goals: By end of reading group should be able to:

- ① understand main tasks ML can be used for
- ② understand selection of most common methods in ML
- ③ be able to implement simple algorithms in python and use prebuilt packages for more complex methods
- ④ have an awareness of the theoretical guarantees and possible failure of different methods

Structure of reading group

Each talk : ~30 mins theory + ~30 mins practical implementation in python

In first talks, the practical section will cover setting up python and basic use of relevant packages.

For rest of term : Part I (weeks 3-8) : Supervised learning
Part II (weeks 9-10) : Unsupervised learning

For the rest of the talk we will explain what supervised and unsupervised learning is, and why they are useful.

what is ML?

In a nutshell, using algorithms which learn from data to solve tasks.

Example (Supervised learning)

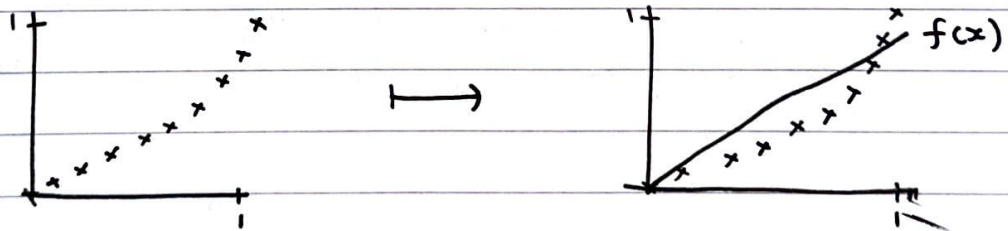
Let X and Y be metric spaces, $f^+ : X \rightarrow Y$ a function. Assume we want to approximate f^+ , which we don't know using data

$$\{(x_1, f^+(x_1)=y_1), \dots, (x_n, f^+(x_n)=y_n)\} \subset X \times Y$$

A supervised-learning algorithm takes this data as an input and returns a function $f : X \rightarrow Y$ for which $d_Y(f^+(x), f(x))$ is small, (at least for x that is close to training data).

Concrete example

$$X = Y = \mathbb{R}^2, f^+(x) = x^2, \text{ data: } \left\{ \left(\frac{i}{N}, f^+\left(\frac{i}{N}\right) \right) \right\}_{i=0}^N$$



least squares regression is a supervised ML algorithm that tries to find $f(x) = ax+b$ to minimise

$$\sum_{i=0}^N |(ax_i + b) - f^+(x_i)|^2$$

This makes the error between f^+ and f small for $x \in [0, 1]$ ("interpolation") but the error will be big for $x \gg 1$ ("extrapolation").

Practical example (see code session)

We'll see a dataset where $X = Y = \mathbb{R}^3$ where we have $N=20$ data pairs $\{(x_1, y_1), \dots, (x_N, y_N)\}$ with

$x_i = (\text{\# chin ups}, \text{\# sit ups}, \text{\# jumps})$

$f^+(x_i) = y_i = (\text{weight}, \text{waist}, \text{pulse})$

of person i in a physiological experiment.

Goal: learn underlying mapping $f^+: X \rightarrow Y$, for at least "realistic" values of $x \in X$ similar to training data.

Remark: In next talk we will consider a more sophisticated approach where we don't assume a one-to-one mapping ~~then~~ $f^+: X \rightarrow Y$.

Another example (classification of images)

The previous examples have all been on regression, predicting a continuous quantity e.g. $y \in Y = \mathbb{R}^3$. A related problem, which is mathematically the same, is predicting a discrete quantity e.g.

$X = \{\text{set of } d \times d \text{ pixelated images of either cats or dogs}\}$
 $\subseteq [0, 1]^{d \times d \times 3}$ (red, green or blue (RGB))

$Y = \{\text{cat}, \text{dog}\}$

$f^+: X \rightarrow Y$

Goal: use data $\{(x_i, y_i)\}_{i=1}^N$ to learn the map that gives each image the correct label.

this is typically done by convolutional neural networks (week 7)

Part I will be dedicated to supervised learning. But there are lots of tasks other than supervised learning, this will be the focus of part II (§2-7 Bach)

The distinction is between supervised and unsupervised with semi-supervised learning in between (will discuss).

Supervised learning (e.g. classification, regression)	Semi-supervised learning	Unsupervised learning (e.g. generative modelling, clustering, dimension reduction, ...)
all data $x \in X$ has a label $y \in Y$	some data have known labels, some don't	none of data $x \in X$ have labels

Summary of unsupervised learning methods

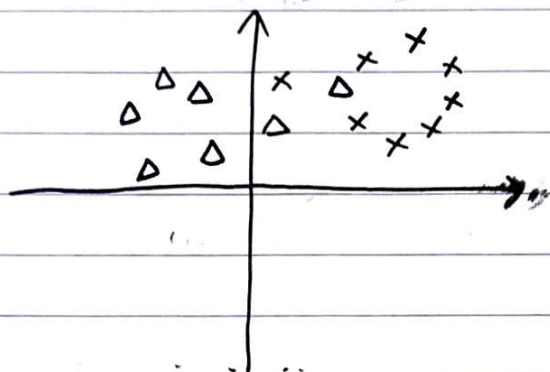
① Generative modelling: Assume $\{x_i\}_{i=1}^N$ are independent samples from some prob. distribution μ^+ on the metric space X . would like to learn a prob. dist. μ on X that is "close" to μ^+ in some distance or space of prob. dist. on X .

e.g. Kullback-Leibler, divergence-area, relative entropy or Wasserstein distance.

Practical ability: large language models (LLMs), let $\{x_i\}_{i=1}^N$ be sentences downloaded from internet, learning μ allows us to generate new sentences

Remark: for a useful LLM, use conditional generation, don't just generate random sentences, generate something conditioned on the user's inputs

- ② Clustering: Given data $\{x_i\}_{i=1}^N$ which we postulate arise from K "clusters", come up with a decision rule to separate the data.



Suppose there are two distinct groups X and Δ , we know there are 2 groups but not the labels X and Δ , just coordinates $x_i \in \mathbb{R}^2$.

want to learn which data gets which label.

- ③ Dimension reduction: Given data $\{x_i\}_{i=1}^N \sim \mu^+$, learn a ~~map~~ map $f: X \rightarrow Z$ with $\dim Z \ll \dim X$ where $f(x) \in Z$ meaningfully preserves some features of $x \in X$.

Example (Autoencoders, week 10)

Given $\{x_i\}_{i=1}^N \sim \mu^+$, learn maps $f: X \rightarrow Z$ and $g: Z \rightarrow X$ such that $g(f(x)) \approx x$ for $x \sim \mu^+$.