

Key for Moodle: MSFML25

Tutorials: Thursdays, 2:15, Ther. 39, B133

Exam: oral (likely), $\sim 10-14$ days after July 26,
2-4 weeks.

Lecture notes uploaded.

ML: Computational method to convert experience \leadsto expertise.

- Experience: past data.
- Expertise: prediction of future outcomes

Std learning tasks: Classification, regression, clustering, dimensionality reduction / manifold learning - find lower dimnl. rep of data while preserving their props.

This course: \hookrightarrow statistical learning theory

1. Thtdl foundations. PAC Learning Framework, Rademacher Complexity, VC-dimension
 2. Analysis of ML methods & algos (w/ apps. of p1).
- SVM & Kernel methods
 - Boosting
 - Logistic regression
 - SGD
 - NNs (DL)
 - May be part 2.

Basic ML Setup.

Input space X , eg \mathbb{R}^n , $[0, 1]^2, \dots$

Output space Y , eg $\{0, 1\}$, \mathbb{R}, \dots

Training data (x_i, y_i) , $i = 1, \dots, m$ (labeled data)
or x_i , $i = 1, \dots, m$ (unlabeled)

Hypothesis class \mathcal{H} : a set of functions $h: X \rightarrow Y$.

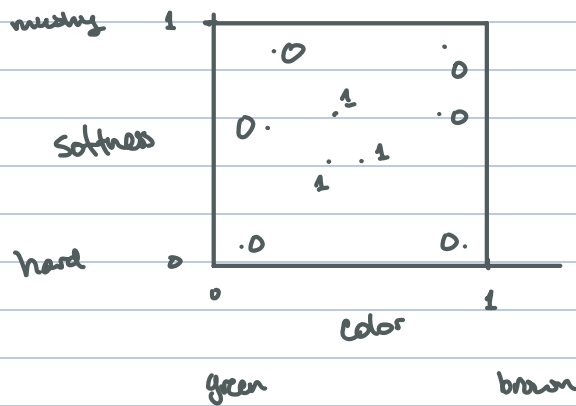
What class of functions to learn? Why do we need to choose it?

Suppose you want to predict whether a papaya is tasty or not.

Based on experience w/ other fruits, we assume that taste depends on color & softness.

Assume we measure both on an interval $[0, 1]$.

$$\leadsto X = [0, 1]^2$$



0 $\hat{=}$ not tasty
1 $\hat{=}$ tasty

→ Try a couple fruits, measure softness & color & label w/ 0 or 1.

Goal: find a function $h: [0, 1] \rightarrow [0, 1]$ s.t.

$$h(a,b) = \begin{cases} 1 & \text{if papaya w/ color } a, \text{ softness } b \\ & \text{is typically tasty,} \\ 0 & \text{if } \rule{1cm}{0.4pt} \text{not tasty.} \end{cases}$$

The pairs $(x_i, y_i) = ((a_i, b_i), y_i)$ are distributed according to some unknown distribution D , representing environ.

Given $(x_1, y_1), \dots, (x_n, y_n) \in [0, 1]^2 \times \{0, 1\}$,
how should we choose h ?

h should satisfy $h(x_i) = y_i$ for (at least almost) all (x_i, y_i) .

There may be many such functions.

Ex. $h(x) = \begin{cases} 1 & \text{if there is a training eg } (x_i, y_i) \text{ s.t. } y_i = 1 \\ 0 & \text{o/w} \end{cases}$

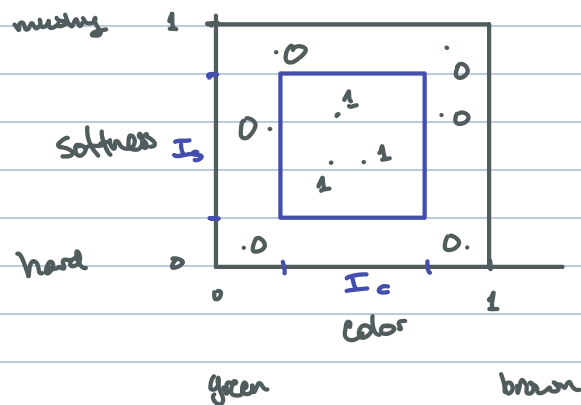
Does not seem realistic. (\rightarrow overfitting).

Then $h(x_i) = y_i \forall i$ (x_i distinct).

Other possibility (more realistic): assume there are intervals $I_c \subset [0, 1]$ and $I_s \subset [0, 1]$ s.t. papayas w/ color value $\in I_c$, softness value $\in I_s$ are (typically) tasty. Other papayas not tasty.

$$h(a, b) = \begin{cases} 1 & \text{if } (a, b) \in I_c \times I_s, \\ 0 & \text{o/w.} \end{cases} \leadsto h = \mathbb{1}_{I_c \times I_s}.$$

Want to define class H from which we choose h , eg H could be the set of characteristic functions of axis-aligned rectangles.



How to choose such a rectangle?

Choosing a suitable class H requires preknowledge about the problem.

In general, H should not be the class of all functions

$$[0, 1]^2 \rightarrow \{0, 1\}$$

1. Supervised.

Learning receives labeled data $(x_i, y_i) \leadsto$ tries to find function

$$h(x_i) \approx y_i,$$

$i = m+1, m+2, \dots$

Example: $x_i = ([0, 255] \cap \mathbb{Z})^{n_1 \times n_2}$, $y_i = \{1, 0\}$ cat in image or not
 $3n_1 \times n_2 \leadsto$ if you make color channels.

2. Unsupervised.

Learner receives unlabeled data $x_i, i=1, \dots, m \rightarrow$ understand structure of data points. They may be contained in a much smaller submanifold/subspace, or cluster.

3. Semisupervised

Learner $\leftarrow (x_i, y_i) ; i=1, \dots, m ;$ unlabeled $x_i, i=m+1, \dots, N$.
Predict labels of unlabeled data; helps learn struct of data.

4. Online learning.

Mult. training rounds. Receive x_i , predict $\hat{y}_i \rightarrow$ receive real label $y_i \rightarrow$ incur loss $l(y_i, \hat{y}_i)$. Goal: minimize cumulative loss

$$\sum_{i=1}^m l(y_i, \hat{y}_i).$$

5. RL

Mixed training/testing phases. Learner \rightarrow actions \rightarrow receives immediate reward/loss. Goal: maximize reward over course of actions.

\hookrightarrow explore vs. exploit dilemma.

\swarrow
get into abt environ,
unexplored action

\searrow known action,
exploit already collected
information

6. Active Learning

Learner can decide on data points x_i to query the label y_i .

Eg, scientific experiments, oil exploration.

Hope: better predictions/less samples.

PAC Learning Framework.

Consider supervised learning scenario for binary classification

X : input, $Y = \{0, 1\}$ output, training data

$$S = ((x_1, y_1), \dots, (x_m, y_m)) \in (X \times Y)^m$$

Assumptions.

1. $x_i, i=1, \dots, m$ drawn iid, according to some unknown prob. dist. D on X .
2. Labels are given $y_i = f(x_i), i=1, \dots, m$ for some map $f: X \rightarrow Y$.

Goal: find f (at least approximately).

More generally (later): assume (x_i, y_i) drawn iid from D on $X \times Y$.

Learner's output: a prediction rule

$$h: X \rightarrow Y,$$

ideally $h=f$.

Learner selects hypothesis from

$$\mathcal{H} \subset \{g: X \rightarrow Y\}$$

and f may or may not be contained in \mathcal{H} .

Assumptions on measurable spaces: for msl spaces X, Y, Z, \dots

1. If space is finite/countably infinite, it is equipped w/ the σ -algebra consisting of all subsets of the space.
2. Otherwise, assume it is a metric space, complete & separable, equipped w/ corresponding Borel σ -alg.
3. For products $X \times Y$, carries the product σ -alg of X and Y .

Generalization Error. For $h: X \rightarrow Y = \{0,1\}$, the generalization error

or risk of h is:

$$R(h) = \mathbb{P}_{x \sim D}(h(x) \neq f(x)) = \mathbb{E}[\mathbb{1}_{\{x | h(x) \neq f(x)\}}]$$

→ note $P(A) = E[\mathbb{1}_A]$

"prob. of a set" $\rightarrow = \int_{\Omega} \mathbb{1}_A dP(\omega)$
 $= \int_A 1 dP(\omega)$

- Above, we assumed a fixed labeling function $f: X \rightarrow Y = \{0, 1\}$ and probability distribution D on X .