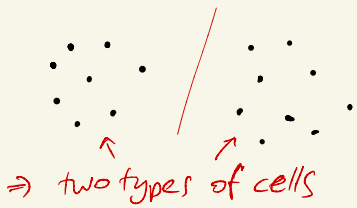# LECTURE 24

## Machine Learning

- What is learning, understanding, attention, experience?
- How do we make technology achieve that?
- Math. models? Based on h.d. probability.

① <u>Unsupervised</u> learning — from own experience (infant). <u>Supervised</u> - from a teacher.
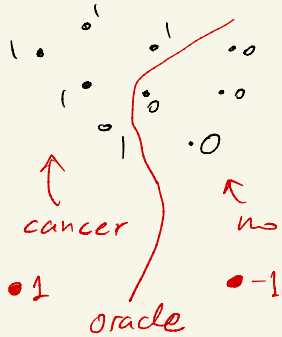
<u>Examples</u> we have seen before:

(a) <u>Unsupervised learning</u>: clustering



⇒ two types of cells

<u>Unlabeled</u> data $\quad x_1, \ldots, x_n \in \mathbb{R}^d$

e.g. $n$ cells, $d$ genes

(B) <u>Supervised learning</u>: classification



cancer $\quad$ no

• 1 $\qquad$ • -1

oracle

<u>Labeled</u> data $\quad (x_1, Y_1), \ldots, (x_n, Y_n) \in \mathbb{R}^d \times \{0,1\}$
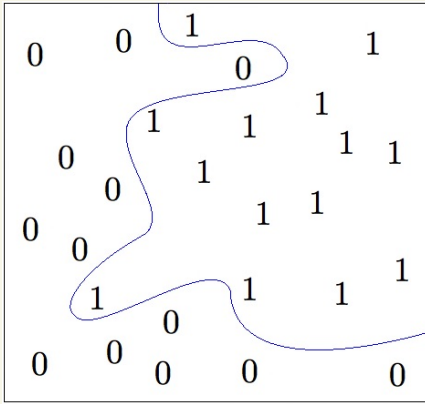
e.g. $n$ people, $d$ symptoms cancer/no "Training data"

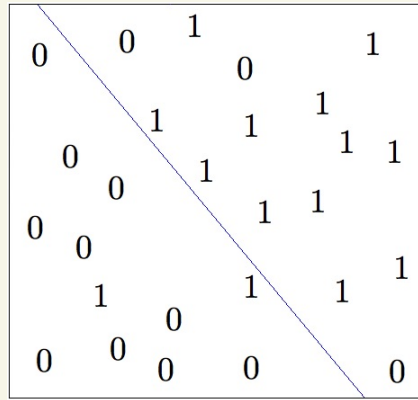Want to build an "<u>oracle</u>"
that makes a diagnosis
for a new patient: $\quad x_{n+1} \longmapsto Y_{n+1}$

# Supervised ML : a general framework

- A pair of random variables (or vectors) $(X, Y) \in \mathcal{X} \times \mathcal{Y}$.

  <span style="color:red">↑ label</span>  <span style="color:red">↑ ↑ & sets</span>

Ex. $(X, Y) \in \mathbb{R}^d \times \{0, 1\}$ as above. Objective reality.

<span style="color:red">↑ symptoms</span> <span style="color:red">↖ cancer/no</span>  $X, Y$ are correlated, ideally strongly.

- The joint distribution of $(X, Y)$ is <u>unknown</u>. We only see:

- <u>Training data</u> $(X_1, Y_1), \ldots, (X_n, Y_n)$: iid copies of $(X, Y)$.

- <u>Goal</u>: predict $Y$ from $X$ as best as we can.

$\Rightarrow$ We want to construct an <u>oracle</u>

$$h: \mathcal{X} \to \mathcal{Y}: \qquad h(x) \approx Y \qquad\qquad (*)$$

to make predictions for new, unseen data: $h(X_{n+1}) = Y_{n+1}$

<span style="color:red">↑ input</span>  <span style="color:red">↑ output</span>

⓪ How do we quantify the "goodness of fit" $(*)$ ?

- We fix a <u>loss function</u> $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$, e.g. $\ell(t-s) = (t-s)^2$, and

  define the <u>risk</u> (a.k.a. <u>test error</u> <span style="color:red">⤵</span>)

$$R(h) := \mathbb{E}\, \ell(h(X), Y) = \mathbb{E}\left(h(X_{n+1}), Y_{n+1}\right)$$

<u>Examples</u>:

(a) quadratic loss  $\ell(t, s) = (t-s)^2$  $\Rightarrow$  $R(h) = \mathbb{E}\left(h(X) - Y\right)^2$

(b) logistic loss $\to$ 

 bounded

(c) hinge loss (svm)

— 2 —

Ⓠ How do we construct an oracle $h$ ?
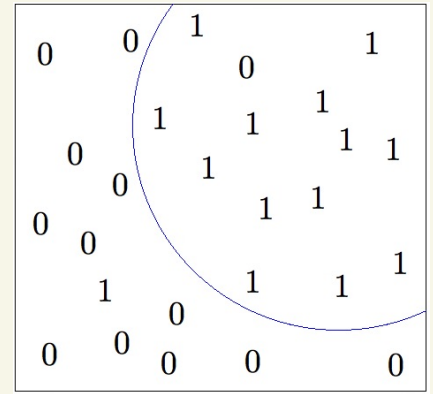
(1) **h too complex:** **overfitting**

(2) **h too simple:** **underfitting**

(3) **h is OK:** **good fit**



OUR STRATEGY:

1. Fix some <u>collection</u> of functions $\mathcal{H}$, called a <u>hypothesis class</u>.
2. Select $h \in \mathcal{H}$ that best fits the training data.

<u>Examples</u>:

(a) $\mathcal{H} = \{ \text{all functions } h: X \to Y \} \Rightarrow$ overfitting (1)

(b) All <u>linear</u> functions: $\mathcal{H} = \{ h(x) = \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R} \}$.

   $\Rightarrow$ linear regression

(c) $\mathcal{H} = \{ h(x) = \text{sign}(\langle w, x \rangle + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \} \Rightarrow$ SVM (2)

(c) $\mathcal{H} = \{ \text{all polynomials } p(x) \text{ of degree} \leq 2 \}$   (3)

(d) $\mathcal{H} = \{ \text{all functions realized by a given neural network architecture} \}$

(e) $\mathcal{H} = \{ h_1, h_2 \} \Rightarrow$ hypothesis testing

No systematic way to choose $\mathcal{H}$.  ("Model selection")

- The best $h \in \mathcal{H}$ is the one that minimizes the risk

$$R(h) = \mathbb{E}\, \ell(h(x), Y).$$

$$h^* := \operatorname*{argmin}_{h \in \mathcal{H}} R(h).$$

- But $R(h)$ can't be computed (can't take $\mathbb{E}$ over the <u>population</u>)

<u>empirical risk</u> (a.k.a. <u>training error</u>)

↓

$$R_n(h) := \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), Y_i).$$

$$h_n^* := \operatorname*{argmax}_{h \in \mathcal{H}} R_n(h)$$

↳ <u>can</u> be computed from training data ↗ (can be NP hard)

<u>Ex</u> (Binary classification), quadratic loss ⟹

$$R_n(h) = \frac{1}{n} \sum_{i=1}^{n} \underbrace{(h(x_i) - Y_i)^2}_{\parallel} = \% \text{ of misclassified training data.}$$

$$\begin{cases} 1 & \text{if } h(x_i) \neq Y_i \\ 0 & \text{otherwise} \end{cases}$$

**Empirical Risk Minimization (ERM) Algorithm**

① Training: for input data $(x_1, Y_1), \ldots, (x_n, Y_n)$; compute $h_n^*$.

② Prediction: on query $X$, output $h_n^*(X)$   "oracle"

How do we measure the quality?

- Generalization error $:= R(h_n^*)$

measures how well the algorithm _generalizes_ to unseen data.

- Examples :

  (a)   $\mathcal{H} = \{ \text{all functions} \}$, $Y = f(x)$.

      $\exists$ a perfect fit to the training data: $h_n^*(x) := \begin{cases} Y_i & \text{if } x = X_i \\ 0 & \text{elsewhere.} \end{cases}$

       training error $R_n(h_n^*) = 0$. (Overfitting)
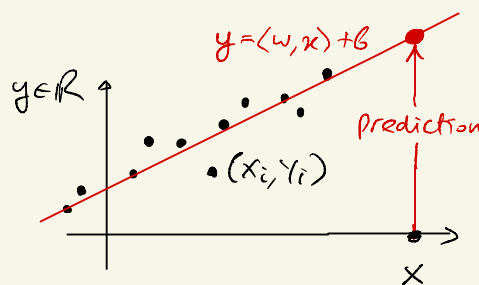
      BUT does NOT generalize well :

           $R(h_n^*)$ is large.      _Memorizes, not generalizes._

  (b)  $\mathcal{H} = \{ \text{all linear functions} \}$, quadratic loss $\Rightarrow$

$$W_n^* = \underset{W \in \mathbb{R}^d,\, b \in \mathbb{R}}{\arg \min} \; \frac{1}{n} \sum_{i=1}^{n} \left( \langle W, X_i \rangle + b - Y_i \right)^2$$

        $=$ linear regression.   OK.



_Our goal_ : Bound the generalization error.

           How does it depend on the "complexity" of $\mathcal{H}$ ?