

LECTURE 34

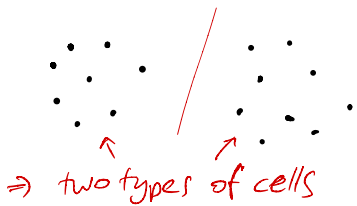
Machine Learning

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

① Unsupervised learning - from own experience (infant). Supervised - from a teacher.

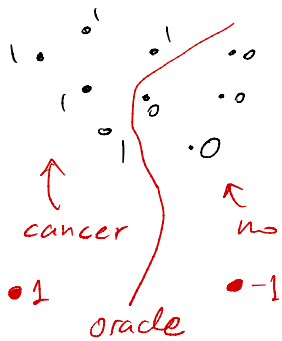
Examples we have seen before:

(a) Unsupervised ML: clustering



Unlabeled data $x_1, \dots, x_n \in \mathbb{R}^d$
e.g. n cells, d genes

(b) Supervised ML: classification



Labeled data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{0, 1\}$
e.g. n people, d symptoms cancer/no "Training data"

want to build an "oracle"

that makes a diagnosis

of a new person: $x_{n+1} \mapsto y_{n+1}$

Supervised ML: a general framework

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

Ex. $(X, Y) \in \mathbb{R}^d \times \{0, 1\}$ as above. Objective reality.
 \uparrow symptoms \uparrow cancer/no X, Y are correlated, ideally strongly.

- The joint distribution of (X, Y) is unknown. We only see:
- Training data $(x_1, y_1), \dots, (x_n, y_n)$: iid copies of (X, Y) .
- Goal: predict Y from X as best as we can.

\Rightarrow We want to construct an oracle

$$h: \mathcal{X} \rightarrow \mathcal{Y}: \quad h(x) \approx Y \quad (*)$$

to make predictions for new, unseen data: $h(x_{n+1}) = y_{n+1}$
 \uparrow input \uparrow output

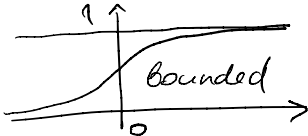
① How do we quantify the "goodness of fit" (*)?

- We fix a loss function $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$,
define the risk (a.k.a. test error)

$$R(h) := \mathbb{E} l(h(x), Y) = \mathbb{E} (h(x_{n+1}), y_{n+1})$$

Examples:

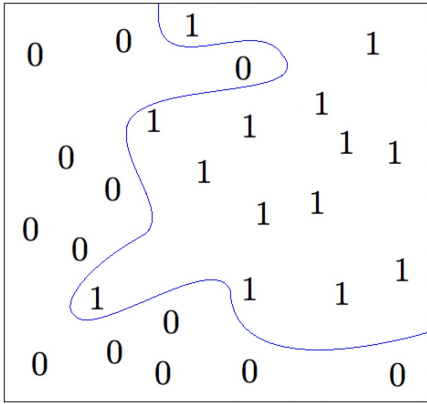
(a) quadratic loss $l(t, s) = (t-s)^2 \Rightarrow R(h) = \mathbb{E} (h(x) - Y)^2$

(b) logistic loss \rightarrow  bounded

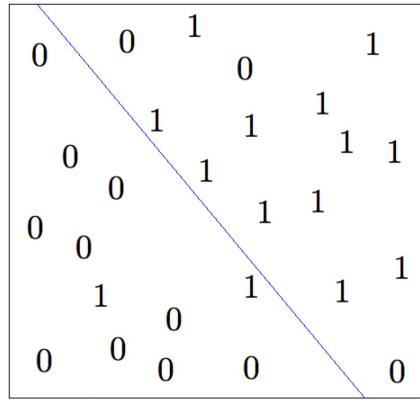
(c) hinge loss (svm)

Q How do we construct an oracle h ?

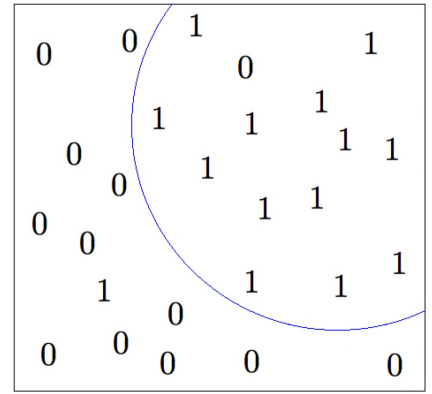
(a) h too complex:
overfitting



(b) h too simple:
underfitting



(c) h is OK:
good fit



OUR STRATEGY:

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

Examples:

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

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

\Rightarrow linear regression

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

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

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

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

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

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

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

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

- But $R(h)$ can't be computed. Estimate it by the empirical risk (a.k.a. training error)

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

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

↑ can be computed from training data ↗ (can be NP hard)

Ex (Binary classification), quadratic loss $\Rightarrow R_n(h) = \frac{1}{n} \sum_{i=1}^n (h(x_i) - y_i)^2 = \%$ of misclassified data.
"1 if $h(x_i) \neq y_i$ "

- Generalization error $R(h_n^*) - R(h^*) \geq 0$
measures how well the algorithm generalizes

- Examples:

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

\exists a perfect oracle $h^* = f$, whose risk $R(h^*) = 0$.

\exists a perfect fit to the training data $h_n^*: X_n \rightarrow Y_n$, so training error $R_n(h_n^*) = 0$. (Overfitting)

BUT does NOT generalize well:

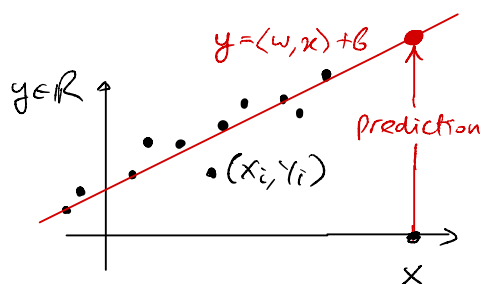
$R(h_n^*)$ is large.

Memorizes, not generalizes.

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

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle + b - y_i)^2$$

= linear regression.



Our goal : bound the generalization error.

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

$$\underline{\text{Lem}} \quad R(h_n^*) - R(h^*) \leq 2 \sup_{h \in \mathcal{H}} |R_n(h) - R(h)|$$

Proof

$$\begin{aligned} R(h_n^*) &\leq R_n(h_n^*) + \varepsilon && (h_n^* \in \mathcal{H}) \\ &\leq R_n(h^*) + \varepsilon && (h^* = \text{minimizer of } R_n) \\ &\leq R(h^*) + 2\varepsilon && (h^* \in \mathcal{H}) \quad \text{QED.} \end{aligned}$$