

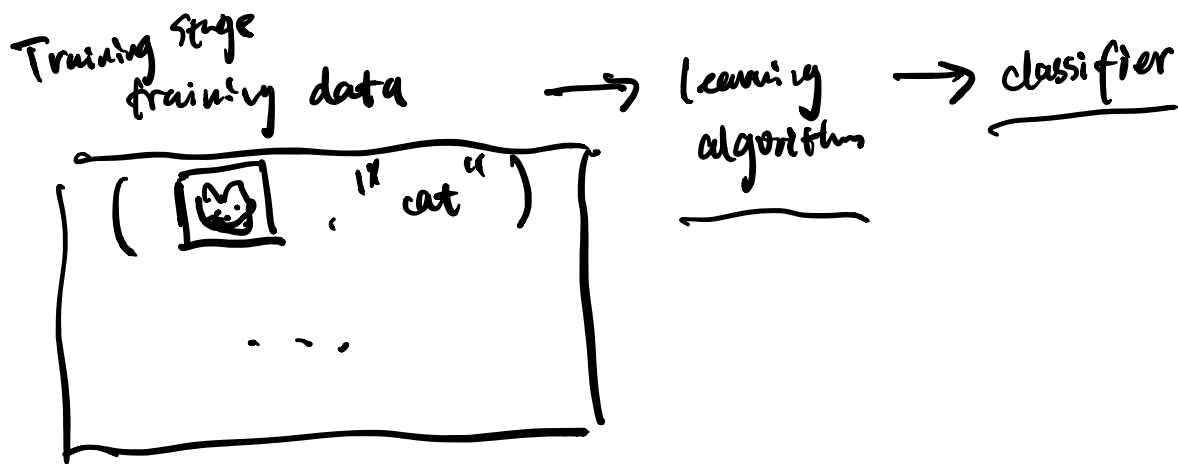
Announcements

1. Calibration HW due Friday (22nd)
2. Scribe note sign-up sheet (add a link in class website)

Part 1 Statistical Learning

- PAC learning (Probably Approximately Correct)

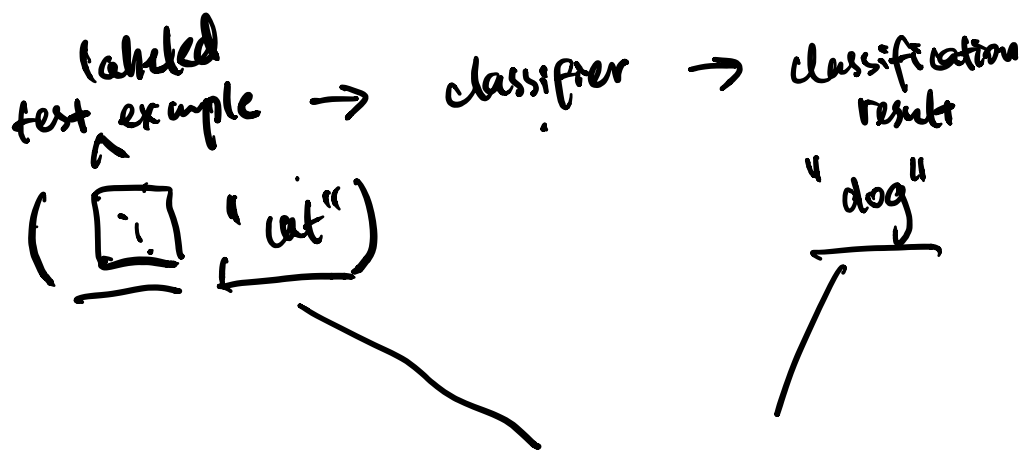
Batch (supervised learning) pipeline



Test stage



Test stage (another view)



match?

Assumption: training & testing examples come from the same distribution.

Qns:

1. How can ^{we} formalize this process mathematically?
2. Can we formally define what a good predictor / good learning algorithm should look like?

PAC Model (Valiant, 1984)

- instance domain \mathcal{X} : e.g. images. $\mathcal{X} = \mathbb{R}^{64 \times 64 \times 3}$
- label space $\mathcal{Y} = \{-1, +1\}$ $+1$: "cat", -1 : "dog"
- data distribution D , supported on $\mathcal{X} \times \mathcal{Y}$.
- training set S : drawn from D .

$(x_1, y_1) \dots (x_m, y_m)$ independent & identically distributed (i.i.d.)
 $\sim D$

- classifier (hypothesis) $h: \mathcal{X} \rightarrow \mathcal{Y}$

 $\rightarrow h \rightarrow$ "cat"

- Hypothesis class (concept class): a structured collection of classifiers.

$\mathcal{H} = \{ \text{all neural nets with ResNet-18 architecture} \}$

- Performance evaluation:

h (classifier)'s generalization error ϵ

$$P_{(x,y) \sim D} (h(x) \neq y) = \text{err}(h, D)$$

h 's training error:

$$\frac{\sum_{i=1}^m \mathbb{I}(h(x_i) \neq y_i)}{m} = \text{err}(h, S)$$

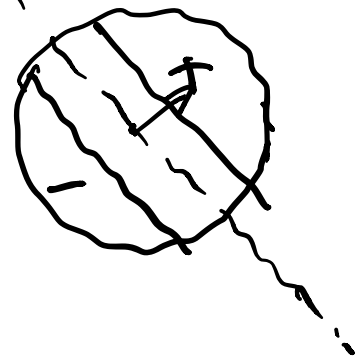
$$= P_{(x,y) \sim \text{unif}(S)} (h(x) \neq y)$$

$$\frac{\sum_{(x,y) \in S} \mathbb{I}(h(x) \neq y)}{|S|}$$

$$\mathbb{I}(A) = \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{if } \dots \text{ false} \end{cases}$$

— Realizable (separable): we are given the promise that there exists a classifier $h^* \in H$, such that $\text{err}(h^*, D) = 0$.

Agnostic: there may or may not exist a classifier in H w/ zero test error.



Def We call Alg A PAC learns hypothesis class \mathcal{H} with sample complexity function: $f: (0,1) \times (0,1) \rightarrow \mathbb{N}$, if for any distribution D realizable with respect to \mathcal{H} , $\epsilon, \delta > 0$, when A receives $m \geq f(\epsilon, \delta)$ ^{i.i.d.} training examples as input, it outputs a classifier \hat{h} , such that with probability $1 - \delta$, $\text{err}(\hat{h}, D) \leq \epsilon$.

Remarks:

1. ϵ : target error rate
 δ : confidence parameter
2. polynomial running time required in Valiant's original paper, but not here.
3. m : sample complexity function; want it to be as small as possible.
 m here is distribution independent
 it may sometimes be better to consider distribution dependent sample complexity results.

$m(z, \delta, D)$

4. If A returns \hat{h} in H ; A is called a proper learning alg.

If A doesn't nec. return \hat{h} ; \hat{h}

improper \hat{h} . . .

Sometimes improper learning is beneficial.

Sample complexity of finite hypothesis classes.

Result: each element can be represented by

$(b \times \# \text{ weights})$ bits.

- H is a finite set.

- Realizability

$(x_1, y_1) \dots (x_n, y_n) \rightarrow \text{Alg} \rightarrow \hat{h}$

~~1~~ The consistency algorithm:

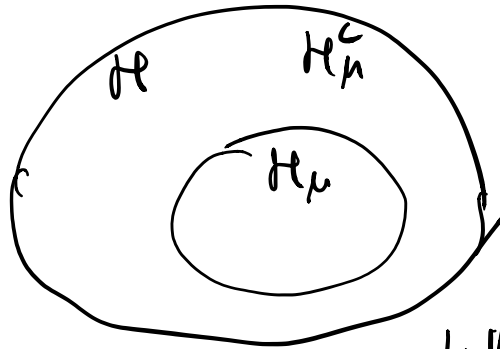
1. Return \hat{h} in H such that it agrees w/ all examples in S , i.e. for all i , $\hat{h}(x_i) = y_i$.

Theorem: For finite H . If the consistency alg. is given m iid training examples from D realizable w.r.t. H , then

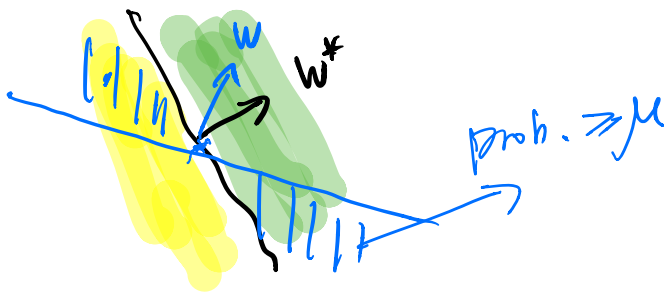
With probability $1-\delta$. Its output \hat{h} is such that

$$\text{err}(\hat{h}, D) \leq \underbrace{\sum_{\mathcal{H}} (m, \delta, |\mathcal{H}|)}_{\text{realizable}} := \frac{\ln|\mathcal{H}| + \ln \frac{1}{\delta}}{m}.$$

In other words. It PAC learns \mathcal{H} with sample complexity

$$f(m, \delta) = \frac{1}{\epsilon} (\ln|\mathcal{H}| + \ln \frac{1}{\delta}).$$


$$\mathcal{H}_\mu = \left\{ h \in \mathcal{H} : \text{err}(h, D) > \mu = \frac{\ln|\mathcal{H}| + \ln \frac{1}{\delta}}{m} \right\}$$



Next class: $P(\mathcal{E}) \geq 1-\delta$,
 show there exists an event $\bar{\mathcal{E}}_\mu$ such that
 for all $h \in \mathcal{H}_\mu$, there will be at least one i .

$$h(x_i) \neq y_i.$$