

GEMS OF TCS

PAC LEARNING

Sasha Golovnev

April 29, 2021

CLASSES OF LEARNING PROBLEMS

- Classification

Binary
Labeled data:

classification: $X \rightarrow \underline{\{0,1\}}$

(x_1, y_1) $\begin{cases} \text{.jpeg} \\ \text{"animal"} \end{cases}$

(x_2, y_2) $\begin{cases} \text{.jpeg} \\ \text{"food"} \end{cases}$

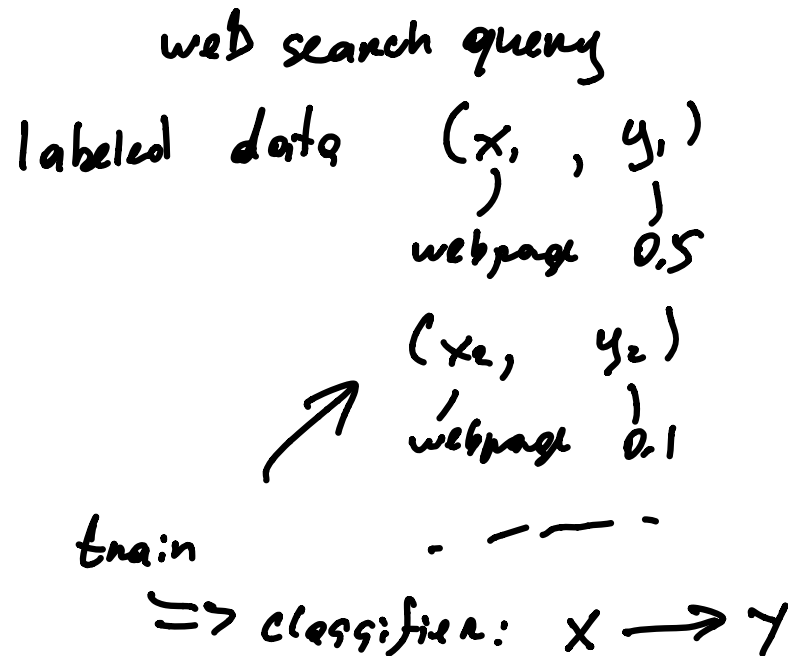
...

(x_m, y_m) $\begin{cases} \text{.jpeg} \\ \text{"food"} \end{cases}$

CLASSES OF LEARNING PROBLEMS

- Classification

- Ranking



CLASSES OF LEARNING PROBLEMS

- Classification
- Ranking
- Regression

classifier: $X \rightarrow \mathcal{R}$

stop values

labeled data: (x_1, y_1)

;

(x_m, y_m)

train

\Rightarrow classifier: $X \rightarrow \mathcal{R}$

CLASSES OF LEARNING PROBLEMS

- Classification

- Ranking

- Regression

- Clustering

labeled data

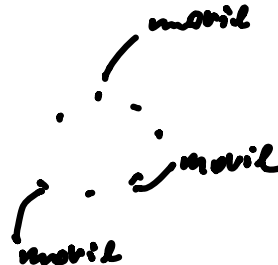
supervised learning

expt. $x_1 \rightarrow y_1$

$x_2 \rightarrow y_2$

unlabeled data

unsupervised learning



- semi-supervised
- online
- active
- reinforcement



CLASSES OF LEARNING PROBLEMS

- Classification
- Ranking
- Regression
- Clustering
- ...

EXAMPLE

Spam filter \equiv binary classification

- Given a set of labeled emails

$(x_1, y_1 \in \{0, 1\})$
|
email \leftarrow spam/non-spam

$(x_2, y_2 \in \{0, 1\})$

.....

EXAMPLE

- Given a set of labeled emails
- Build a classifier that predicts spam/non-spam labels for incoming emails

SETUP

- Partition labeled data into three sets:
 - training sample
 - validation sample
 - test sample

SETUP

- Partition labeled data into three sets:
 - training sample
 - validation sample
 - test sample
- Identify relevant features

- sender
- loans, bank account, prince, business
- nickname, real name
- links to shady webpages
-

good features

usually \Rightarrow good classification

$\Rightarrow v \in \mathbb{R}^n$

SETUP

- Partition labeled data into three sets:
 - training sample
 - validation sample
 - test sample
- Identify relevant features
- Train on training sample

Ex. of our classifier:

given features $V \in \mathbb{R}^n$

classifier: $V_1 \cdot 2 - V_2 \cdot 0.1 + V_3 \cdot 5 + V_4 \cdot 4 + \dots$

spam

$\geq \Theta$
 $< \Theta$

If $\Theta = \text{large} \Rightarrow$ conservative classifier

If $\Theta = \text{small} \Rightarrow$ aggressive classifier non-spam

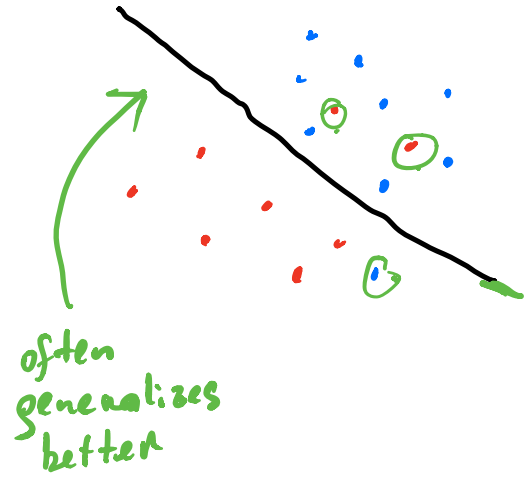
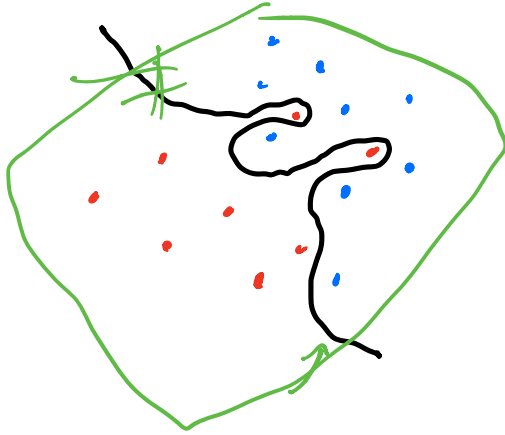
SETUP

- Partition labeled data into three sets:
 - training sample
 - validation sample
 - test sample
- Identify relevant features
- Train on training sample
- Tune parameters using validation sample
For example, choose θ

SETUP

- Partition labeled data into three sets:
 - training sample
 - validation sample
 - test sample
- Identify relevant features
- Train on training sample
- Tune parameters using validation sample
- Evaluate using test sample

$$2 \cdot v_1 - 0.1v_2 + 5v_3 + 4v_4 \geq 5.5$$



often
generalizes
better

WHAT CAN BE LEARNED?

- What can be learned?
- What cannot be learned?
- How many samples do we need to learn?

WHAT CAN BE LEARNED?

- What can be learned?
- What cannot be learned?
- How many samples do we need to learn?
- Framework of PAC learning (L. Valiant, 1984)

DEFINITIONS

- X —set of all possible instances/examples

X = set of emails
= set of feature vectors

DEFINITIONS

- X —set of all possible instances/examples
- \mathcal{D} —target distribution over X

DEFINITIONS

- X —set of all possible instances/examples
- \mathcal{D} —target distribution over X
- c —target concept
- \mathcal{C} —concept class

\mathcal{C} = class of all concepts
ex.: \mathcal{C} = all functions from feature vectors $\rightarrow \{0, 1\}$

$c \in \mathcal{C}$

labeled data

$$x_1 \in X \quad y_1 = c(x_1)$$

$$x_2 \in X \quad y_2 = c(x_2)$$

$$x_m \in X \quad y_m = c(x_m)$$

DEFINITIONS

- X —set of all possible instances/examples
- \mathcal{D} —target distribution over X
- c —target concept
- C —concept class
- H —set of concept hypotheses

e.g., $H =$ linear functions from
Features $\rightarrow \{0,1\}$

$$v_1 \cdot 3 + v_2 \cdot 0.1 + v_3 \cdot 5 + \dots \geq 100$$

DEFINITIONS

- X —set of all possible instances/examples
- D —target distribution over X
- c —target concept
- C —concept class
- H —set of concept hypotheses
- **Goal:** given training set, select h $\in H$ that approximates c well

For a random email ($x \in X$ sampled from D)
 $h(x)$ $= c(x)$ with high probability

ERRORS

Generalization Error

For hypothesis h , target concept c , and target distribution D :

$$R(h) = \Pr_{x \sim D} [\underline{h(x)} \neq \boxed{c(x)}].$$

We want to minimize $R(h)$

ERRORS

Generalization Error

For hypothesis h , target concept c , and target distribution D :

$$R(h) = \Pr_{x \sim D} [h(x) \neq c(x)]$$

— don't know how to compute

Empirical Error

x_1 $y_1 = c(x_1)$
 \vdots
 x_m $y_m = c(x_m)$ *training*

For hypothesis h , target concept c , and sample $S = (x_1, \dots, x_m)$:

$$\hat{R}(h) = \frac{|\{x_i : h(x_i) \neq c(x_i)\}|}{m}$$

— know to compute

AVERAGE ERROR

$$\mathbb{E}_S[\hat{R}(h)] = R(h).$$

$S \sim D^m$

Empirical error

Generalization error

PAC LEARNING

PAC (Probably Approximately Correct)

Concept class C is PAC-learnable if there exists learning algorithm s.t.

- for all $c \in C$, $\epsilon > 0$, $\delta > 0$, all distributions D :
sample S
learning
alg $\Rightarrow h_S$
$$\Pr_{S \sim D^m} [R(h_S) \leq \epsilon] \geq 1 - \delta,$$

$$\epsilon = 0.01 \quad \delta = 0.01$$

PAC LEARNING

PAC (Probably Approximately Correct)

Concept class C is PAC-learnable if there exists learning algorithm s.t.

- for all $c \in C, \epsilon > 0, \delta > 0$, all distributions D :

$$\Pr_{S \sim D^m} [R(h_S) \leq \epsilon] \geq 1 - \delta,$$

approximately *probably*

- for random samples of size

$$m \leq \text{poly}\left(\frac{1}{\epsilon} \frac{1}{\delta} n\right).$$

- Probably: confidence $1 - \delta$
- Approximately correct: accuracy $1 - \epsilon$

⌞
PAC

(two features representing email)
 $x \in \mathbb{R}^2$ EXAMPLE

C = set of axis-aligned rectangles

Is C PAC-learnable?
We'll prove it is

Learning alg:

takes $(x_1 \in \mathbb{R}^2, y_1 = c(x_1))$

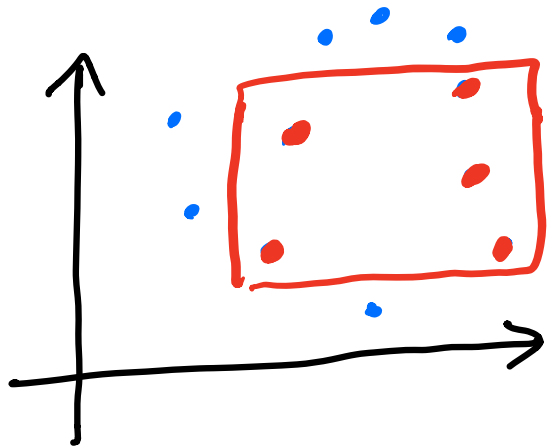
$(x_2 \in \mathbb{R}^2, y_2 = c(x_2))$

$(x_m \in \mathbb{R}^2, y_m = c(x_m))$

returns rectangles h

s.t. w.p. $1 - \delta$

h & c are ϵ -close

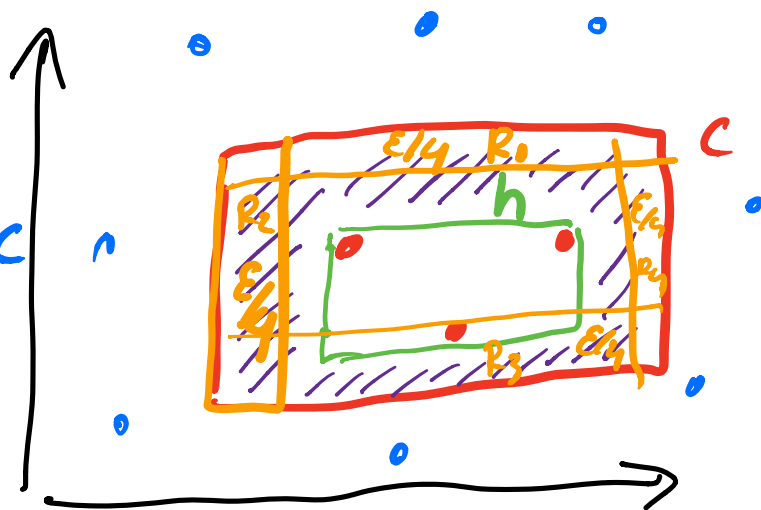


$$m = \text{poly}\left(\frac{1}{\epsilon}, \frac{1}{\delta}\right)$$

m labeled points

we don't know C

Alg smallest axis-aligned rectangle that contains all red points.



WTS: Alg PAC learns C

//// - we'll make errors
 $h(x) \neq C(x)$

IF h is not ϵ -close $C \Rightarrow$
 $\text{area}[C \setminus h] \geq \epsilon \Rightarrow$

h misses R_1 OR R_2 OR R_3 OR R_4
(All of them touch $h \Rightarrow h$ is ϵ -close to C)

$$\Pr[h \text{ is not } \epsilon\text{-close to } c] \leq$$

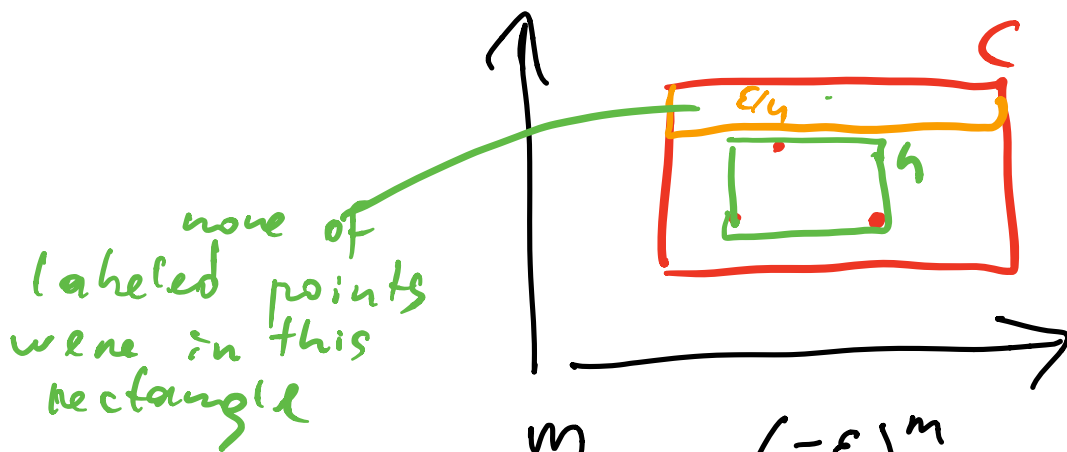
$$\Pr[h \text{ misses } R_1] +$$

$$\Pr[h \text{ misses } R_2] +$$

$$\Pr[h \text{ misses } R_3] +$$

$$\Pr[h \text{ misses } R_4]$$

$$= 4 \cdot \Pr[h \text{ misses rectangle of area } \epsilon/4]$$



$$= 4 \cdot (1 - \epsilon/4)^m \leq 4 \cdot \left(e^{-\frac{\epsilon}{4}}\right)^m = 4e^{-m\epsilon/4}$$

Taylor series $\forall x: \underline{1+x} \leq e^x$

$$\delta = 4e^{-m\epsilon/4}$$

$$m = \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right) \Rightarrow$$

h will be ϵ -close to c w.p. $1-\delta$

PAC-learn C

$$m = \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right) = \underline{\text{poly}\left(\frac{1}{\epsilon}, \frac{1}{\delta}\right)}$$

$$\epsilon = 0.01 \quad \delta = 0.02 \Rightarrow$$

some concrete m of labeled points suffice to learn rectangle