

MATH/COMP 562

Probability Review & Notation

PAC Learning

Ref

Mohri Ch 2

Shalev-Schwartz

2.3.1 & Ch 3

2022



Notation: Prob. vs Measure theory

Measure Theory
Measur. Space
 (X, \mathcal{M}, μ)

Measurable Set

Measurable function
 $f(x)$

Integral of f
 $\int_X f(x) d\mu(x)$

probability measure

Prob
Sample space
 (Ω, \mathcal{B}, P)

= Event =

Random Variable
 X

Expectation / Expected Value /
Mean of X
 $\mathbb{E}[X]$

Distribution $P(x)$ or $\mathcal{D}(x)$

ML
 $\Omega = X$ no \mathcal{B} mentioned
 $P \Rightarrow p(x)$ or $\mathcal{D}(x)$.

never. but x_i sampled iid.
from $P(x)$

$S^m = \{x_1, \dots, x_m\}$
dataset

random variable:
 X or f

Typical distribution
 $p(x)$ generates data

$\hat{p}(x) = \frac{1}{m} \sum x_i$
"empirical distribution"

Note

As $m \rightarrow \infty$
 $\hat{p}_m \rightarrow p$

in some sense.

not a function

yes as distribution.

discrete r.variable X

$$\Omega = \{\alpha_1, \alpha_2, \dots, \alpha_6\}$$

$$p(\alpha_i) = \frac{1}{6}$$

$$X(\alpha_i) = i$$

roll dice.

$$E[X] = \frac{1+2+\dots+6}{6}$$

$$= 2.1$$

$$\text{Var}(X) = \dots$$

$$S^M = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6\}$$

$$p(i) = \begin{cases} \frac{2}{12} & i=1, 2, 3, 6 \\ \frac{1}{12} & i=4 \\ \frac{3}{12} & i=5 \end{cases}$$

continuous r.v.

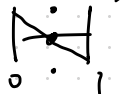
$$\Omega = [0, 1]$$

$$\alpha \in \Omega$$

Borel sets, generated by intervals (a, b)

$$p(S=(a,b)) = b-a \quad \text{uniform}$$

$$X(\alpha) = \frac{1}{2} - \alpha$$



$$\mu = E[X] = \frac{1}{2}$$

$$\text{Var}(X) = E[(X-\mu)^2]$$

$$= \int \alpha^2 dp(\alpha)$$

$$= \frac{\alpha^3}{3} \Big|_0^1 = \frac{1}{3}$$

Convergence in distribution

Discrete $\hat{p}_m(x) = \frac{1}{600} \begin{cases} 170 & x=1 \\ 96 & \\ \vdots & \\ 131 & x=6 \end{cases} \quad m=600$

Continuous $X = [0, 1]$

$p(x) = dx$ usual measure

$p([a, b]) = b - a$

$(a, b) \subset [0, 1]$

$$\hat{p}_m(x) = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}$$

Given (a, b)

$$\hat{p}_m(a, b) \rightarrow p(a, b), \quad \text{for all } (a, b) \subset [0, 1]$$

ML Notation

$(\mathcal{X}$ domain
 X r.v.)

$\mathcal{X} = [0,1]^d$ vector images.

$S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ dataset of cat images & labels

$y_i = \begin{cases} +1 & \text{if cat: pict. cute} \\ 0 & \text{o.w.} \end{cases} = y(x_i)$

x_i in database

$y(x) =$ true label of cuteness.

$\hat{p}(x) = \frac{1}{m} \sum \delta_{x_i}$ empirical distribution.

Trained Model / function / hypothesis

f or $h = A(S^m)$

\uparrow
model

\uparrow
hypothesis

\nwarrow
training alg

$p(x) =$ prob density on $[0,1]^d$ of cat images

$f(x) \in \{0,1\}$

random variable.
(because depends on S^m)

Expected loss

$$L(f) = \mathbb{E}(l_{0-1}(f, y)) = \int_{\mathcal{X}} l_{0-1}(f(x), y(x)) dp(x)$$

Mohri Ch2

X domain
 $Y = \{0, 1\}$

$Y : X \rightarrow Y$
target unknown

$h : X \rightarrow Y$
 $h \in \mathcal{H} = \{h \mid h : X \rightarrow Y\}$
hypothesis class

$$l_{0-1}(y, y') = \begin{cases} 1 & y \neq y' \\ 0 & y = y' \end{cases}$$

Generalization error $R(h) = L_{D, C}(h)$
 $= \mathbb{E}_{x \sim D} [l_{0-1}(h(x), c(x))] = \mathbb{P}_{x \sim D} [h(x) \neq c(x)]$

or Risk or Test error

Empirical (training error)

$$S = S_m = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

$$y_i = y(x_i)$$

$$\hat{R}_S(h) = \hat{L}_S(h) = \frac{1}{m} \sum_{i=1}^m l_{0-1}(h(x_i), y(x_i))$$

$$\left(= \mathbb{E}_{x \sim \hat{D}} [l_{0-1}(h(x), y(x))] \right)$$

Defn: probably approximately correct (PAC) how many samples

a concept class \mathcal{C} , is PAC learnable if for every $c \in \mathcal{C}$ there exist a learning algorithm such that, given any $\epsilon > 0$, $\delta > 0$ error tolerance there exists $m = m(\epsilon, \delta)$ s.t.

if you have m samples (generated i.i.d. from \mathcal{D})

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\underbrace{L_{\mathcal{D}}(h_S)}_{\text{gen loss}} \leq \underbrace{\epsilon}_{\text{error}} \right] \geq \underbrace{1 - \delta}_{\text{prob tol}}$$

over different S_m

EG coin toss

$$\epsilon = 0.02$$

$$\delta = 0.05$$

want

$$\left| \frac{\# \text{heads}}{n} - \text{prob}(\text{heads}) \right| < \epsilon$$

with prob $\geq 1 - \delta$

Learn / measure
by averaging

E.G. election
survey

52% vote for X

$$\pm 3\%$$

$$\epsilon = 0.03$$

$$19/20$$

$$\delta = 0.05$$

MATH Background

MATH 1

Union Bound

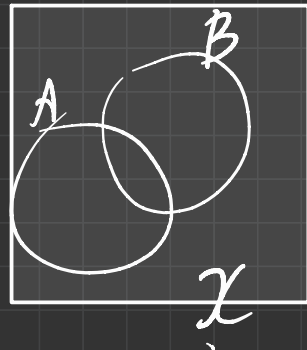
Given distribution \mathcal{D} on \mathcal{X}

$$A \subset \mathcal{X} \quad B \subset \mathcal{X}$$

pl.)

$$D(A \cup B) \leq D(A) + D(B)$$

$$D\left(\bigcup_{i=1}^m A_i\right) \leq \sum_{i=1}^m D(A_i)$$



MATH 2

proof

$$\textcircled{+} (1-\varepsilon)^m \leq \exp(-m\varepsilon) \quad \varepsilon \in [0, 1)$$

$$1-x \leq e^{-x} \quad x \in (0, 1)$$

$$(1-x)^m \leq e^{-mx} \quad x \in (0, 1)$$

$$\exists x \Rightarrow (1-\varepsilon)^m \leq e^{-m\varepsilon}$$

$$\left(\begin{array}{l} e^{-x} = 1 - x + \frac{x^2}{2} \dots \\ \underline{1-x \leq e^{-x}} \end{array} \right)$$

IDEA of proof.

Simplex problem

$C = \{ \text{intervals } [0, t] \subset [0, 1] \}$

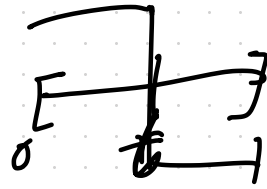
PAC learnable:

$$h(x) = \begin{cases} 1 & x \leq \hat{t} \\ -1 & \text{o.w.} \end{cases}$$

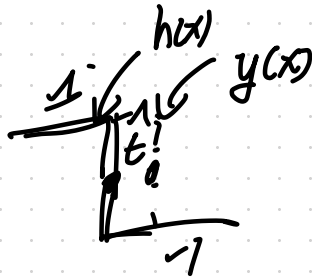
$$\hat{t} = \max(x_i | y_i = +1)$$

$$R(h) = L(h)$$

$$= \mathbb{E} [L_{0.7}(h_{\text{opt}}, y(x))] = t - \hat{t}$$

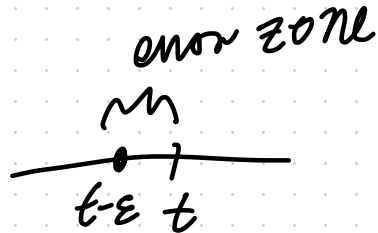


$$y(x) = \begin{cases} 1 & x \leq t \\ -1 & \text{o.w.} \end{cases}$$



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

Fix $\epsilon > 0$. (Assume $t > \epsilon$ on. fin.)



For a single sample x_i
prob miss error zone
 $= 1 - \epsilon$

For m sample, prob all miss $= \mathbb{P}(R(h_S) > \epsilon)$
 $= (1 - \epsilon)^m$
 $\leq \exp(-m\epsilon)$ by Math 2

Thus want

$$\mathbb{P}() \leq \delta \Rightarrow e^{-m\epsilon} \leq \delta$$
$$\Rightarrow -m\epsilon \leq \log \delta$$

$$\Rightarrow m \geq \frac{1}{\epsilon} \log \frac{1}{\delta} \quad \checkmark$$

Equivalently Gen. bound:

$$\delta = \exp(-m\epsilon)$$

SO with prob $\geq 1 - \delta$

$$R(h_S) \leq \frac{1}{m} \log \frac{1}{\delta}$$



HW

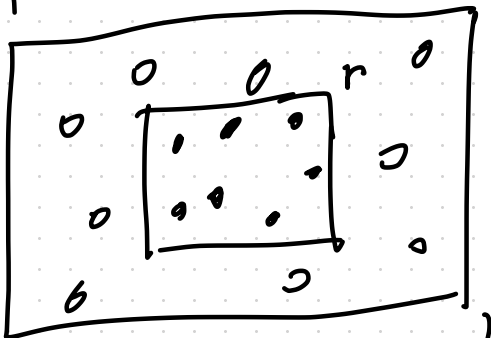
EX 2.4

Learnin

Axis-Aligned Rectangles

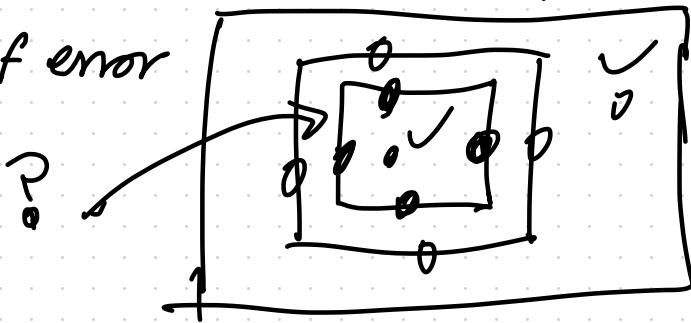
$$C = \{ \text{all axis aligned rect} \}$$

$$y(x) = \begin{cases} 1 & x \in C \\ 0 & x \notin C \end{cases}$$



Alg Define $h =$ smallest rectangle contains positive examples.

Window of error



prove (Math)
 frac ~~size~~ of uncorrect
 error $\rightarrow 0$
 as $m \rightarrow \infty$
 with high prob.