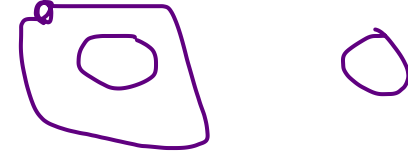hnt3lm eh mn el slides de?
1. hal e7na n2dr nt3lm general function a2dr beha a3ml correct prediciton le data ana mshufthash bona2n 3la el data set elly etmrnt 3leha wala laa, lw ah b3ml keda ezay.
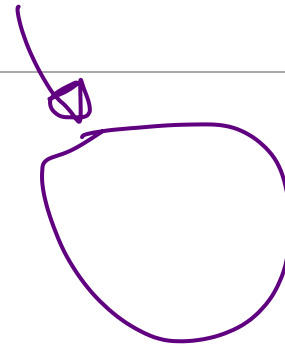2. y3ny a hoeffdeings inequallity, w ezay btshtghl 3la el single hypothsis, and multiple hypotheses.

always remember that our goal is to know -> is learning from data feasible?

# Lecture 2: The Learning Theory

DINA ELREEDY

CMP402B-SPRING 2022

# Learning Puzzle

As we sajd in the last lecture, we have a target function which is completely unknown, and we are trying to learn from data, in other word, we learn from experience, not by applying analysis.
so the obvious question here is, how can a limited data set allow me to correctly generate g function which allows me to accuratly estimate f which is the unknown target function.

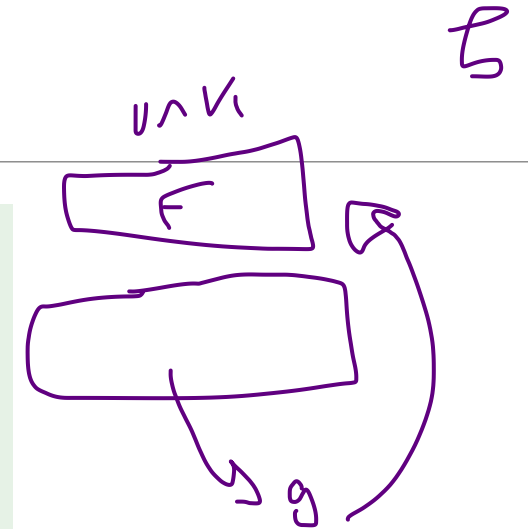Actually we can apply many analytical approaches:
1. if symmetric, then we may say it is 1.
2. if the top left corner is black, then it should be -1
3. etc...
however, we will see that there are no enough data to precisly judge the test case.

$f = -1$

$f = +1$

$f = ?$

We know what we have already seen, but that is not learning, that is memorizing -> basmga

Does the data set D tell us anything outside of D that we didn't know before? if yes -> learning
if no -> them learning is not feasible.

And we will discuss in the following slides that this is not a special case, usually small data set can not give me a general rule to be able to judge correctly.

$u \wedge v_1$

$\xi$

F

g

# Is learning feasible?-Example

X={0,1}^13

Y={0,1}

$2^3 \left\{ \begin{array}{ccc|c} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ & & \vdots & 0 \\ & & & \vdots \\ 1 & 1 & 1 & 0 \end{array} \right\}$

in how many ways we can have distinct y vector?
since repeation is allowed,
since we have 8 samples,
since each sample have 2 options.
therefore # of ways = 2 ^ 8.
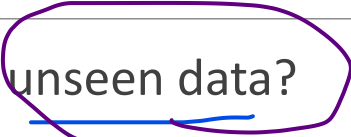
Hypothesis set H:All Boolean functions X→Y

|H|=?

|H|=$2^8$=256

f: Target function

SnP

# Is learning feasible?-Example

this is not a complete table.

| x | y | g | f1 | f2 | f3 | f4 | f5 | f6 | f7 | f8 |
|---|---|---|----|----|----|----|----|----|----|----|
| 0 0 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 0 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 0 1 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 0 1 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 100 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 101 | ? | ? | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| 110 | ? | ? | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| 111 | ? | ? | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |

Observed Training data

Testing data

Same values

so does this mean that learning from the data is doomed? if yes, then why we have invented this course y3ny? :D

Nothing to tell us which f is the true function??

# Is learning feasible?

- Can we generalize to unseen data?

- Yes, it is feasible to learn the unknown function f in a probabilistic way.

# Analogy $N$

30./.

[3]'

•Assume we have a bin of infinitely many red and green marbles.

•Assume we pick a random sample of N marbles out of the bin.   with replacement.

$N = 10$

**B I N**

**S A M P L E**

[3]
[10]

●●●●●●●●●●
$V =$ **fraction of red marbles**

• $\mu = V$??

•Can $V$ say anything about $\mu$?

$X$
$\mu =$ **probability of red marbles**

kol ma bnkbr el sample size, kol ma el v btb2a a2rb lel u
el mfrod en el v tb2a representative w tb2a t2reban keda bt3br 3n el genral mean.

example el 1k dollars lw gbt 10 mrat head
forst enk t5sr hya
P(1 = T) * P ( 2 = T) .....

$\Sigma$

# Hoeffdeing's Inequality

$\mu$: Frequency of red marbles in the bin

$\nu$: Fraction of red marbles in a random sample of size N.

3auzak tukhrug b 3 ma3lomat mn hena:
1. kol ma btkbr el 7gm bta3 el sample, kol ma e7tmalyt enk tkhrug bra el boundry bt2l exponentially.

2. kol ma btsm7 be nsbt error akbur, kol ma brdu el probability bt2l exponentially.

3. el equation de esmaha Hoeffdeing's inequality, w hya hadafha enha t2olak en kol lama bnakhud samples akbur w ngeb el mean bta3hom ($\nu$), nesbt en el $\nu$ ykon b3ed 3n el u bnsba akbur mn e hya nesba olayela.

$$P(|\nu\text{-}\mu| > \varepsilon) \leq 2e^{-2\varepsilon^2 N} \quad \text{For any } \varepsilon > 0$$

e7tmaleyt enk tkhrog bara el bound.

factor in known values, fkeda helw.

- What happens as sample size N increases?
  - $\nu$ has a better approximation of $\mu$.

- Trade-off between N, $\varepsilon$, and the bound.

tb  eh el hadaf y3ny? el hadaf eny a2dr a3ml prediction lel data points elly baraa el D data set -> predict the values of the test sets.

# Analogy to the learning problem

- Bin: Input space X

- Sample:(Observed) Training set D

- Each marble is a point x that belongs to the input space X

- Red marbles: Misclassified points by hypothesis h $h(x) \neq f(x)$

- Green marbles: Correctly classified points by hypothesis h $(h(x)=f(x))$

- $\nu$: Fraction of misclassified data points by hypothesis h in traning data. **(Ein(h))**

- $\nu$ is the in-sample error of the hypothesis h **(Ein(h))**

- $\mu$: Fraction of misclassified data points by hypothesis h in the input space. **(Eout(h))**

- $\mu$ is the out-of-sample error of the hypothesis h **(Eout(h))**

# In-sample versus Out-of sample error

$Ein(h) = fraction\ of\ D\ where\ f\ and\ h\ disagree$

$= \frac{1}{N}\sum_{n=1}^{N}[\![h(xn) \neq f(xn)]\!]$ ✓

Expected mean.

$Eout(h) = P[h(xn) \neq f(xn)]$

True mean

# Hoeffdeing's Inequality

- Hoeffdeing's Inequality to the learning problem:

$$P(|Ein(h)\text{-}Eout(h)|>\varepsilon)\leq 2e^{-2\varepsilon^2 N} \qquad \text{For any } \varepsilon>0$$

The probability that $h$'s performance on points not in the training set deviates from $h$'s performance on points in the training set is low.

# Learning Setup-updated

- Important Note:

A fundamental assumption made

in this analysis is that the sample

is taken randomly and training and

testing data have the same distribution.



UNKNOWN TARGET FUNCTION
$f: X \to Y$

PROBABILITY DISTRIBUTION

$P$ on $X$

TRAINING EXAMPLES
$(x_1, y_1), \dots, (x_N, y_N)$

$x_1, \dots, x_N$

LEARNING ALGORITHM
$A$

FINAL HYPOTHESIS
$g \approx f$

HYPOTHESIS SET
$H$

Source: Learning from data book

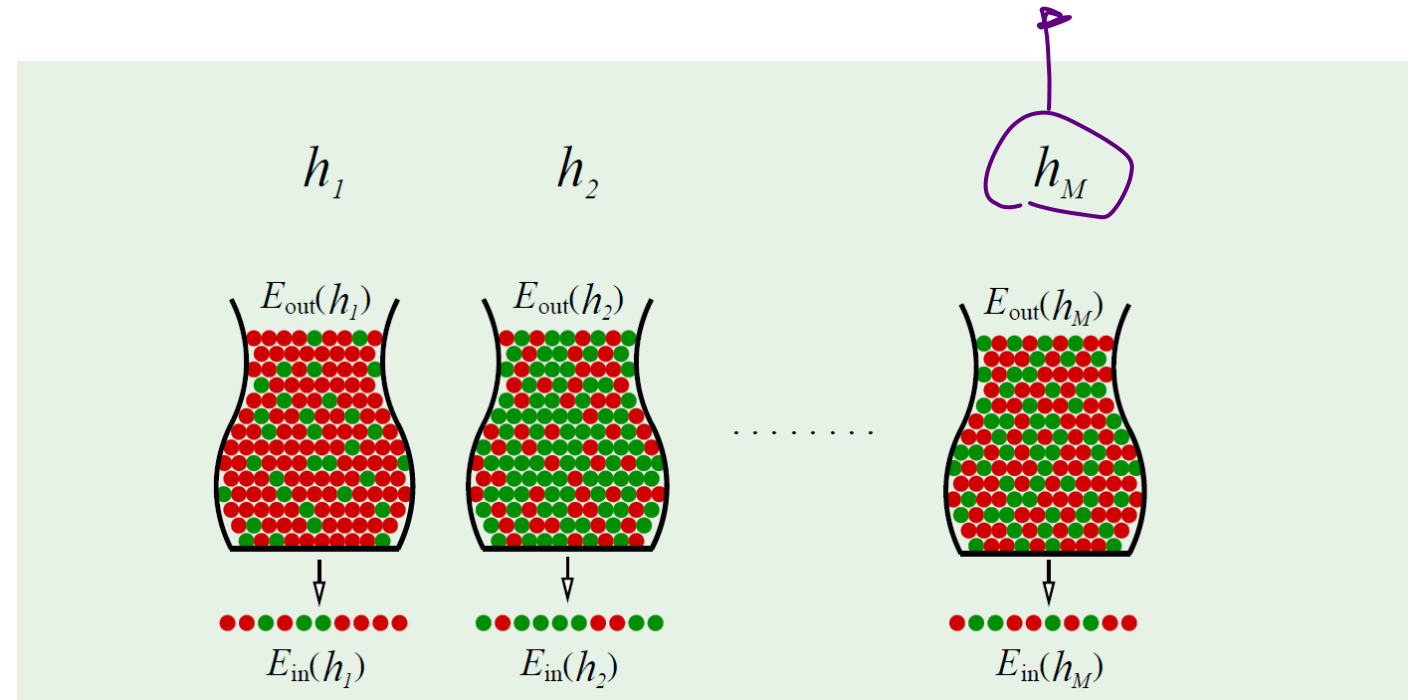usually we don't have a control over v, however, usually we have a set of Hypotheses H, from which we select one hypothesis, which has the least error.

it sounds harder than what it is actually, fkr feha zy hagat el neuarl network el khdnaha, el set of Hypotheses de hya kol el possible values bta3t el weights, w msh e7na bno3od net3lm l7d ma nwsl le weight vector ydena a2l error? hwa da el final hyposis, kol el ablo homa another hypotheses gowa el set of H.

kol el fat, bfterad lw 3ndy hypothesis wahda.

# Multiple Hypotheses ?

- In learning, we have multiple hypotheses to pick from.

- Finite Hypothesis set H={$h_1$,$h_2$,…$h_M$}

- H is analogous to multiple bins.

- The final hypothesis g is selected

by the learning algorithm.

# Hoeffdeing's Inequality for final hypothesis g

$$\mathbb{P}\big[\,|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\,\big] \;\leq\; \mathbb{P}\big[\quad |E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| > \epsilon$$

$$\textbf{or } |E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| > \epsilon$$

$$\ldots$$

$$\textbf{or } |E_{\text{in}}(h_M) - E_{\text{out}}(h_M)| > \epsilon\,\big]$$

$$\leq \sum_{m=1}^{M} \mathbb{P}\big[|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon\big]$$

very upper bound, because we have neglected all the intersections, and assumed that all of the hypothesis, are independent, and that is not true, however, this equation still holds, but it is a very upper bound.

# Hoeffdeing's Inequality for final hypothesis g

$$\mathbb{P}[\ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\ ] \leq \sum_{m=1}^{M} \mathbb{P}\left[|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon\right]$$

$$\leq \sum_{m=1}^{M} 2e^{-2\epsilon^2 N}$$

constant.

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

As M (the hypothesis set size) increases, the bound gets looser.

kol ma 3dd el hypotheses zad, kol ma ana kbrt el bound, we de haga msh kwysa, l2n ana hdfy awsl le bound soghyr.

# Feasibility of Learning

Two main questions:

1. **Can we make sure that Eout(g) is close enough to Ein(g)?**
   Hoeffdeing's Inequality.

2. **Can we make Ein(g) small?**

   Learning models