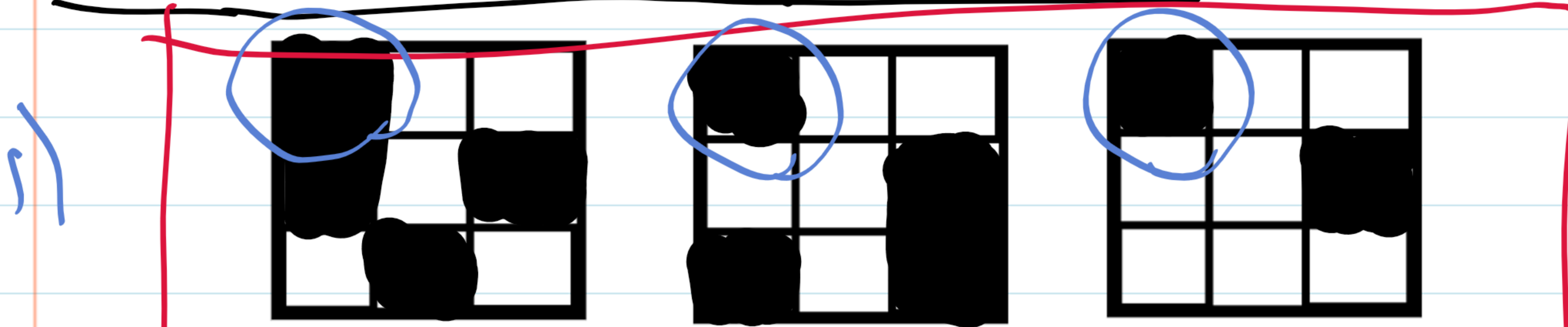
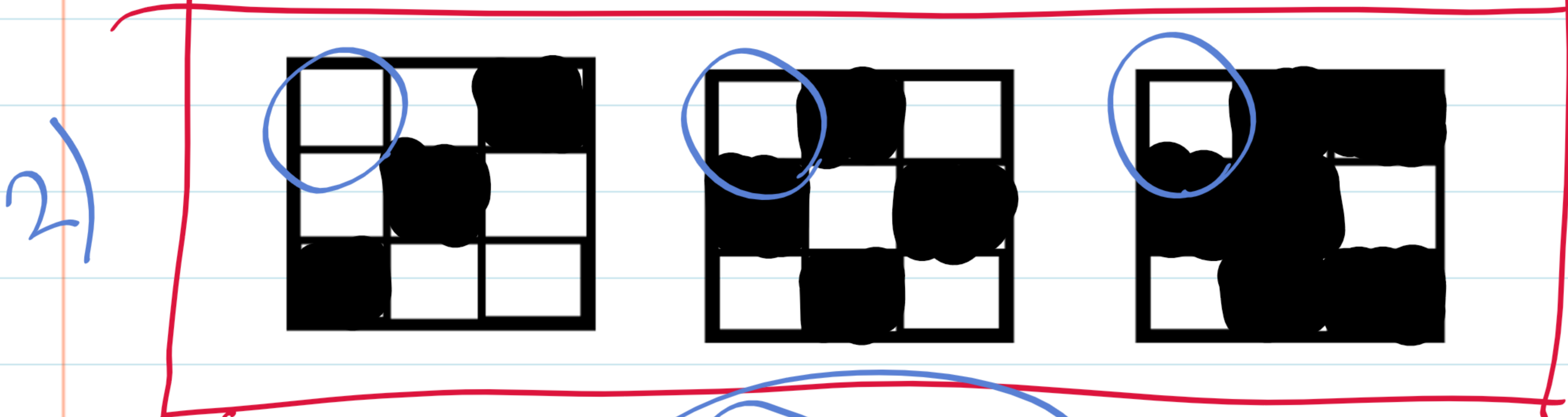


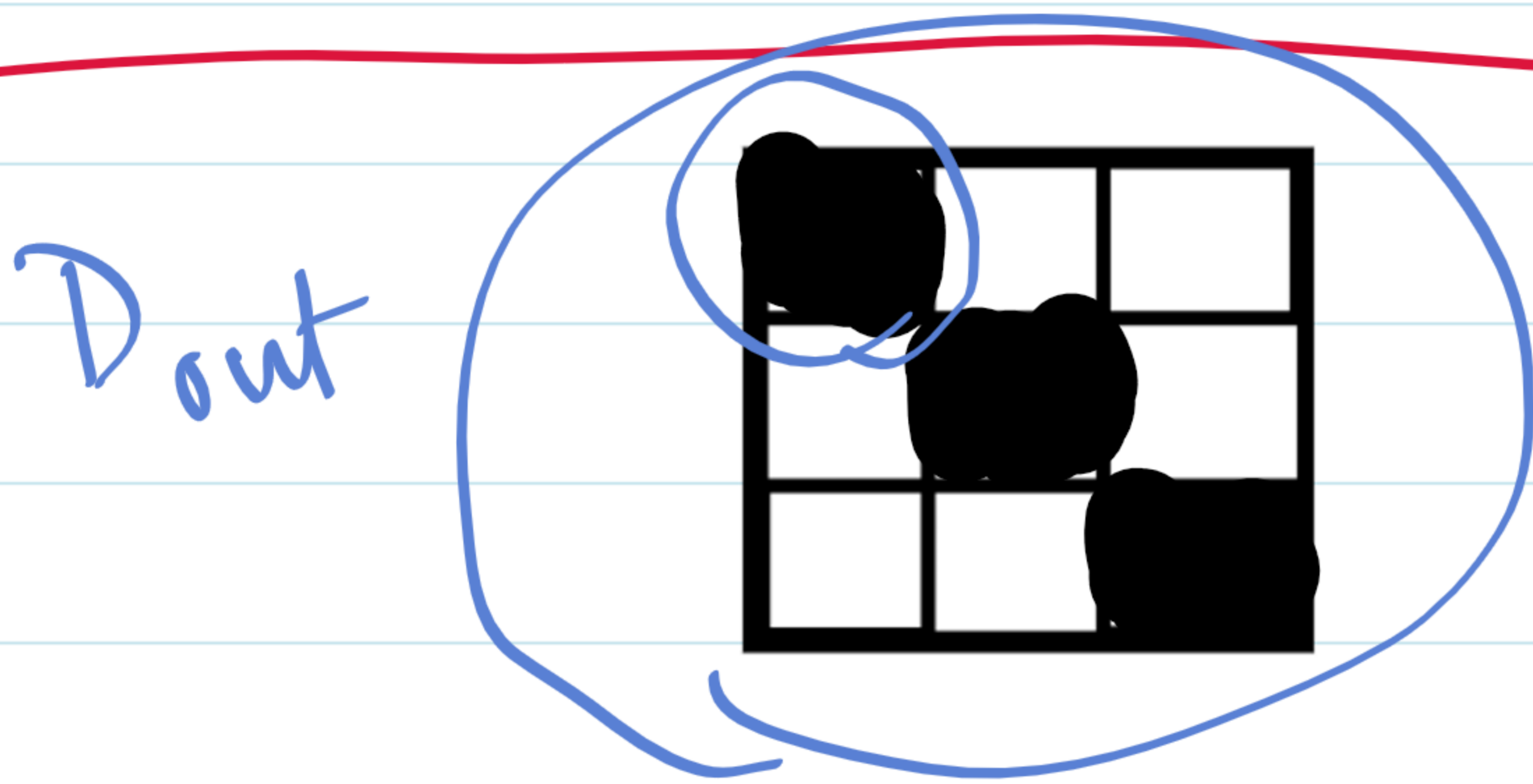
A visual learning problem :



$$f(x) = -1$$



$$f(x) = +1$$



$$f = ??$$

9-bit vector (3x3 black & white array)

Possibility 1 : $f = +1$, due to symmetry

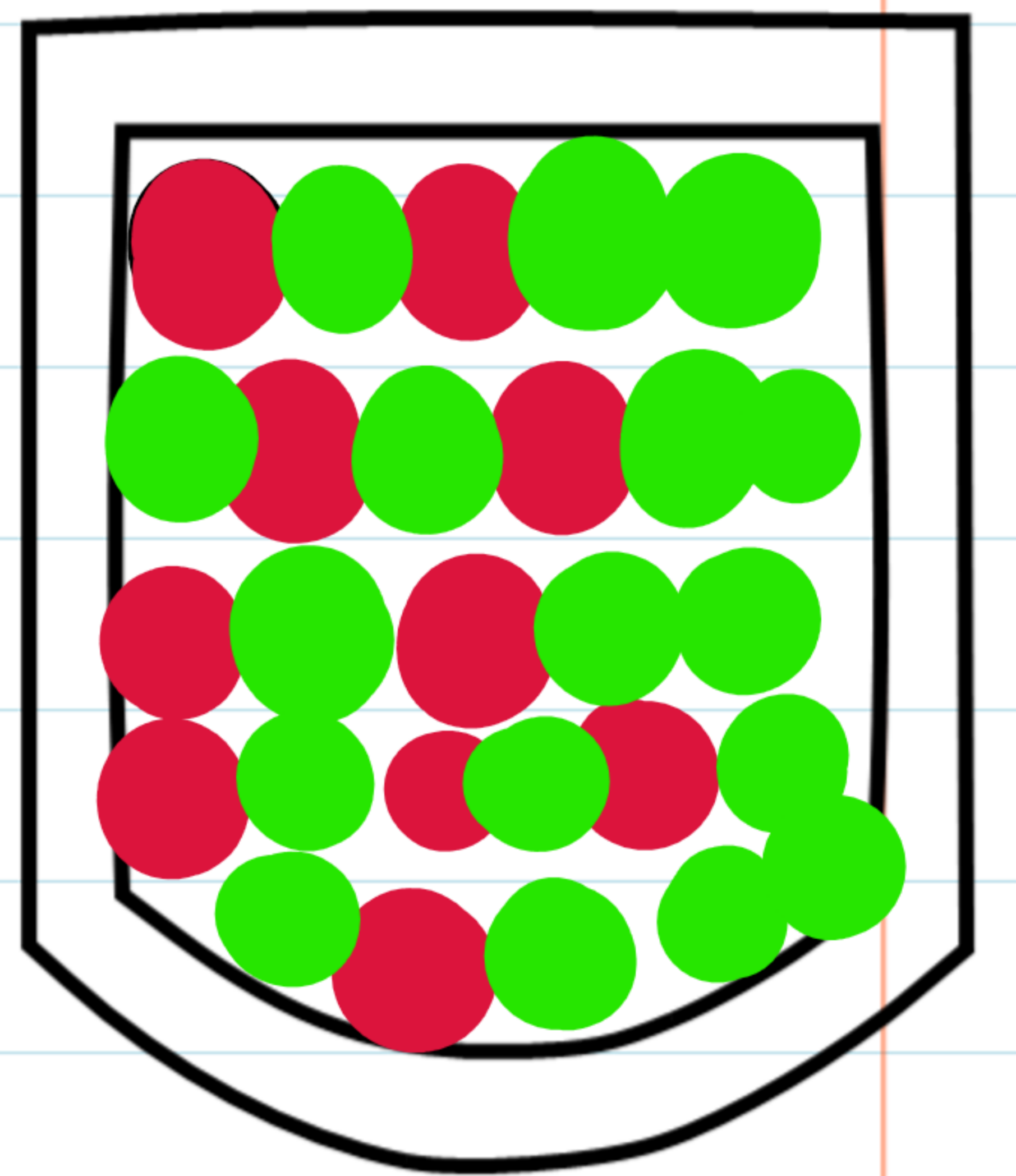
Possibility 2 : $f = -1$, top left corner

Consider a bin with red and green marbles.

$$P[\text{Picking a red marble}] = \mu$$

$$P[\text{Picking a green marble}] = 1 - \mu$$

Bin



→ Value of μ is unknown. Then we pick N marbles independently.

→ Sample = { }

→ The fraction of red marbles in sample = ν .

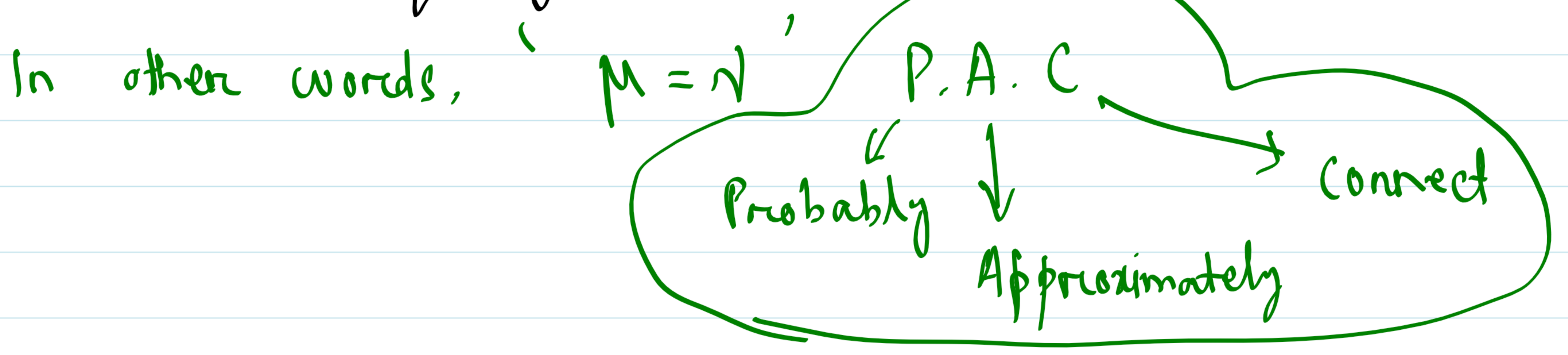
Q: Does ν say something about μ ?

In a big sample (large N), \bar{v} is probably close to μ (with in ϵ)

Formally,
$$P[|\bar{v} - \mu| > \epsilon] \leq 2e^{-\frac{2\epsilon^2 N}{\sigma^2}}$$

Hoeffding inequality

Sample frequency
Bin frequency



Trade off: N , ϵ and the bound.

Connection to learning:

Bin: The unknown is a number μ .

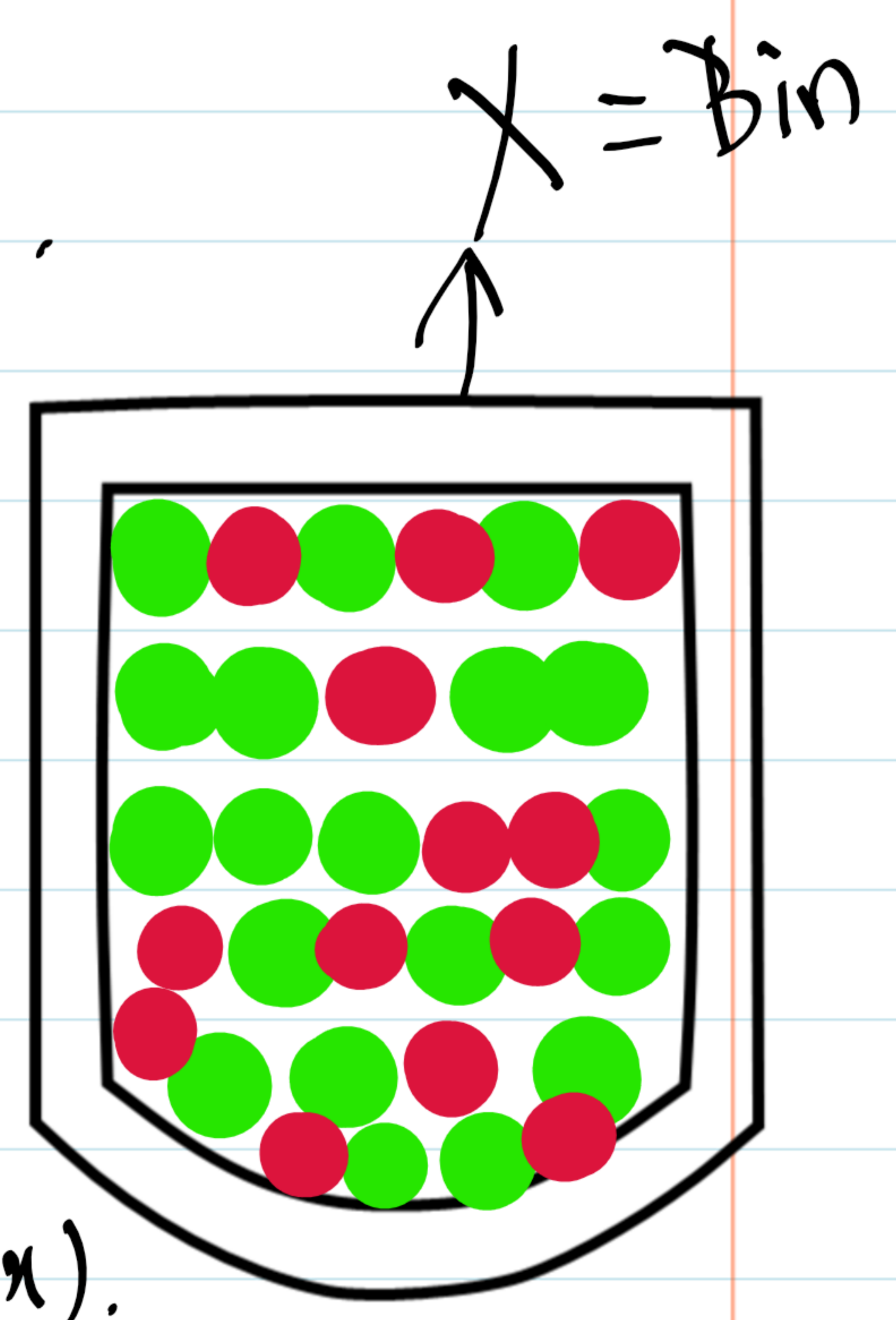
Learning: The unknown is an entire function

$$f: X \rightarrow Y$$

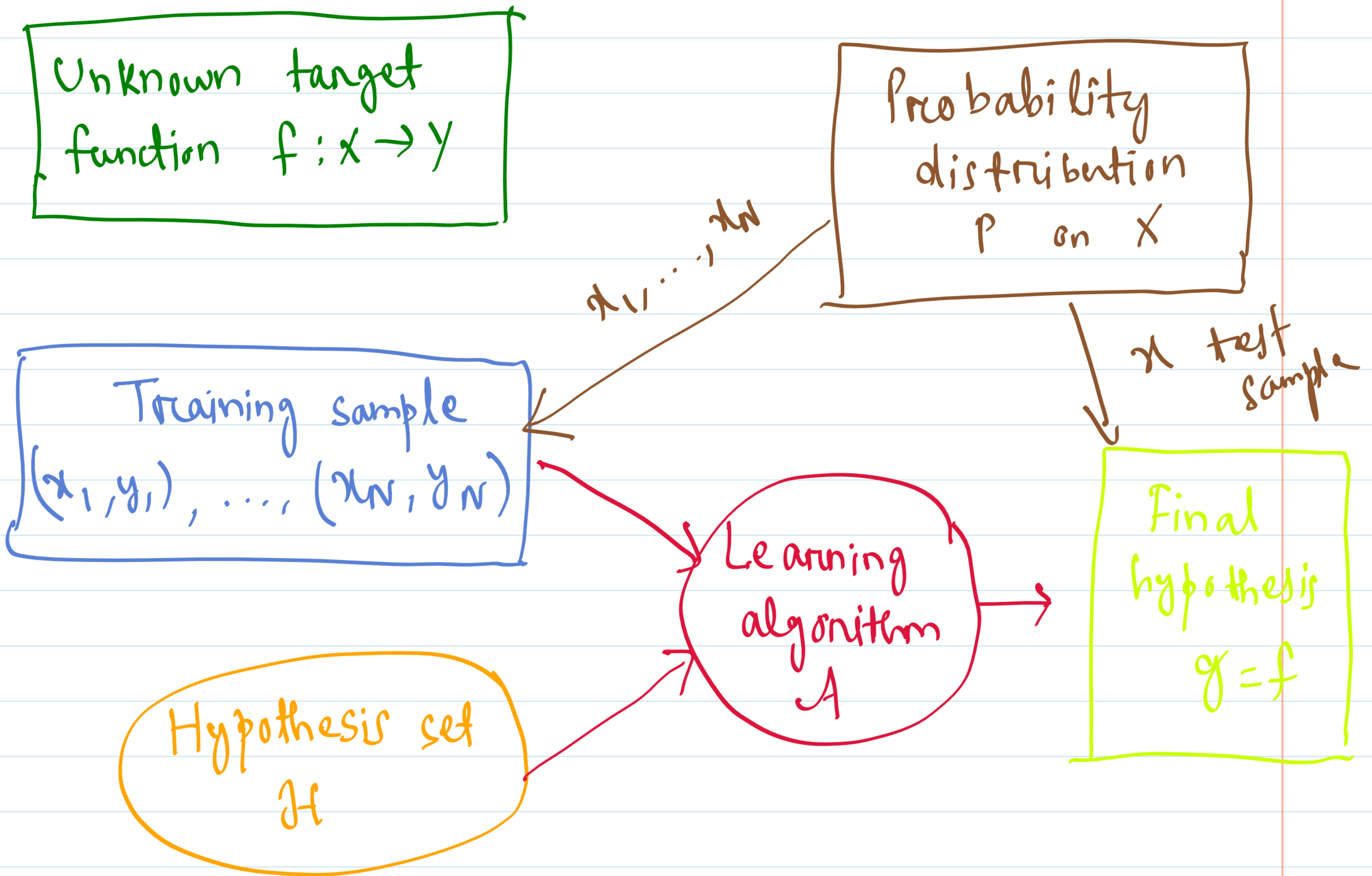
Every marble '0' is a point $x \in X$.

Take any single hypothesis $h \in \mathcal{H}$ & compare it to f at each point $x \in X$.

- ; hypothesis got it right $h(x) = f(x)$
- ; hypothesis got it wrong $h(x) \neq f(x)$.



- P can be anything
- I don't need to know what P is.



→ Both μ and ν depend on h .
 ν is 'in-sample' $E_{in}(h)$ & μ is 'out-sample' $E_{out}(h)$.

The Hoeffding inequality:

$$P\left[|E_{in}(h) - E_{out}(h)| > \epsilon\right] \leq 2e^{-2\epsilon^2 N}$$

Probability inequalities:

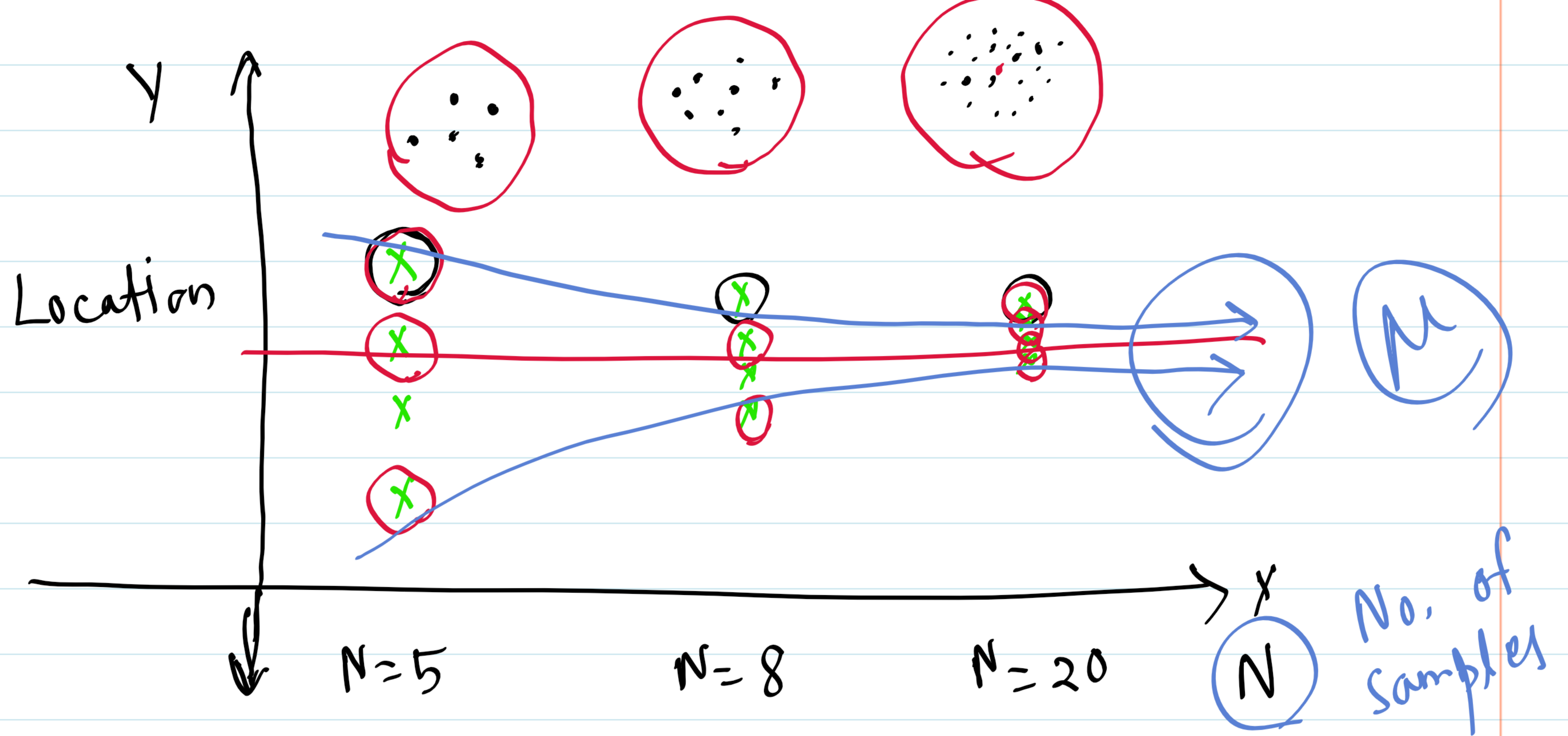
→ Take x_1, \dots, x_N in 1-D, which are i.i.d r.v.,

$$E[x_1] = \dots = E[x_N] = \mu \quad (\text{population mean})$$

→ Empirical average $\nu = \frac{1}{N} \sum_{n=1}^N x_n$.

→ How close ν is to μ .

→



→ Markov inequality: $P[X \geq \epsilon] \leq \frac{E[X]}{\epsilon}$.

→ Chebyshev inequality: Let X_1, \dots, X_N be i.i.d. with $E[X_n] = \mu$, $\text{Var}[X_n] = \sigma^2$, $\bar{X} = \frac{1}{N} \sum_{n=1}^N X_n$,

Then $P[|\bar{X} - \mu| > \epsilon] \leq \frac{\sigma^2}{N\epsilon^2}$.

→ Weak law of large number (WLLN)

Let X_1, X_2, \dots, X_N be i.i.d. r.v. with common μ .

Let $M_N = \frac{1}{N} \sum_{n=1}^N X_n$. Then for any $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P[|M_n - \mu| > \epsilon] = 0$$

→ Strong law of large number (SLLN)

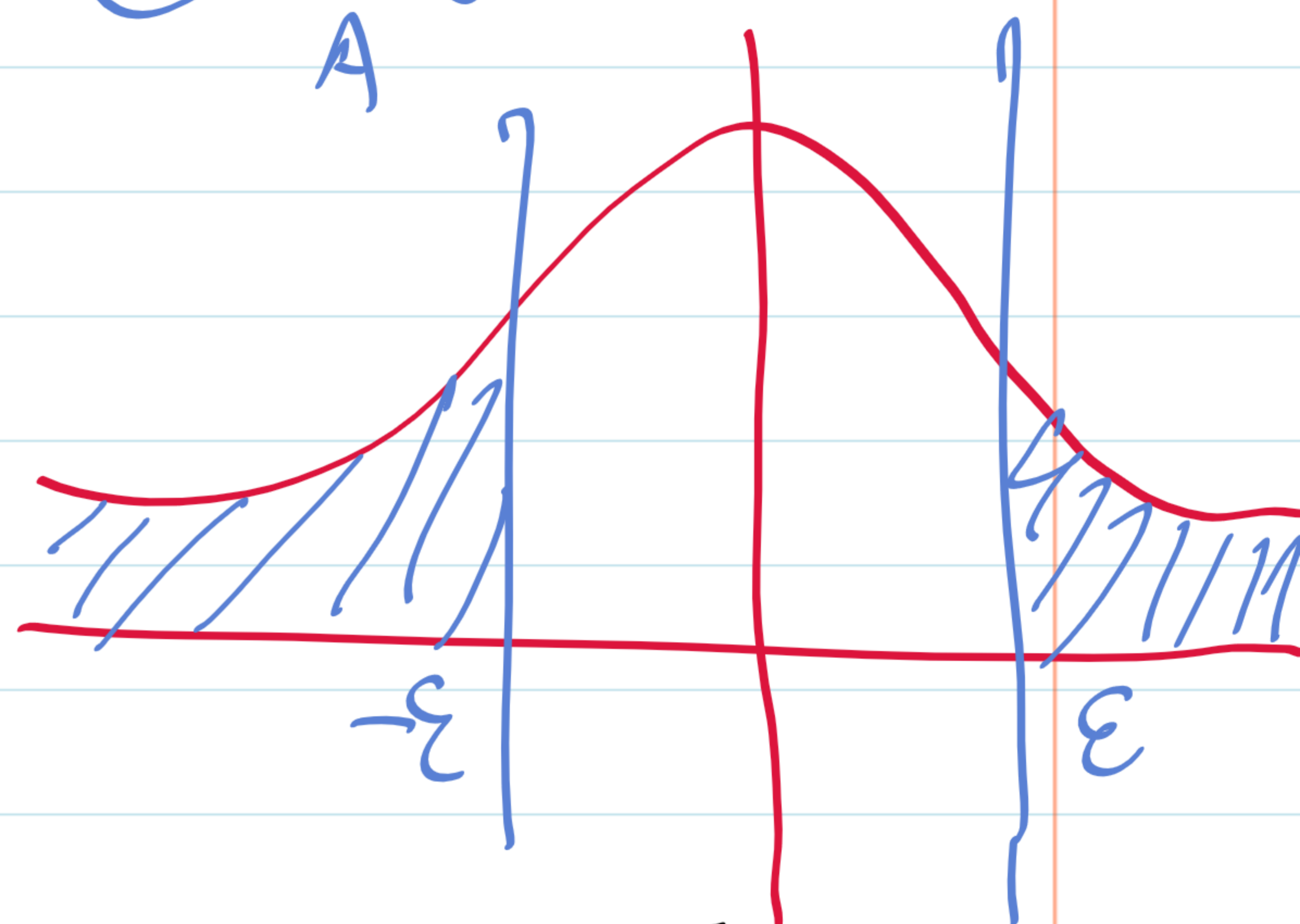
$$P\left[\lim_{n \rightarrow \infty} |M_n - \mu| > \epsilon\right] = 0$$

Let us revisit the bad event:

$$P[|v - \mu| \geq \varepsilon] = P[v - \mu \geq \varepsilon \text{ or } v - \mu \leq -\varepsilon]$$

$$\leq \underbrace{P[v - \mu \geq \varepsilon]}_A + \underbrace{P[v - \mu \leq -\varepsilon]}_A$$

$$\leq 2 \underbrace{A}_{\text{Chernoff}} \leq e^{-2\varepsilon^2 N}$$



The e-trick:

$$P[v - \mu \geq \varepsilon] = P\left[\frac{1}{N} \sum_{n=1}^N x_n - \mu \geq \varepsilon\right]$$

$$= P\left[\sum_{n=1}^N (x_n - \mu) \geq \varepsilon N\right]$$

(e-trick) $= P\left[e^{s \sum_{n=1}^N (x_n - \mu)} \geq e^{s \varepsilon N}\right]$

(Markov) $\leq \frac{E\left[e^{s \sum_{n=1}^N (x_n - \mu)}\right]}{e^{s \varepsilon N}}$

$$\leq \left(\frac{E\left[e^{s(x_n - \mu)}\right]}{e^{s\varepsilon}} \right)^N$$

If we let $z_n = x_n - \mu$. Then

$$E\left[e^{s(x_n - \mu)}\right] = M_{z_n}(s) = \text{MGF of } z_n \text{ (moment generating function)}$$

Hoeffding lemma :

If $a \leq x_n \leq b$, then $\mathbb{E} \left[e^{s(x_n - \mu)} \right] \leq e^{\frac{s^2(b-a)^2}{8}}$.

Let $a = 0, b = 1$.

$$\mathbb{P}[\bar{x} - \mu \geq \epsilon] \leq \left(\frac{\mathbb{E} \left[e^{s(x_n - \mu)} \right]}{e^{s\epsilon}} \right)^N$$

$$\leq e^{\frac{s^2 N}{8} - s\epsilon N} \quad \forall s > 0 \quad (\text{function of } s)$$

minimize : s^*

$$f(s) = \frac{s^2 N}{8} - s\epsilon N$$

$$\frac{d}{ds} \left(\frac{s^2 N}{8} - s\epsilon N \right) = \frac{N}{4} s - N\epsilon = 0$$

$$\Rightarrow \underline{s^* = 4\epsilon}$$

$$s_0 \left[\mathbb{P}(\bar{x} - \mu \geq \epsilon) \right] \leq e^{\frac{N}{8} s^{*2} - s^* \epsilon N}$$

$$= e^{\frac{N}{8} (16\epsilon^2) - 4\epsilon^2 N} = e^{-2\epsilon^2 N}$$

$$\mathbb{P}(|\bar{x} - \mu| \geq \epsilon) \leq 2e^{-2\epsilon^2 N}$$

Hoeffding inequality.

Compare between Hoeffding & Chebyshev :



$$\left. \begin{array}{l} \text{Chebyshev : } \mathbb{P}(|\bar{x} - \mu| \geq \epsilon) \leq \frac{\sigma^2}{N\epsilon^2} \\ \text{Hoeffding : } \mathbb{P}(|\bar{x} - \mu| \geq \epsilon) \leq 2e^{-2\epsilon^2 N} \end{array} \right\} \underline{\underline{\mathbb{P}(|\bar{x} - \mu| \geq \epsilon) \leq 2e^{-2\epsilon^2 N}}}$$

Equivalent to: For probability at least $1-\delta$, we have

$$\underbrace{\mu - \epsilon}_{\text{testing error (unknown)}} \leq \underbrace{\nu}_{\text{training error}} \leq \mu + \epsilon$$

→ Chebyshev: $\delta = \frac{\sigma^2}{N\epsilon^2} \Rightarrow \epsilon = \frac{\sigma}{\sqrt{SN}}$

Hoeffding: $\delta = 2e^{-2\epsilon^2 N} \Rightarrow \epsilon = \sqrt{\frac{1}{2N} \log \frac{2}{\delta}}$

Example: Chebyshev: For probability at least $1-\delta$,

we have $\mu - \frac{\sigma}{\sqrt{SN}} \leq \nu \leq \mu + \frac{\sigma}{\sqrt{SN}}$.

Hoeffding: For probability at least $1-\delta$, we have

$$\mu - \sqrt{\frac{1}{2N} \log \frac{2}{\delta}} \leq \nu \leq \mu + \sqrt{\frac{1}{2N} \log \frac{2}{\delta}}$$

You ask: I have data x_1, \dots, x_N . I want to estimate μ . How many data points N do I need?

Other person: How much δ can you tolerate?

You say: I only have limited no. of data points. Then how good my estimate is? (ϵ)

Other person: How many data points N do you have.

Example: Let $\delta = 0.01$ (1% error), $N = 10000$
 $\sigma = 1$ (n.v.). How much tolerance you
can afford?

$$\epsilon = \frac{\sigma}{\sqrt{\delta N}} = 0.1$$

(Chebyshev)

$$\epsilon = \sqrt{\frac{1}{2N} \log \frac{2}{\delta}} = 0.016$$

(Hoeffding)

↓
lot
tighter

Change your question: probability of error
 $\delta \neq 0.01$,

and $\epsilon = 0.01$, $\sigma = 1$, then

$$N \geq \frac{\sigma^2}{\epsilon^2 \delta} = 1,000,000 \quad \text{and}$$

(Chebyshev)

$$N \geq \frac{\log \frac{2}{\delta}}{2\epsilon^2} \approx \boxed{26,500}$$

(Hoeffding)

PAC

PAC learning framework

Recall the eqⁿ: $P \left[|E_{in}(h) - E_{out}(h)| > \epsilon \right] \leq 2e^{-2\epsilon^2 N}$.

One can bound $E_{out}(h)$ using $E_{in}(h)$.

$E_{out}(h)$ is something that you don't know

$E_{in}(h)$ which is something that you know.

→ RHS is independent of h and $p(x)$. (universal bound)

→ This works for any \mathcal{A} , \mathcal{H} , f & $p(x)$.

→ $\delta = 2e^{-2\epsilon^2 N}$ confidence: $1 - \delta$.

→ $\epsilon = \sqrt{\frac{1}{2N} \log \frac{2}{\delta}}$ accuracy: $1 - \epsilon$.

→ One can write: $P \left[|E_{in}(h) - E_{out}(h)| > \epsilon \right] \leq \delta$,

which is equivalent to

$$P \left[|E_{in}(h) - E_{out}(h)| \leq \epsilon \right] > 1 - \delta.$$

Probably Approximately Correct (PAC) framework,

Probably: Quantify error using probability

$$P \left[|E_{in}(h) - E_{out}(h)| \leq \epsilon \right] \geq 1 - \delta$$

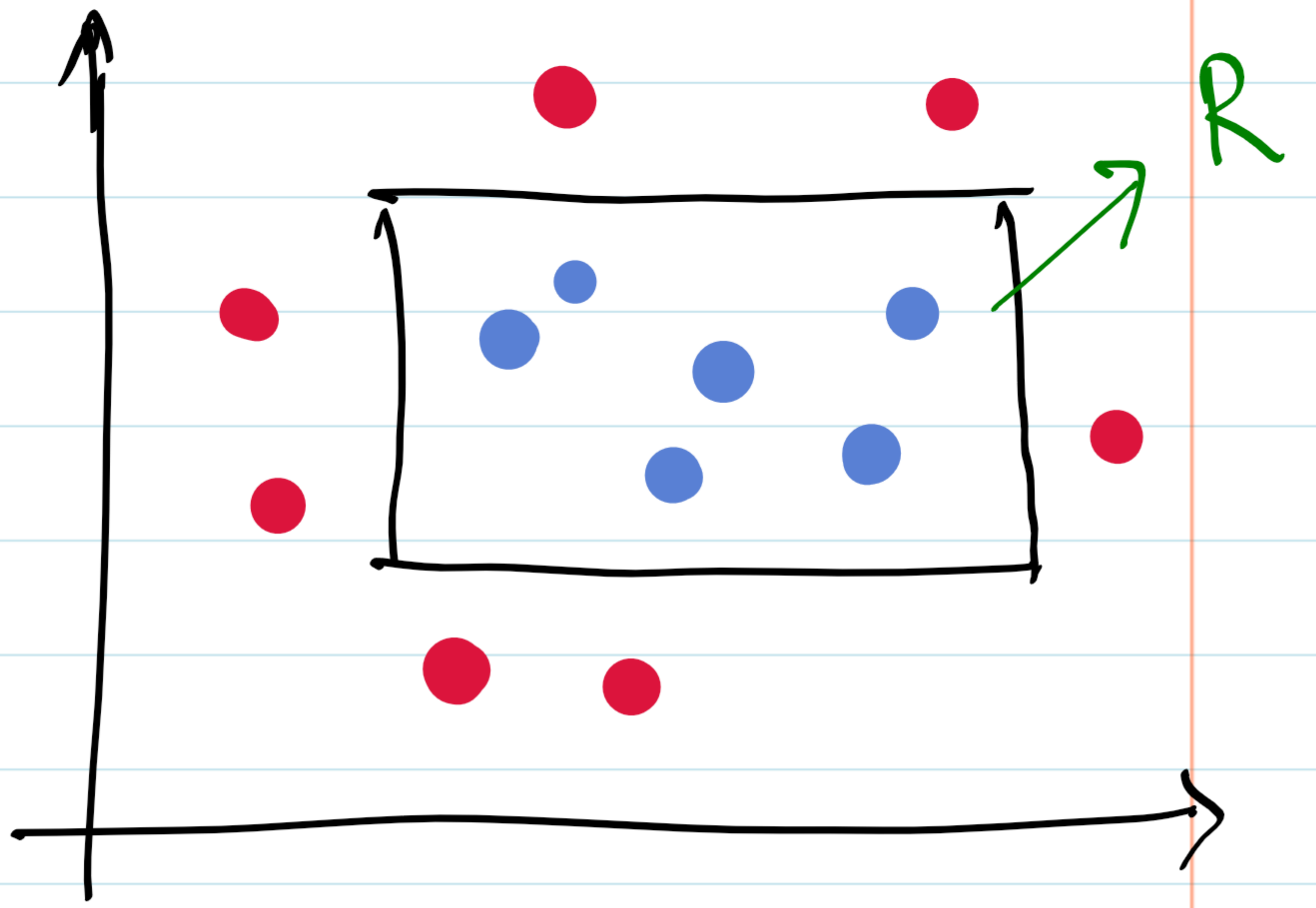
Approximately correct: In-sample error is an approximation of the out-sample error.

$$P \left[|E_{in}(h) - E_{out}(h)| \leq \epsilon \right] \geq 1 - \delta$$

PAC-learnable: If you can find an algorithm \mathcal{A} such that for any ϵ and $\delta > 0$, there exists an N which can make the above inequality holds, then we say that the target function is PAC-learnable.

Example: Rectangular classification:

- Consider a set of 2D data points,
- The target function R is a rectangle.



- Inside R : blue } Data is separable
- Outside R : red }

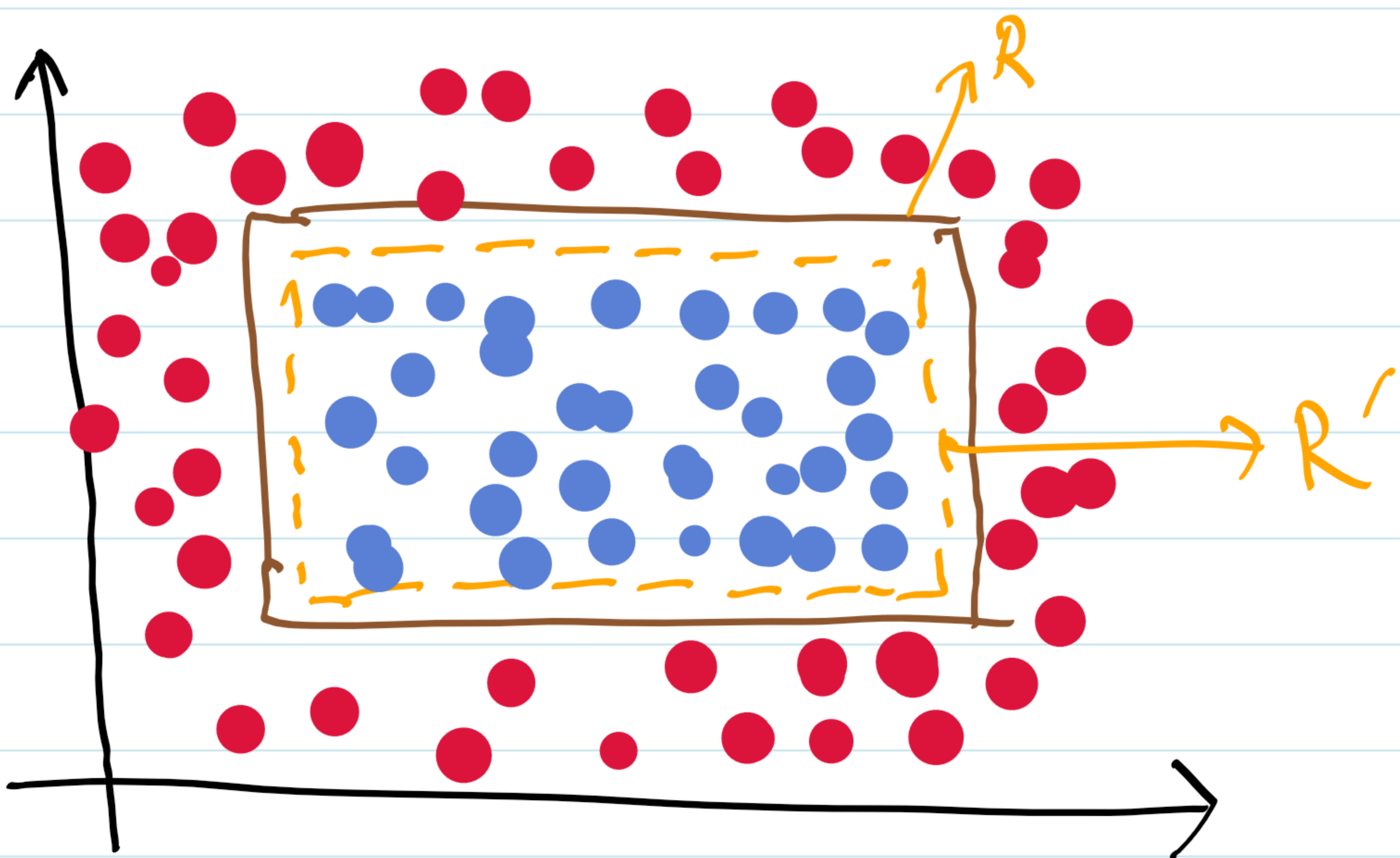
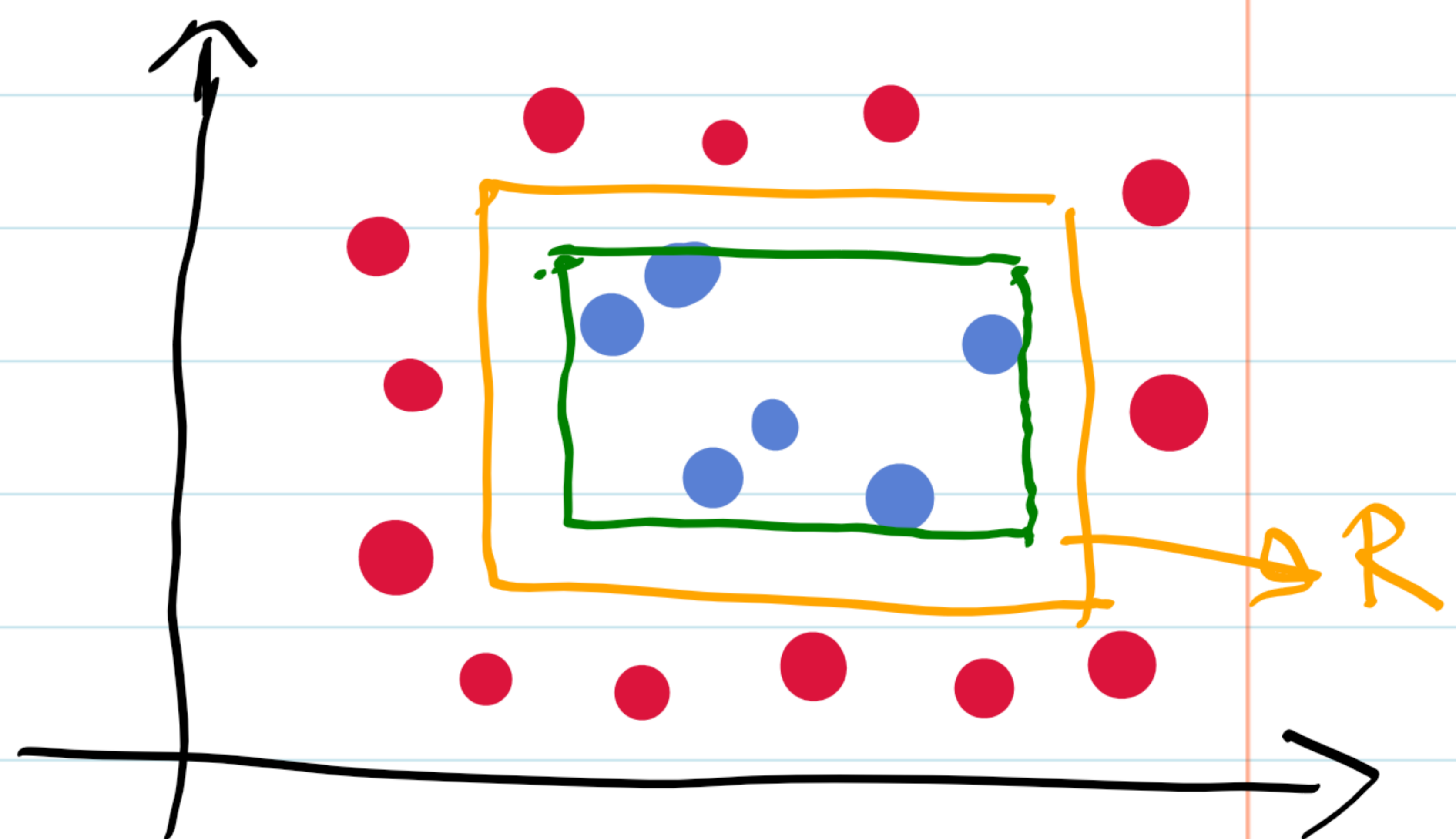
Question: Is this problem PAC-learnable?

- Mathematically we need to propose an algorithm A which takes the training data and returns R' , such that for any $\epsilon > 0$ & $\delta > 0$, there exist an N (which is a func. of ϵ & δ) with

$$P[|E_{in}(R') - E_{out}(R')| > \epsilon] \leq \delta.$$

- Proposed algorithm: Give me the set of data points, find the tightest rectangle that covers the blue dots.

Intuition: As N grows, we can find R' which is getting closer and closer to R .



So for any $\epsilon > 0$, $\delta > 0$ it is possible that as long as N is large enough, we will be able to make training error close to testing error.

Proof - Let R be the target func. Fix $\epsilon > 0$.

Let $P[R]$ denote the probability mass of the region defined by R . i.e., the probability that a point randomly drawn according to distribution P falls within R .

We assume $P[R] > \epsilon$.

We can define 4 rectangular regions π_1, π_2, π_3 & π_4 along the sides of R , each with probability at least $\epsilon/4$.

Let l, r, b & t be the four real values defining

$$R: R = [l, r] \times [b, t].$$

π_4 is defined by $\pi_4 = [l, s_4] \times [b, t]$, with

$$s_4 = \inf \left\{ s : P[[l, s] \times [b, t]] \geq \epsilon/4 \right\}.$$

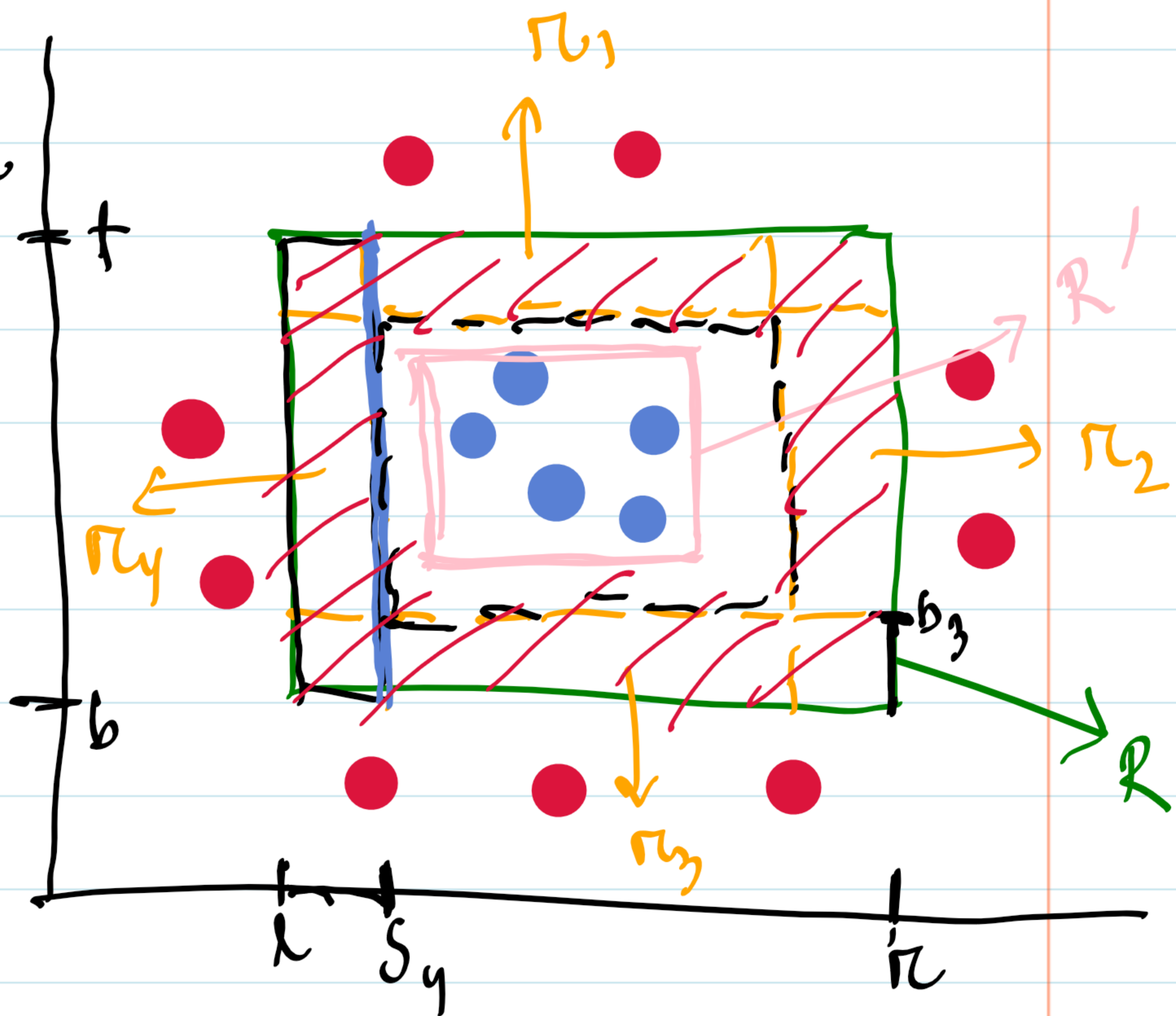
$\bar{\pi}_4$ = excluding the right most, $P[\bar{\pi}_4] \leq \epsilon/4$.

Observe here if R' meets all of these four regions $\pi_i, i=1, 2, 3, 4$, because R' is a rectangle, it will have one side in each of these regions.

$$E_{\text{out}}(R') \leq \epsilon = \left(4 \cdot \frac{\epsilon}{4}\right)$$

$$A_i = \left\{ R' \cap \pi_i \neq \emptyset \forall i \right\}, \quad B = \left\{ E_{\text{out}}(R') \leq \epsilon \right\}$$

$$\begin{aligned} \bigcap_{i=1}^4 A_i \subset B &\Rightarrow B^c \subset \left(\bigcap_{i=1}^4 A_i \right)^c \\ &= \bigcup_{i=1}^4 A_i^c \end{aligned}$$



$$\Rightarrow P(B^c) \leq P\left(\bigcup_{i=1}^4 A_i^c\right)$$

$$\Rightarrow P(E_{\text{out}}(R') > \epsilon) \leq P\left(\bigcup_{i=1}^4 \{R' \cap \pi_i = \emptyset\}\right)$$

$$\leq \sum_{i=1}^4 P(\{R' \cap \pi_i = \emptyset\})$$

$$(P[\pi_i] \leq \epsilon/4)$$

$$\leq 4 (1 - \epsilon/4)^m$$

no. of data points

$$\left(1 - x \leq e^{-x} \text{ for } x \in \mathbb{R}\right)$$

$$\leq 4 \exp(-m\epsilon/4)$$

Thus for any $\delta > 0$, we have

$$P[E_{\text{out}}(R') > \epsilon] \leq \delta$$

$$4 \exp(-\epsilon m/4) \leq \delta$$

$$\Rightarrow m \geq \frac{4}{\epsilon} \log \frac{4}{\delta}$$

$$P[|E_{\text{in}}(R') - E_{\text{out}}(R')| > \epsilon] \leq \delta.$$

Thus for any $\epsilon > 0$ & $\delta > 0$, if the sample size m is greater than $\frac{4}{\epsilon} \log \frac{4}{\delta}$, then

$$P[E_{\text{out}}(R') > \epsilon] \leq \delta. \quad (\text{PAC-learnable})$$

Guarantee VS Possibility

Basically we see the difference between the deterministic and probabilistic learning.

Deterministic: Can D tell us something about f outside of D ?

Certain

NO

Probabilistic: "Can D tell us something possibly about f outside of D ?" - YES

One hypothesis VS final hypothesis:

In Hoeffding, $P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$.
hypothesis h is fixed.

We need to choose h before we look at the data set not after.

When we need to choose g from h_1, \dots, h_M , we need to repeat Hoeffding M times.

The factor 'M':

We can $|E_{in}(g) - E_{out}(g)| > \epsilon$

$\Rightarrow |E_{in}(h_1) - E_{out}(h_1)| > \epsilon$

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

or ...

$|E_{in}(h_M) - E_{out}(h_M)| > \epsilon$

possible cases

$$P(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq \sum_{m=1}^M P(|E_{in}(h_m) - E_{out}(h_m)| > \epsilon)$$

\rightarrow If $A \Rightarrow B$, $P(A) \leq P(B)$

\rightarrow Union bound $P(A \text{ or } B) \leq P(A) + P(B)$.

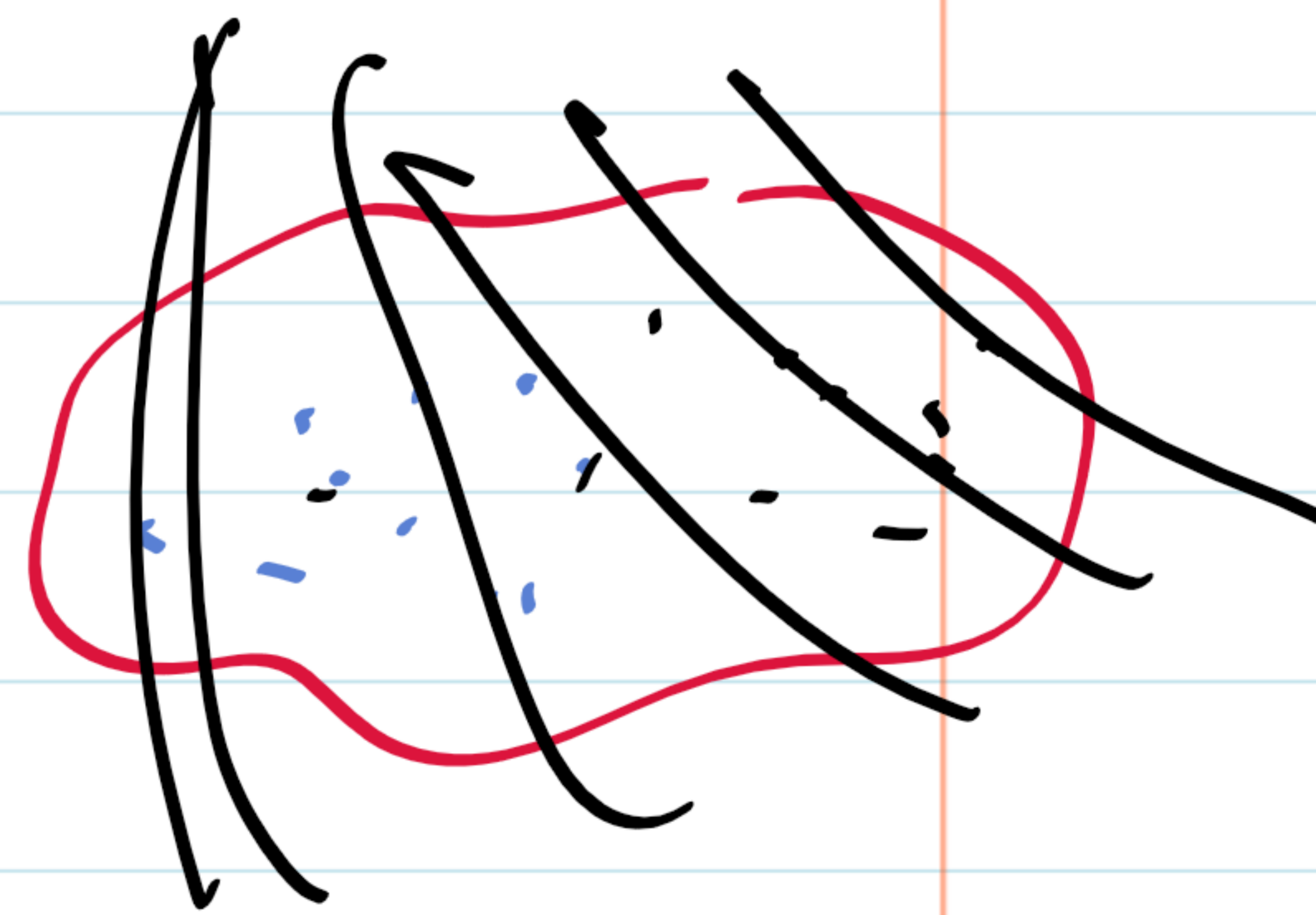
Hoeffding: $P(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq 2Me^{-2\epsilon^2 N}$

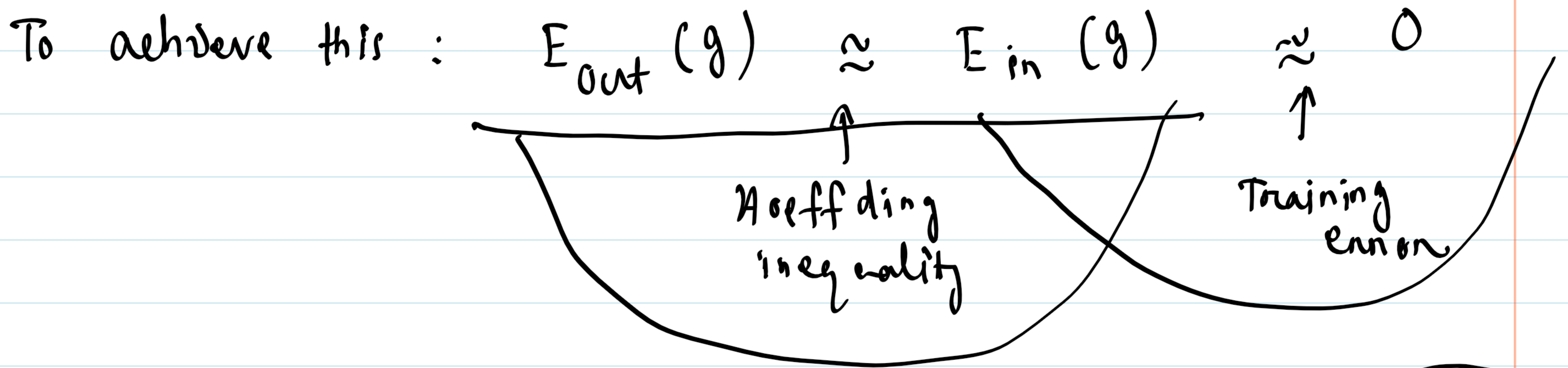
M is a constant.

Bad: M can be large, or even ∞ .

Good: We can bound this M .

Learning goal: $E_{out}(g) \approx 0$.

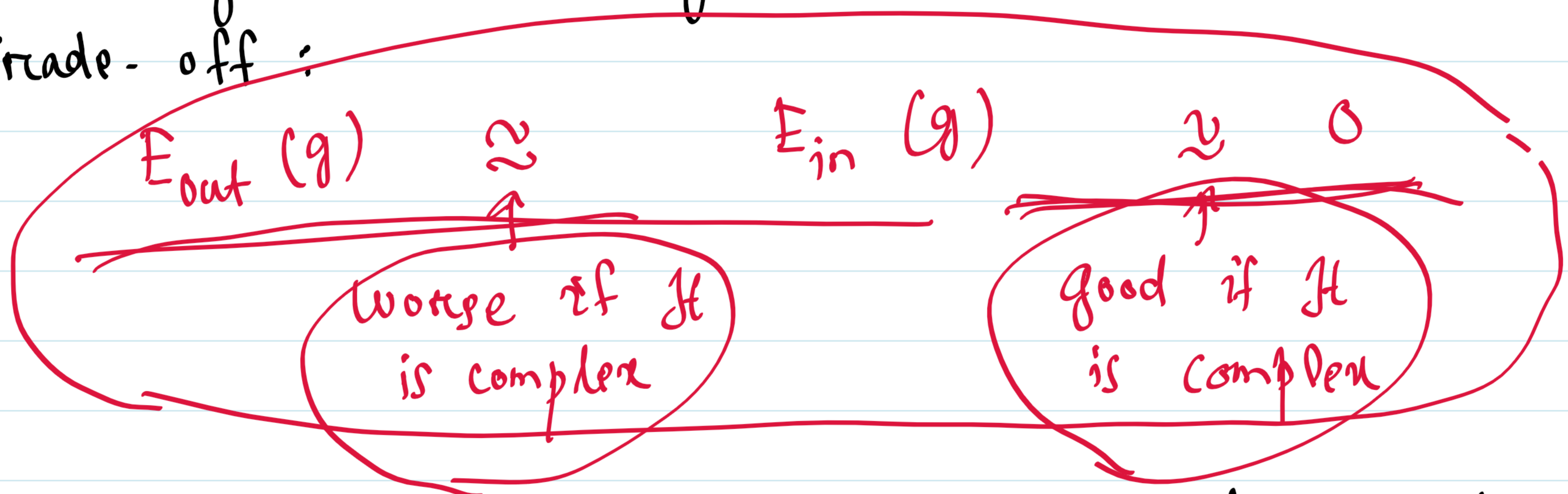




Complex \mathcal{H} :

Hoeffding : $P [|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2M e^{-2\epsilon^2 N}$

- If \mathcal{H} is complex, M will be large, so the approximation will worsen.
- If \mathcal{H} is complex, you have more options during training. So training error is improved.
- Trade-off :



- We can not use a very complex model. Simple model generalise better.

Complex f (target func) :

- good : Hoeffding is not affected by f .
- Bad : If f is complex, then it will be very hard to train, so training error cannot be small.
- Trade-off :

