

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

get rid of
give bound

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

Equivalently, with probability $1 - \delta$,

$$E_{in}(g) - \epsilon \leq E_{out}(g) \leq E_{in}(g) + \epsilon$$

$$\text{So } \delta = 2M e^{-2\epsilon^2 N} \Rightarrow \epsilon = \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}}$$

$$E_{in}(g) - \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}} \leq E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}}$$

Generalization bound

- N is training sample. (the more, the better)
- δ : probability tolerance level, $1 - \delta$: "confidence".
Small δ : 'Conservative', So needs to have large N to compensate for $\log 1/\delta$.
- M : Model complexity, Large M : Complicated model.
So we need large N to compensate for $\log M$.

Upper limit:

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}}$$

→ E_{out} can not be worse than $E_{in} + \epsilon$.

→ Performance guarantee: $E_{in} + \epsilon$ is the worst you will have.

Lower limit:

$$E_{in}(g) - \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}} \leq E_{out}(g).$$

→ E_{out} can not be better than $E_{in} - \epsilon$.

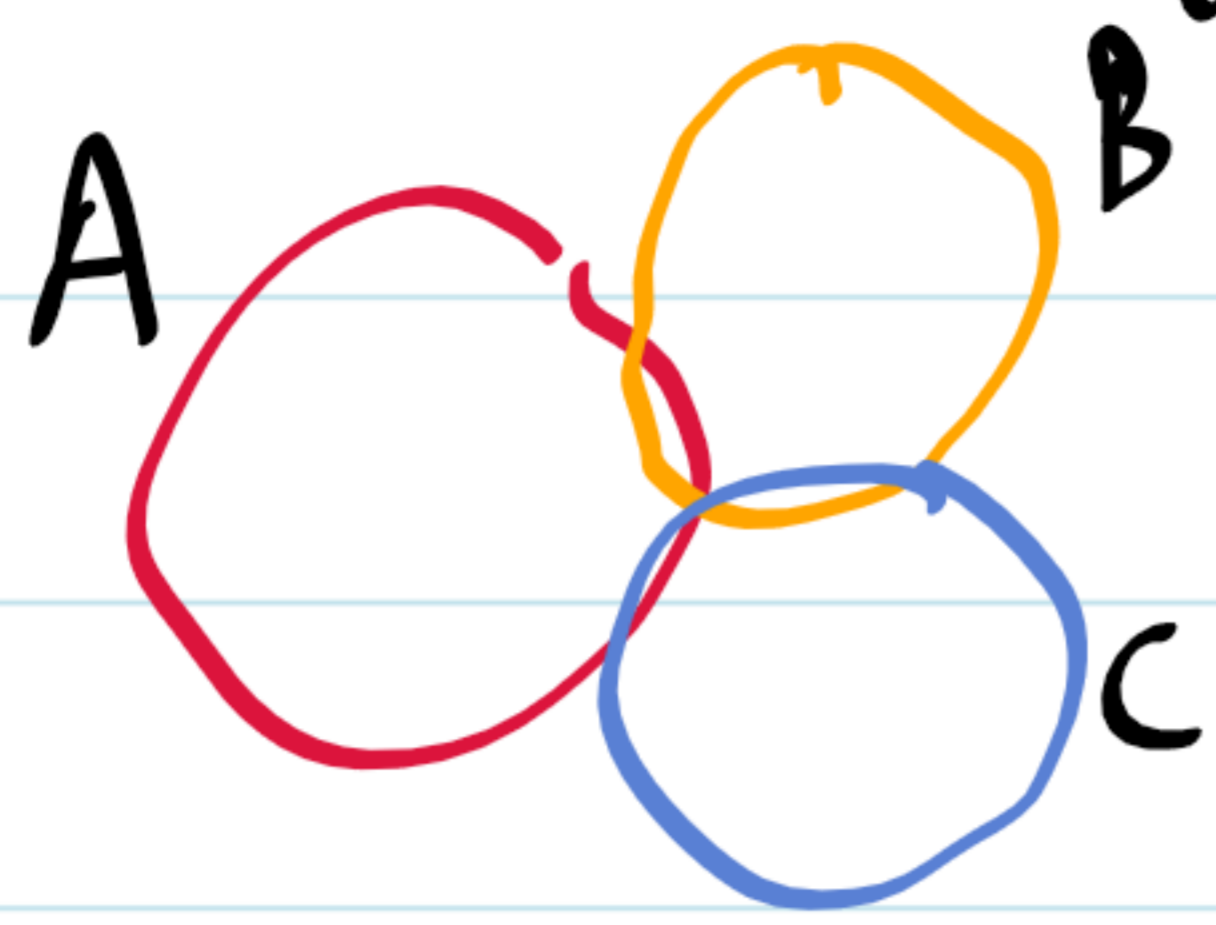
$N, M, 1 - \delta$

Growth function :

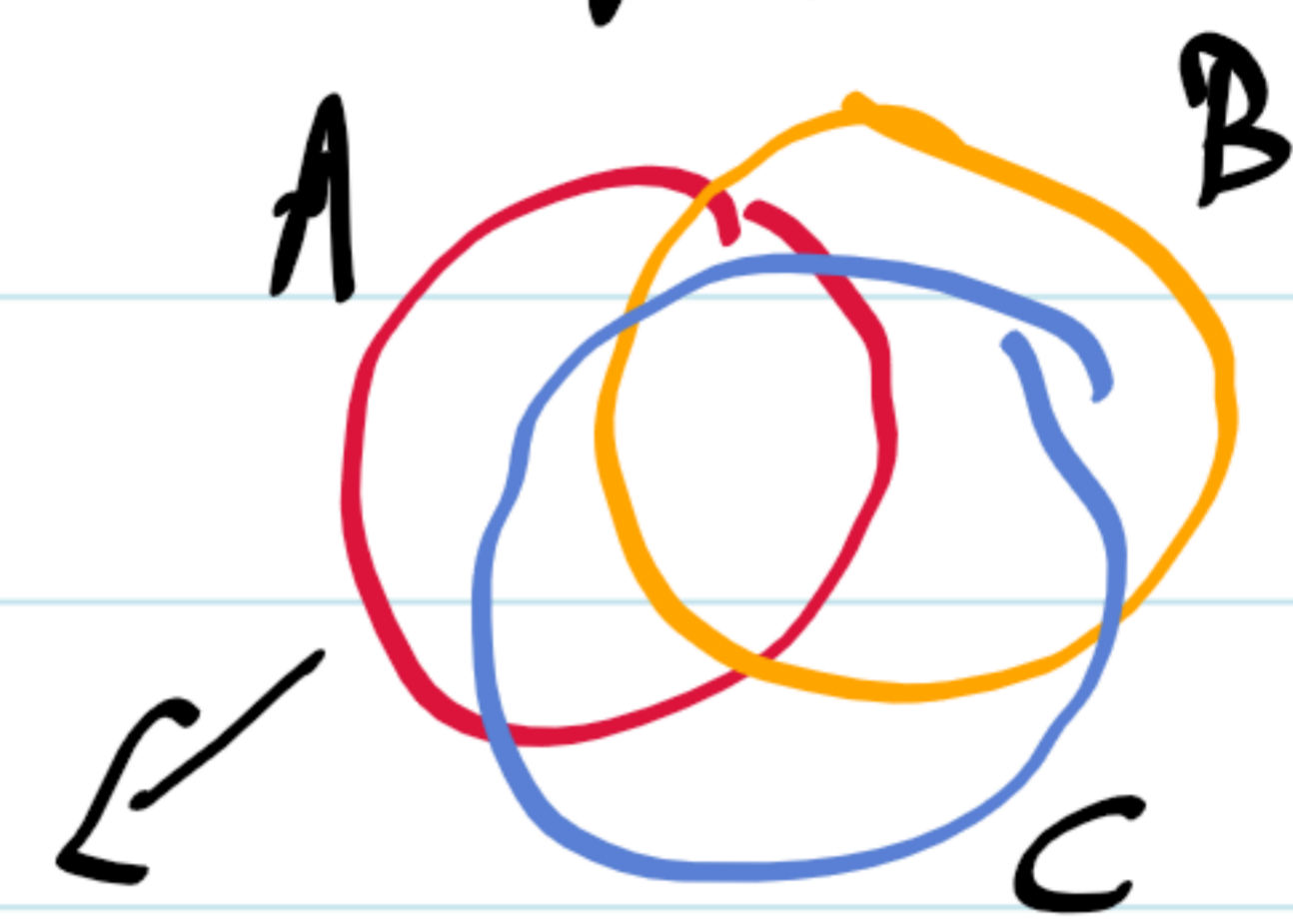
Let B_m be the (bad) event that " $|E_{in}(h_m) - E_{out}(h_m)| > \epsilon$ "

Then $P[B_1 \text{ or } B_2 \text{ or } \dots \text{ or } B_M] \leq P[B_1] + P[B_2] + \dots + P[B_M]$.

Two cases : weakly overlapping



Strongly overlapping



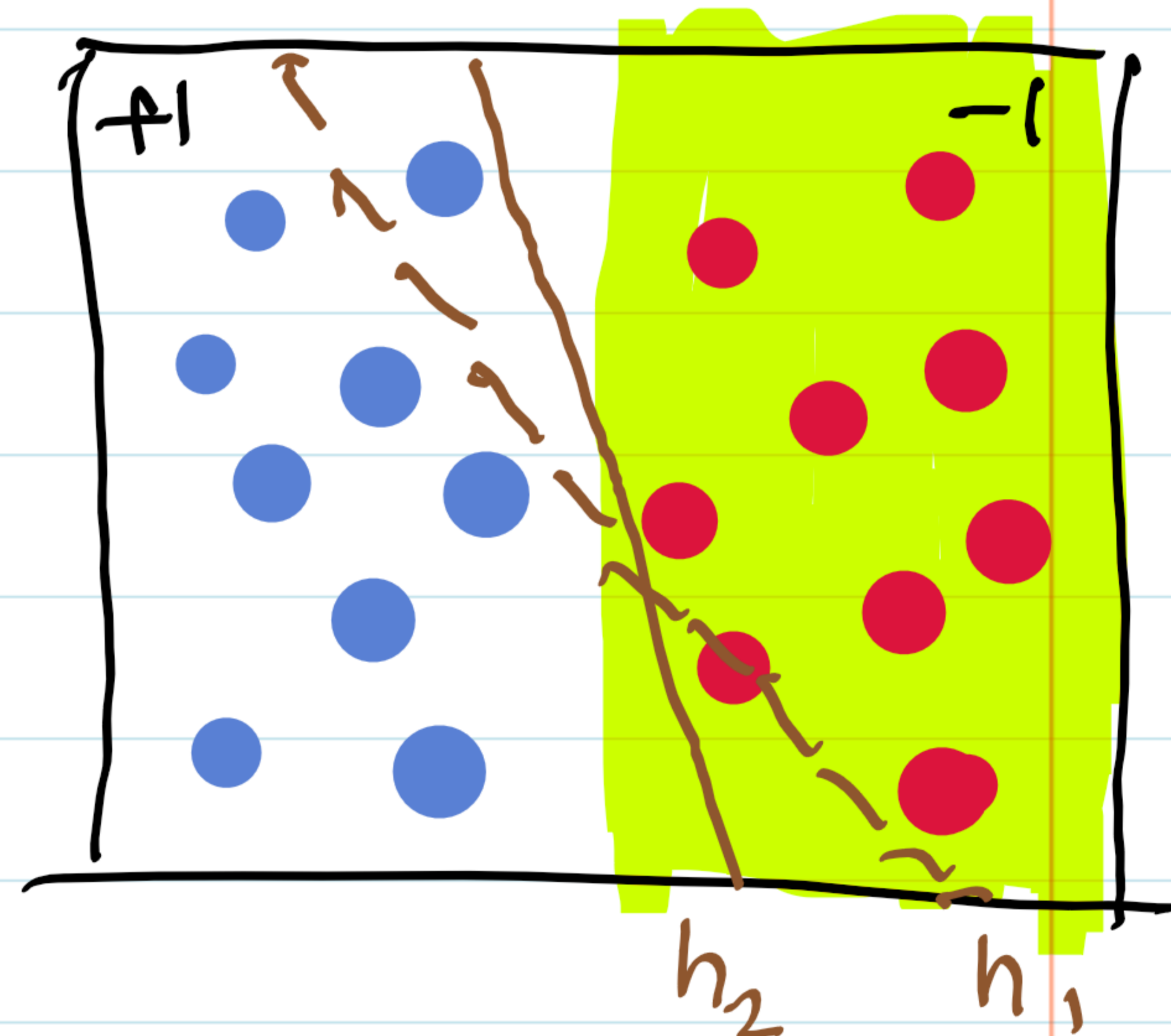
gives a gross over-estimate.

→ If h_1 is very similar to h_2 , then the two events ' $|E_{in}(h_1) - E_{out}(h_1)| > \epsilon$ ' and ' $|E_{in}(h_2) - E_{out}(h_2)| > \epsilon$ ' are likely to coincide for most data sets.

→ Replace M by an effective number '?

ΔE_{out} = change in the $+1$ & -1 area

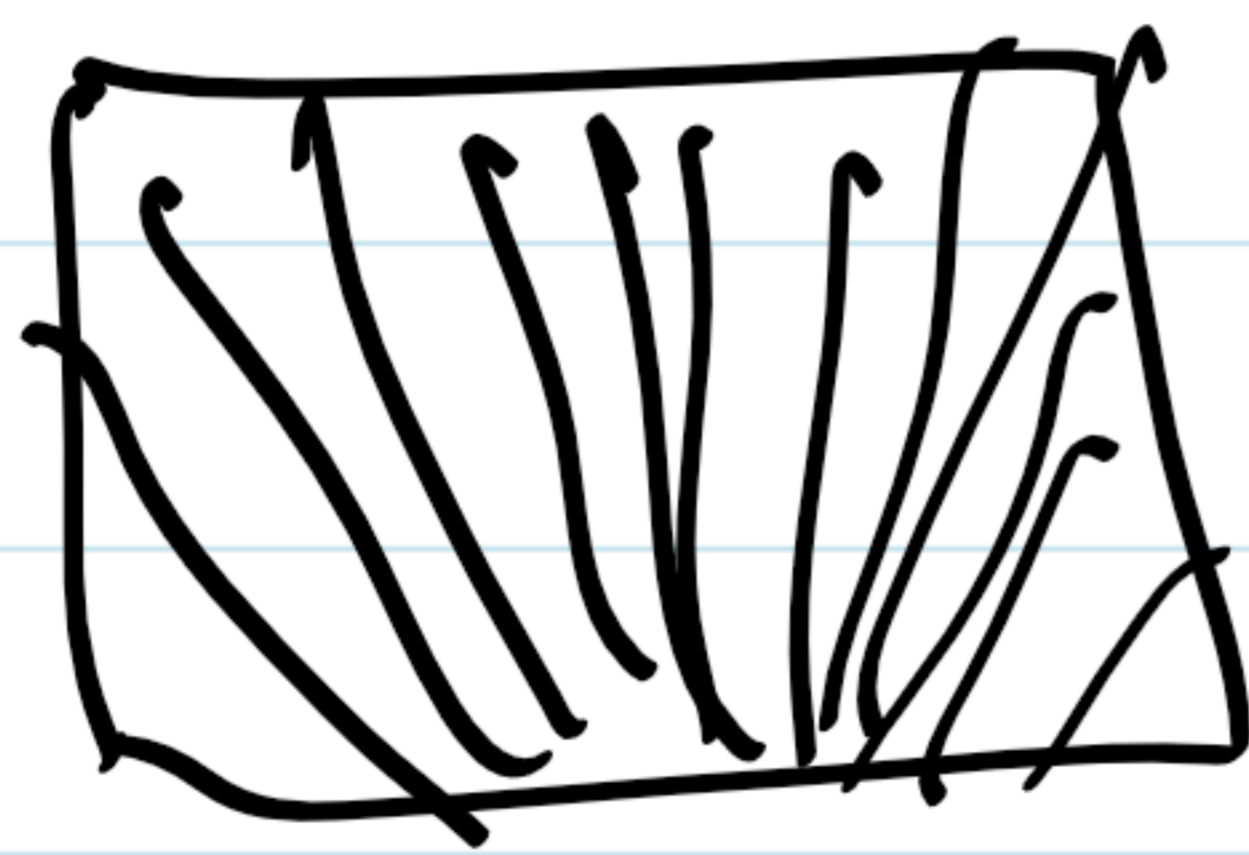
ΔE_{in} = change in the labels of training samples.



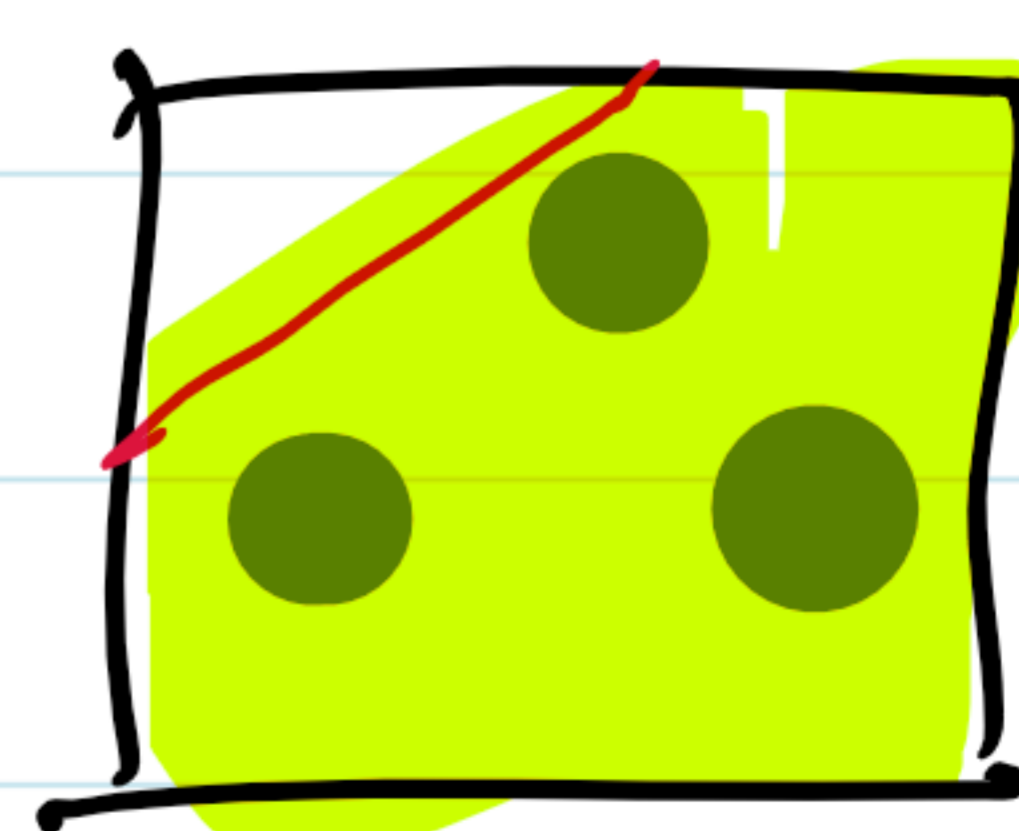
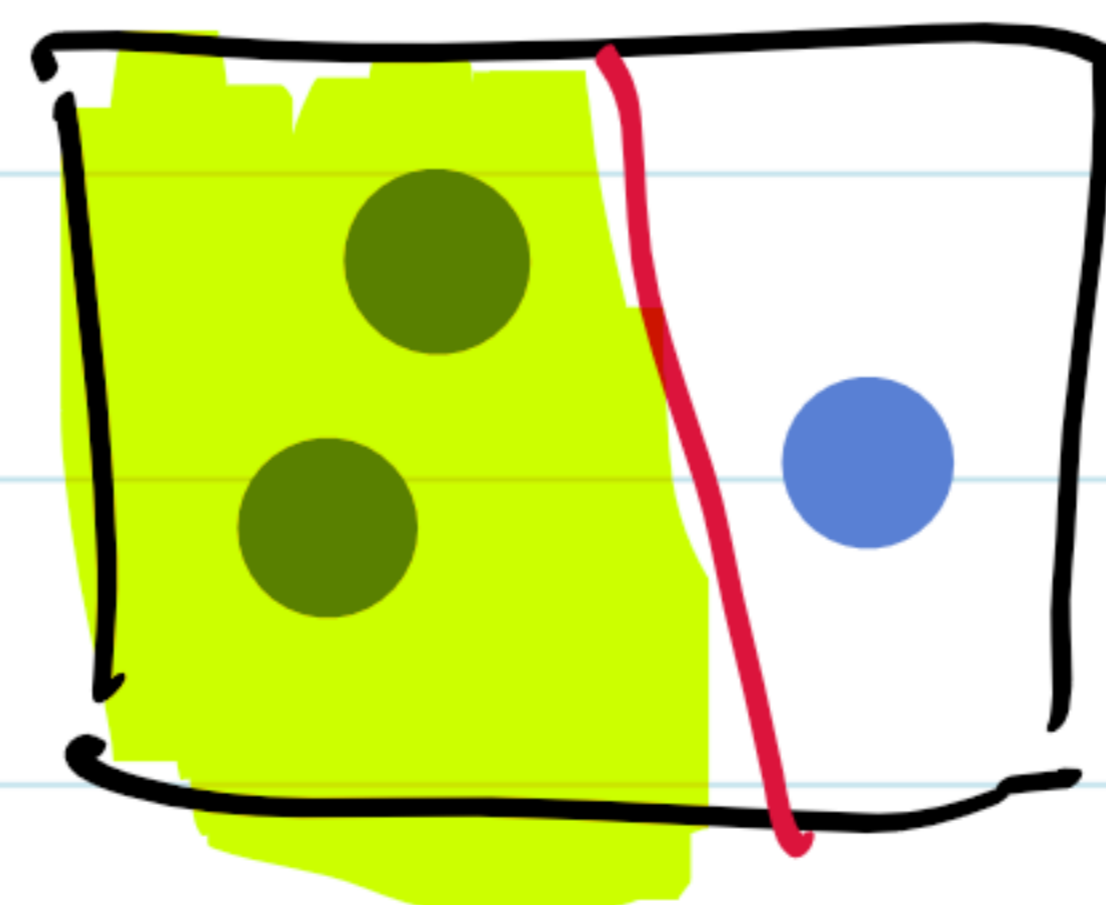
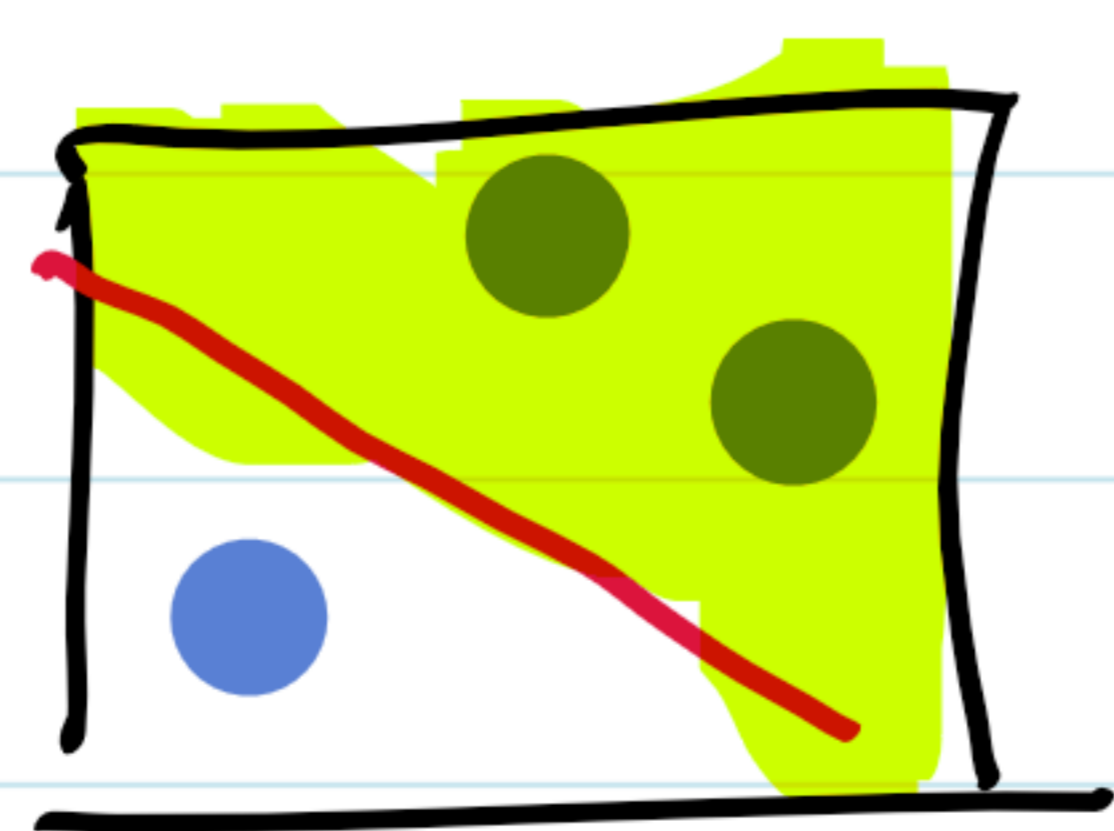
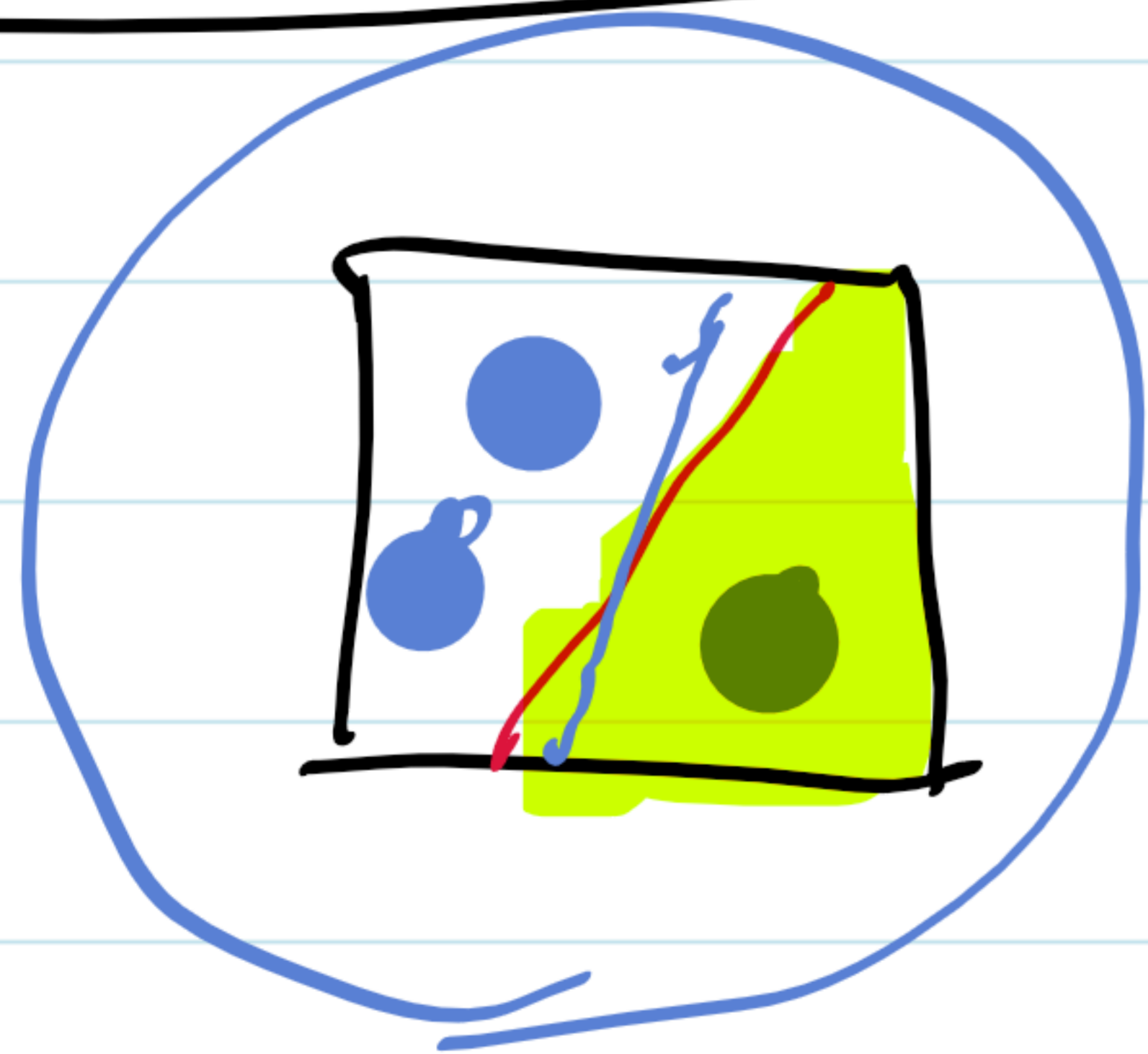
We should expect

$$P[|E_{in}(h_1) - E_{out}(h_1)| > \epsilon]$$

$$\approx P[|E_{in}(h_2) - E_{out}(h_2)| > \epsilon]$$



Restrict to training sets only :

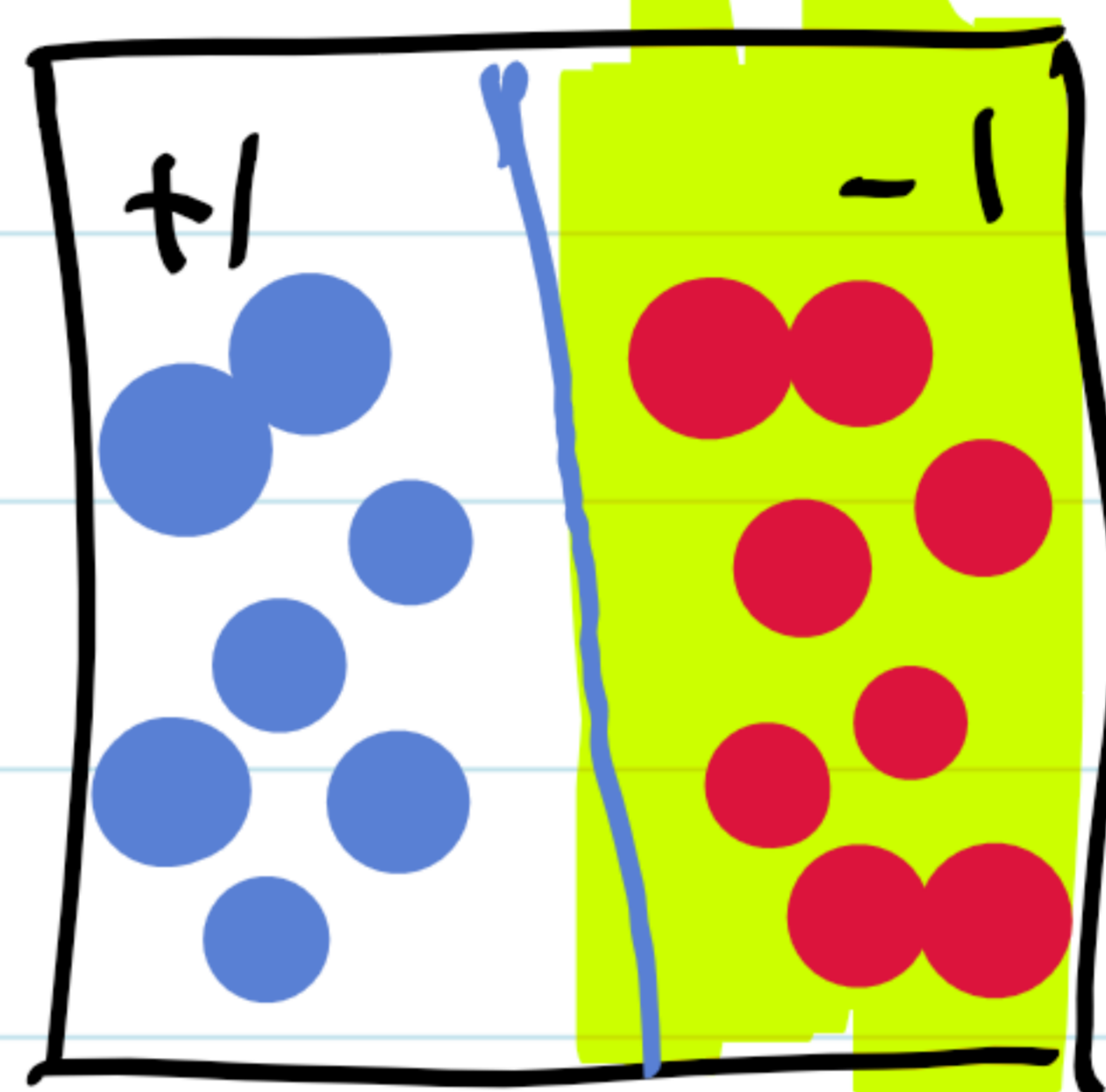
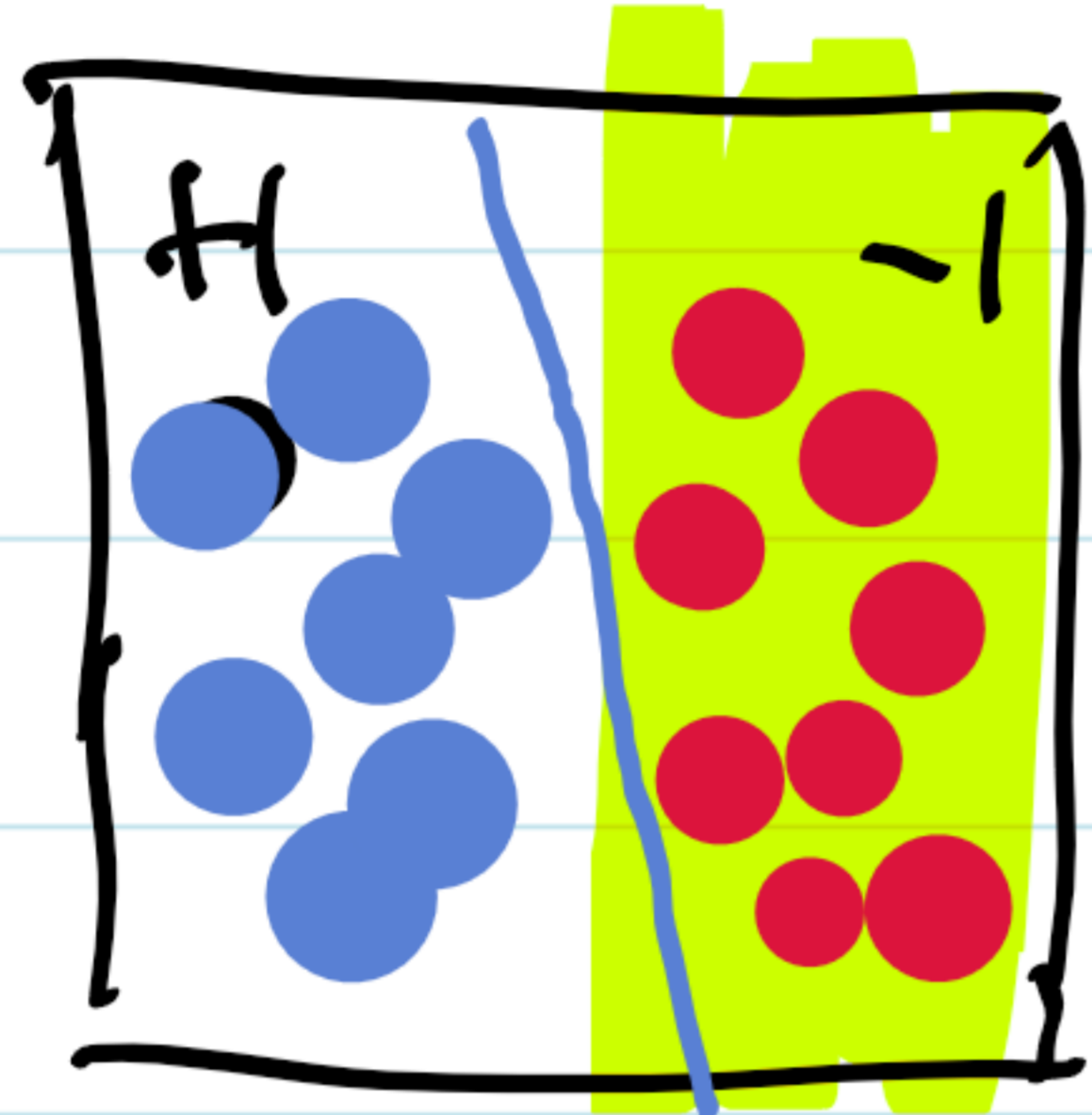


Idea : Just look at the training samples. Don't care until a training sample flips its sign.

Dichotomies : Means no. of mini-hypothesis.

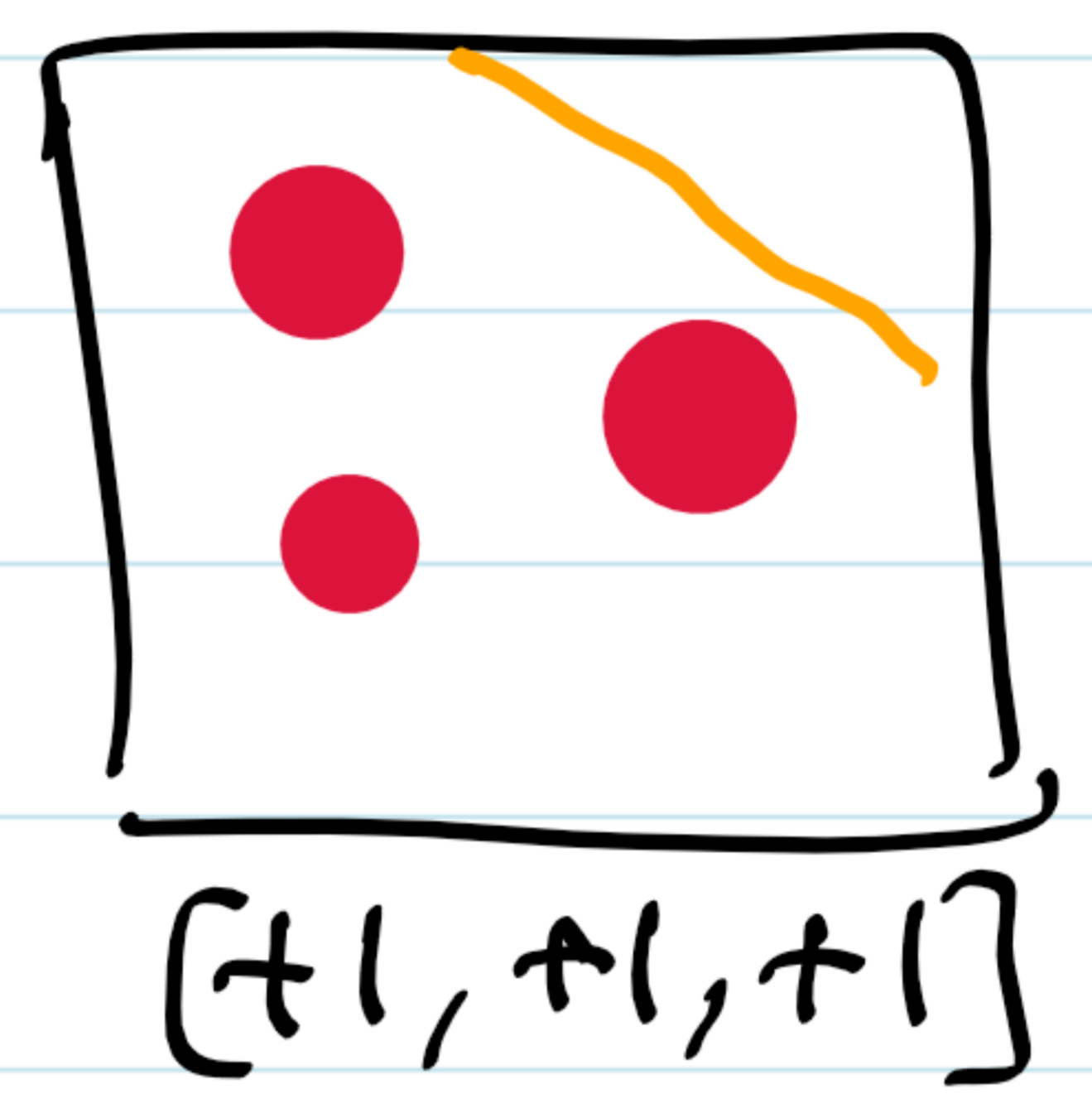
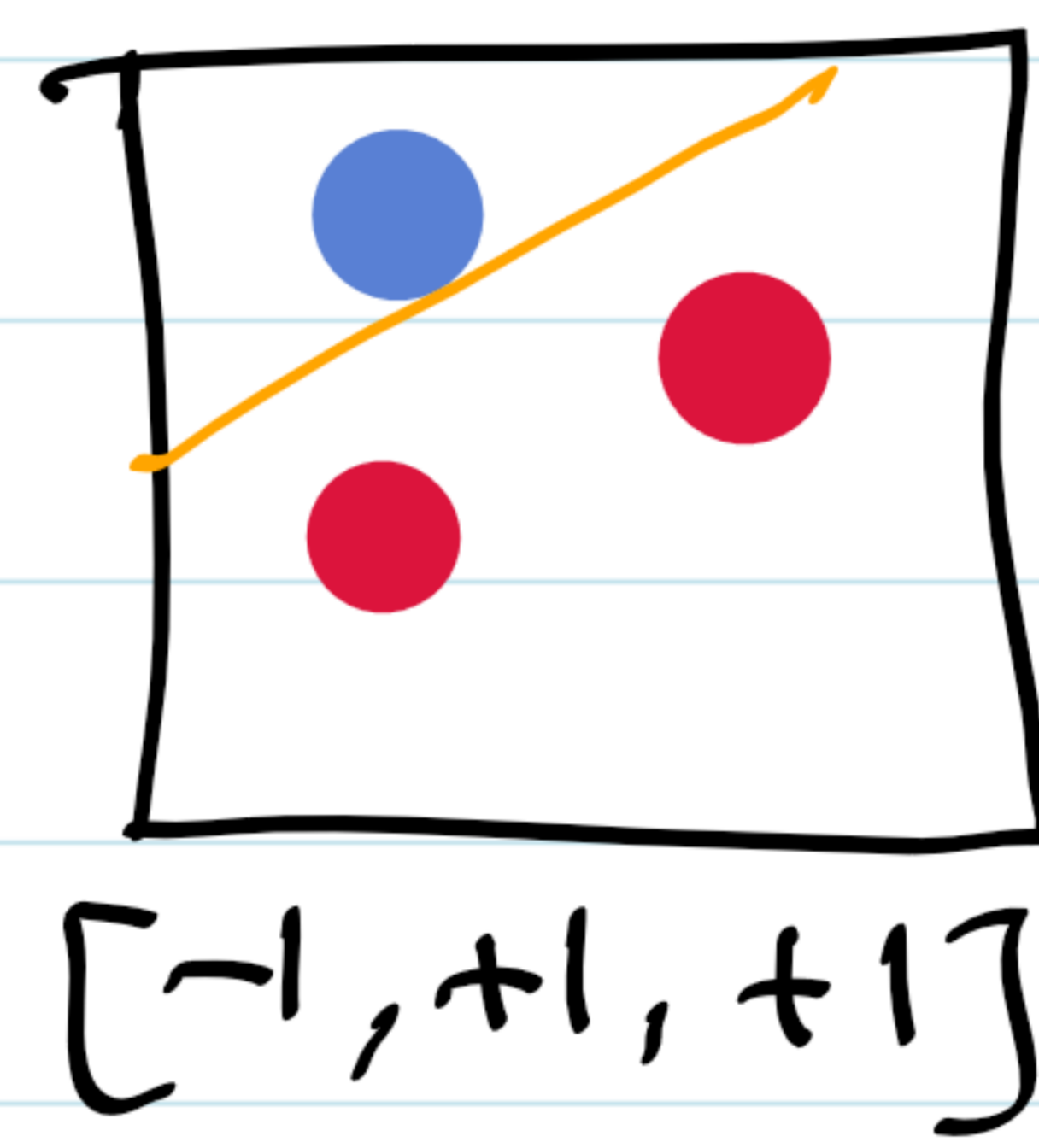
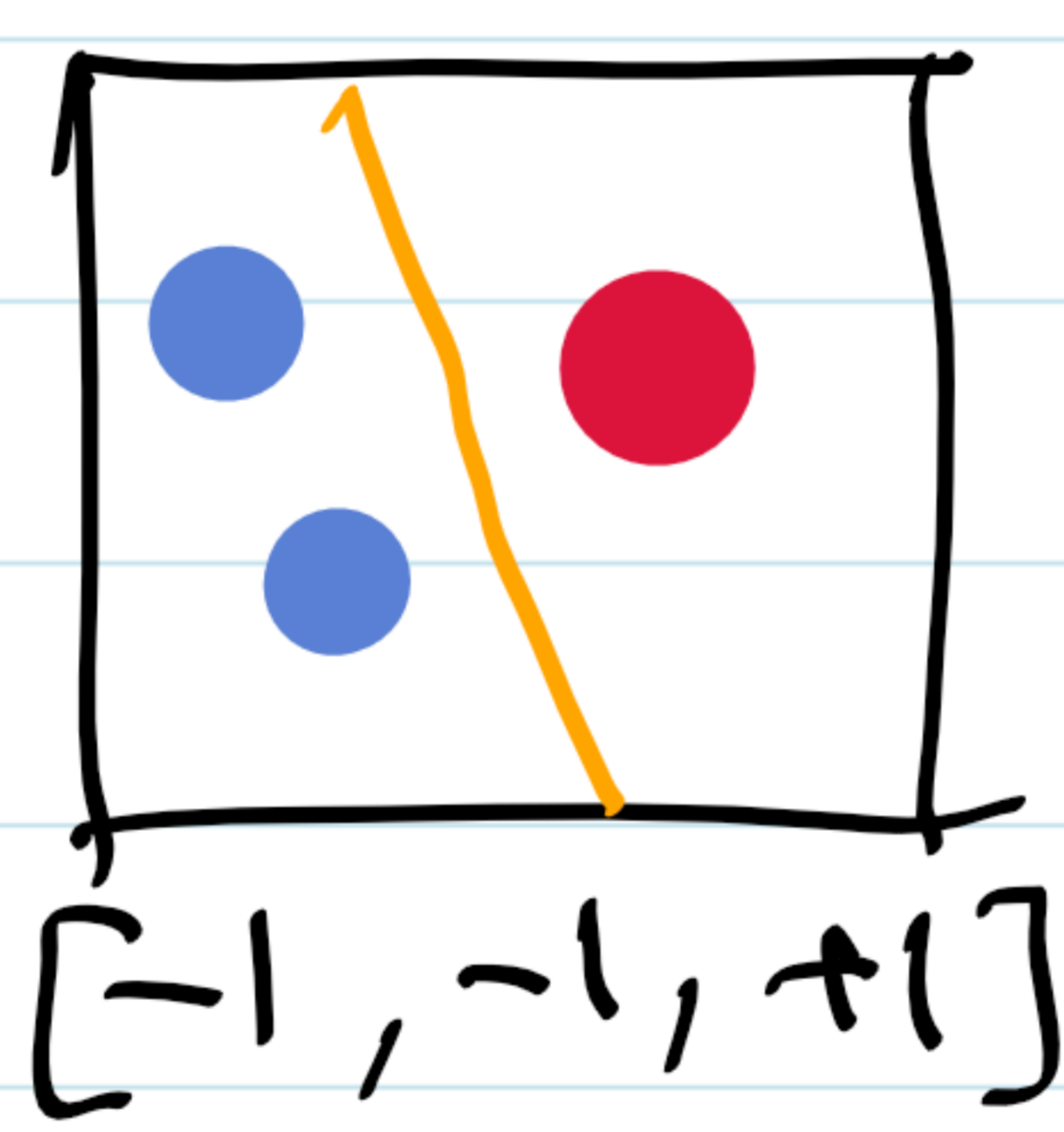
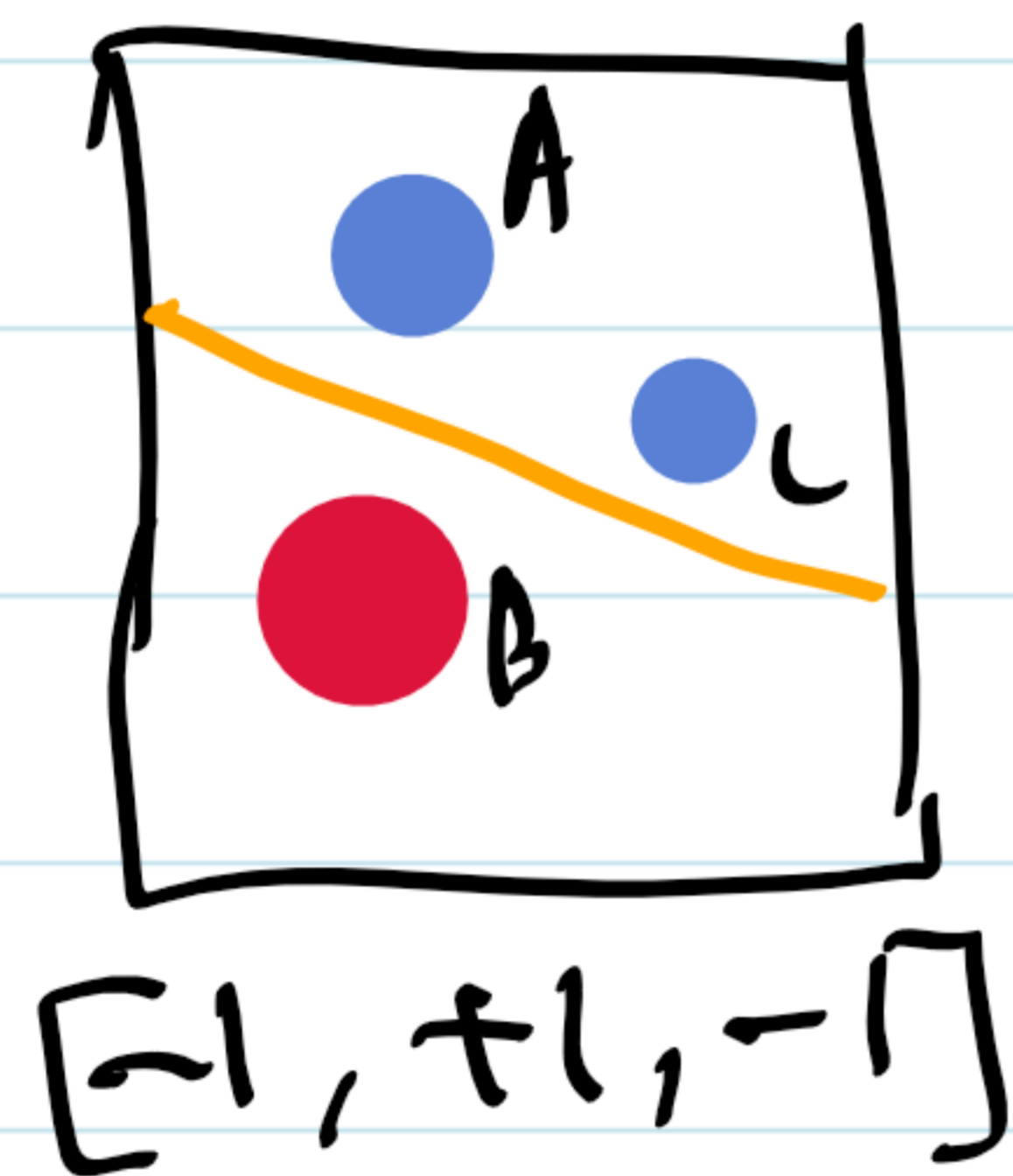
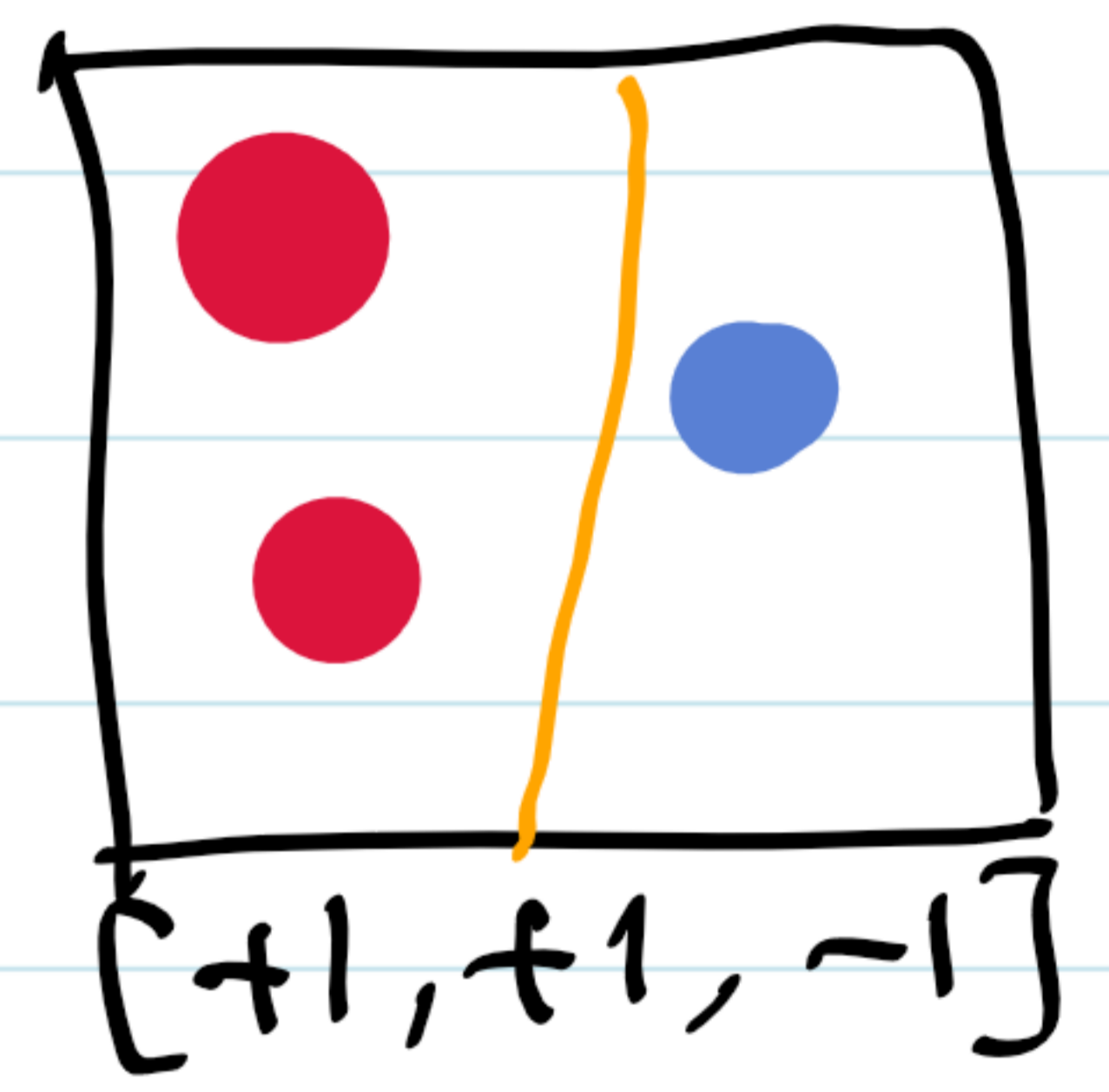
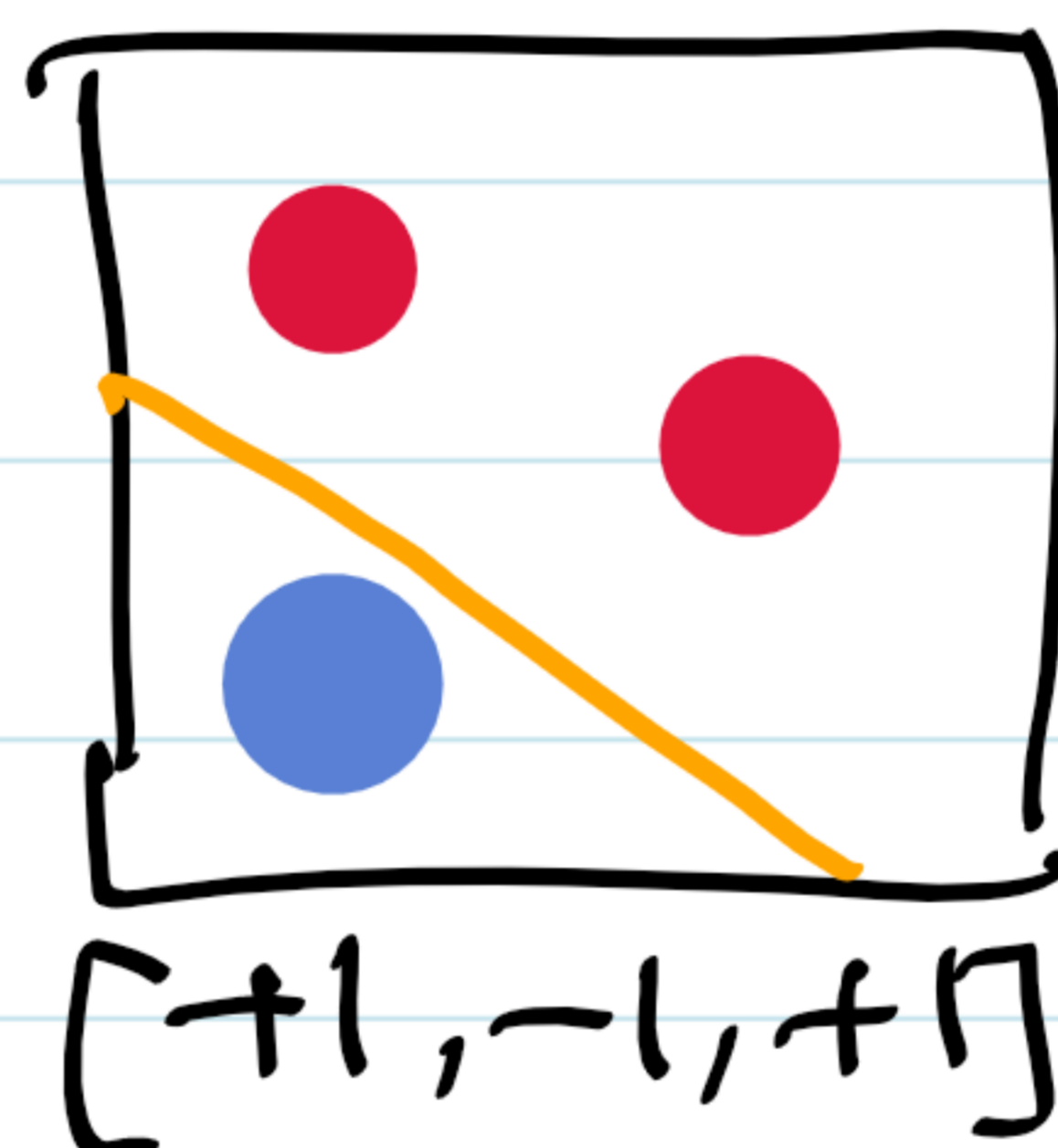
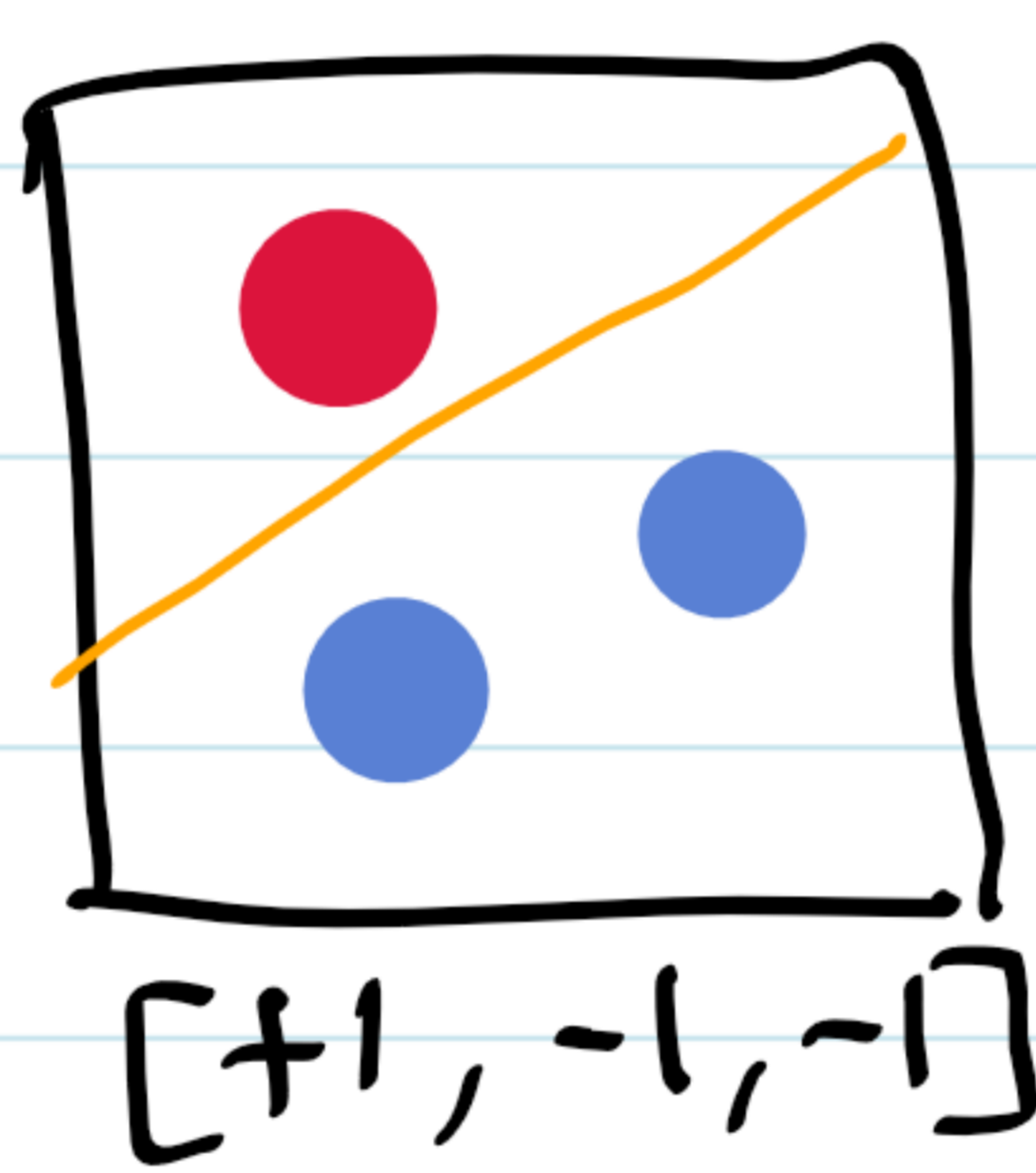
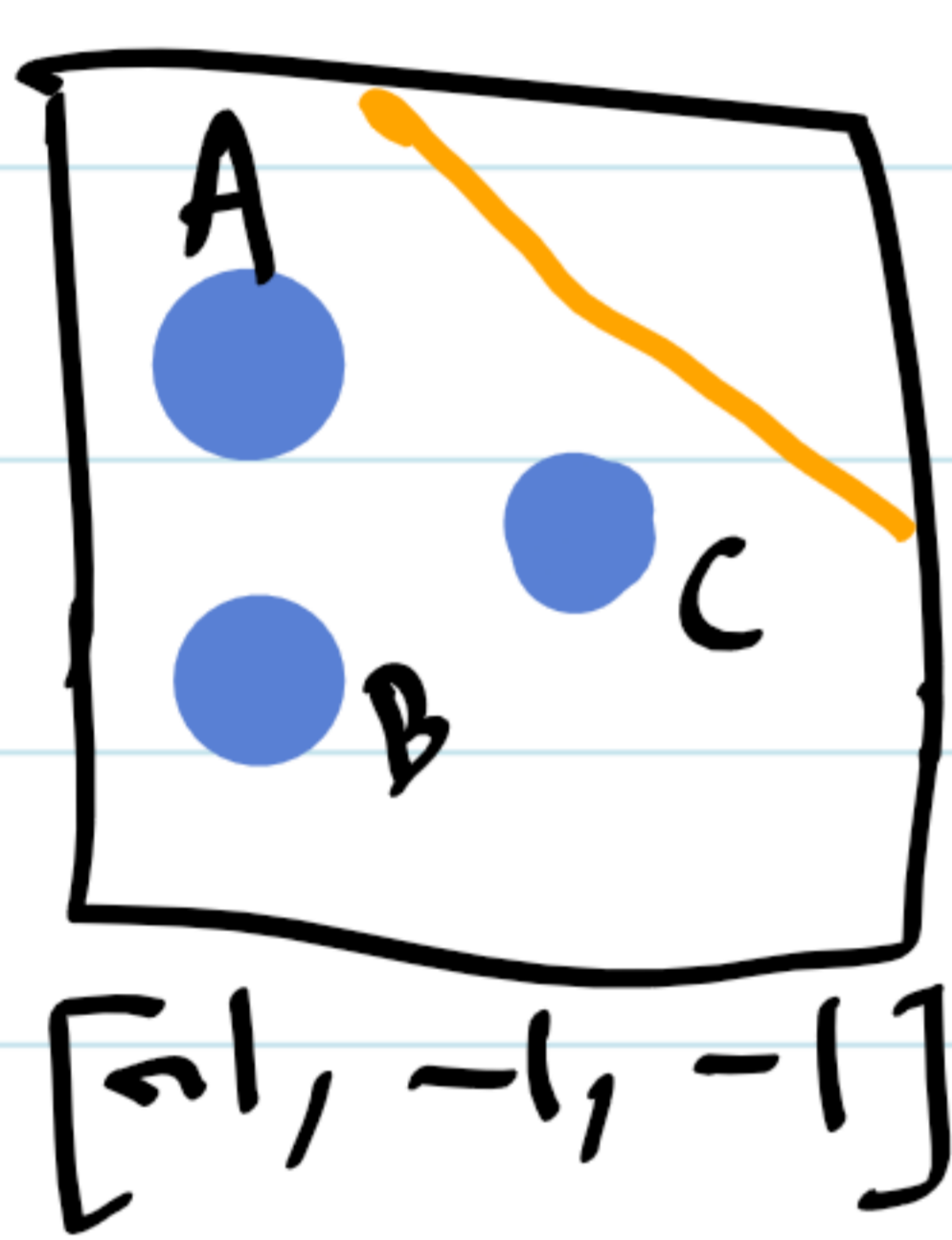
| Hypothesis | Dichotomy |
|---|--|
| $\rightarrow h: X \rightarrow \{+1, -1\}$ $\rightarrow \#$ population samples \rightarrow No. can be infinite | $\rightarrow h: \{x_1, \dots, x_N\} \rightarrow \{+1, -1\}$ \rightarrow for training samples only \rightarrow No. is at most 2^N . |

\rightarrow observe : Different hypothesis, same dichotomy.

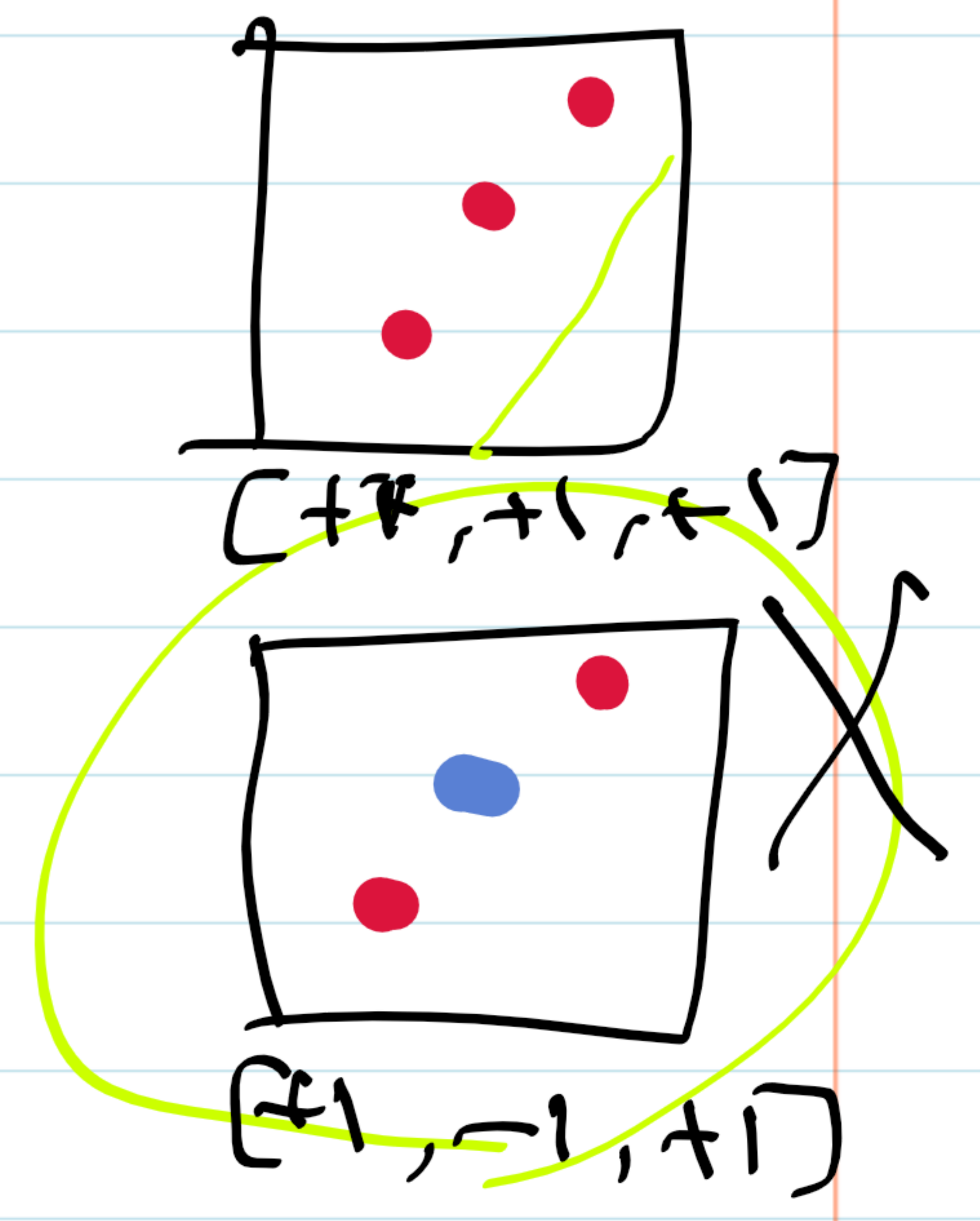
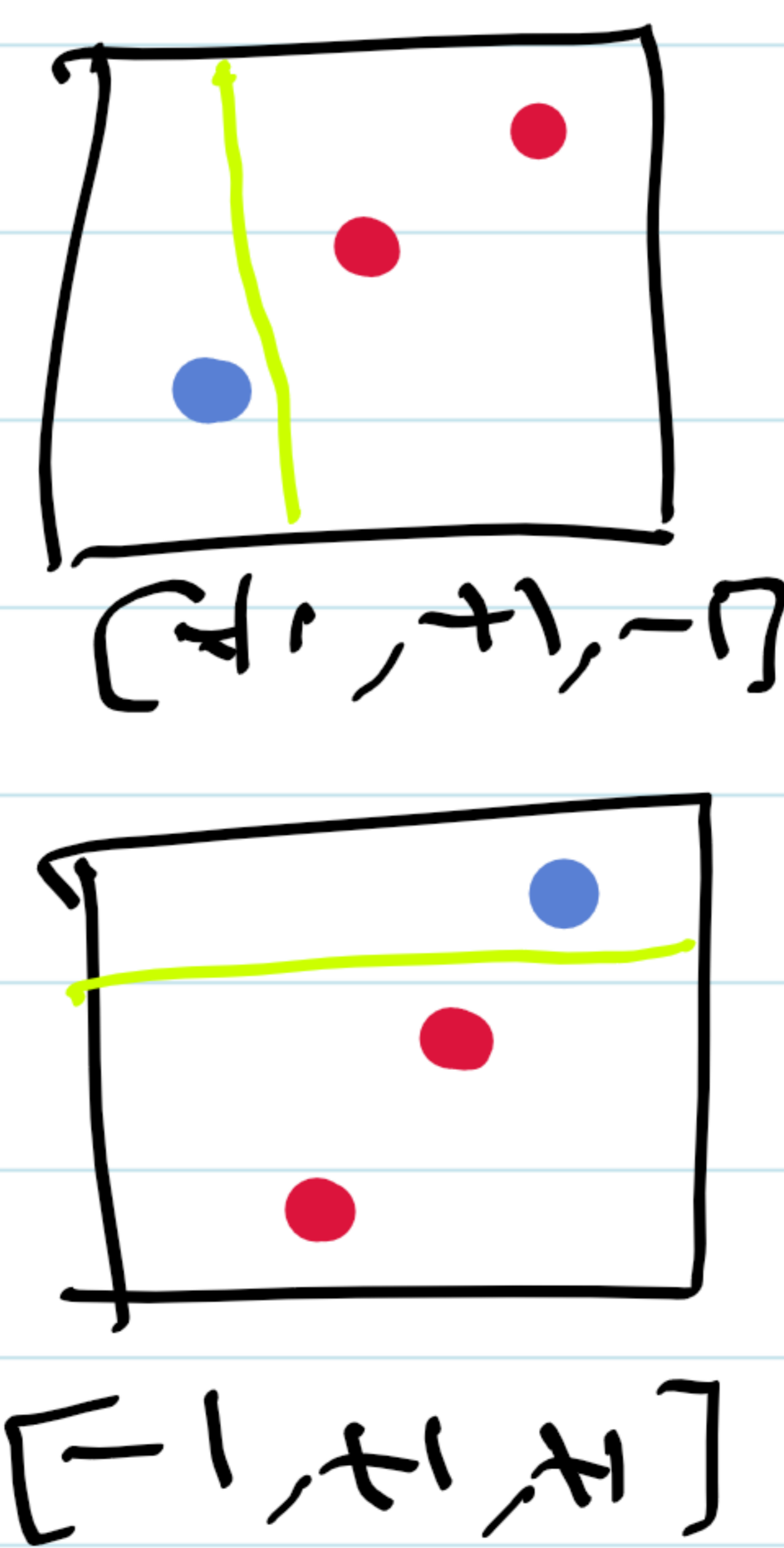
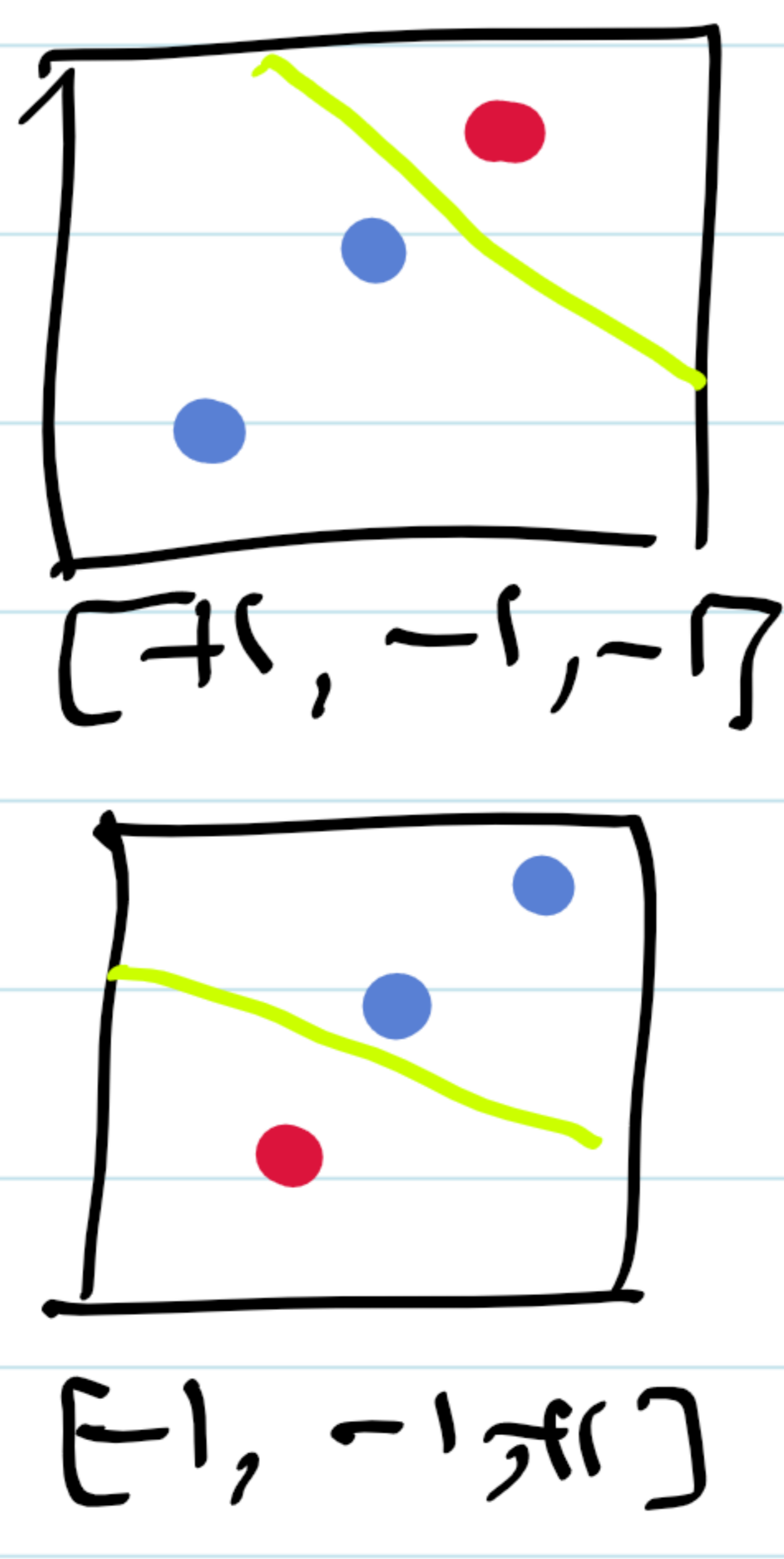
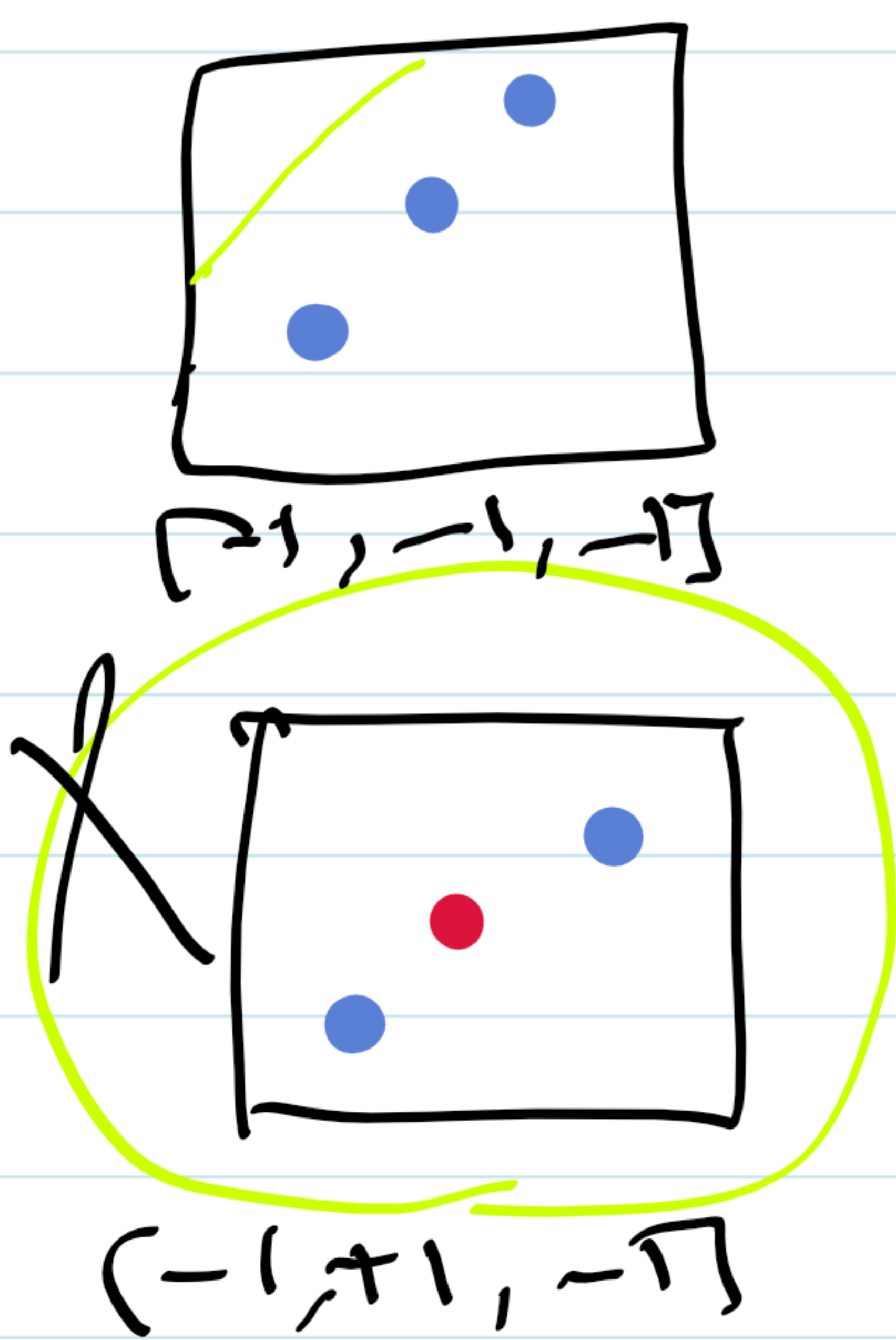


Defⁿ : Let $x_1, \dots, x_N \in X$. The dichotomies generated by \mathcal{H} on these points are

$$\mathcal{H}(x_1, \dots, x_N) = \{ (h(x_1), \dots, h(x_N)) \mid h \in \mathcal{H} \}$$



Not allowed :



Candidate to replace M :

Define growth function

$$m_{\mathcal{H}}(N) = \max_{x_1, \dots, x_N \in X} |\mathcal{H}(x_1, \dots, x_N)|,$$

here $|\cdot|$ denotes the cardinality (no. of elements) of a set.

In words, $m_{\mathcal{H}}(N)$ is the maximum number of dichotomies that can be generated by \mathcal{H} on any N points.

- Large $m_{\mathcal{H}}(N)$ = more complicated \mathcal{H} .
- $m_{\mathcal{H}}(N)$ depends on \mathcal{H} and N , but not on the algorithm and not on the probability distribution $p(x)$.
- To compute $m_{\mathcal{H}}(N)$: Consider all possible choices of N points x_1, \dots, x_N from X . Then pick the one that gives you the most dichotomies.
- For any \mathcal{H} , $\mathcal{H}(x_1, \dots, x_N) \subseteq \{-1, +1\}^N$.
So $m_{\mathcal{H}}(N) \leq 2^N$.

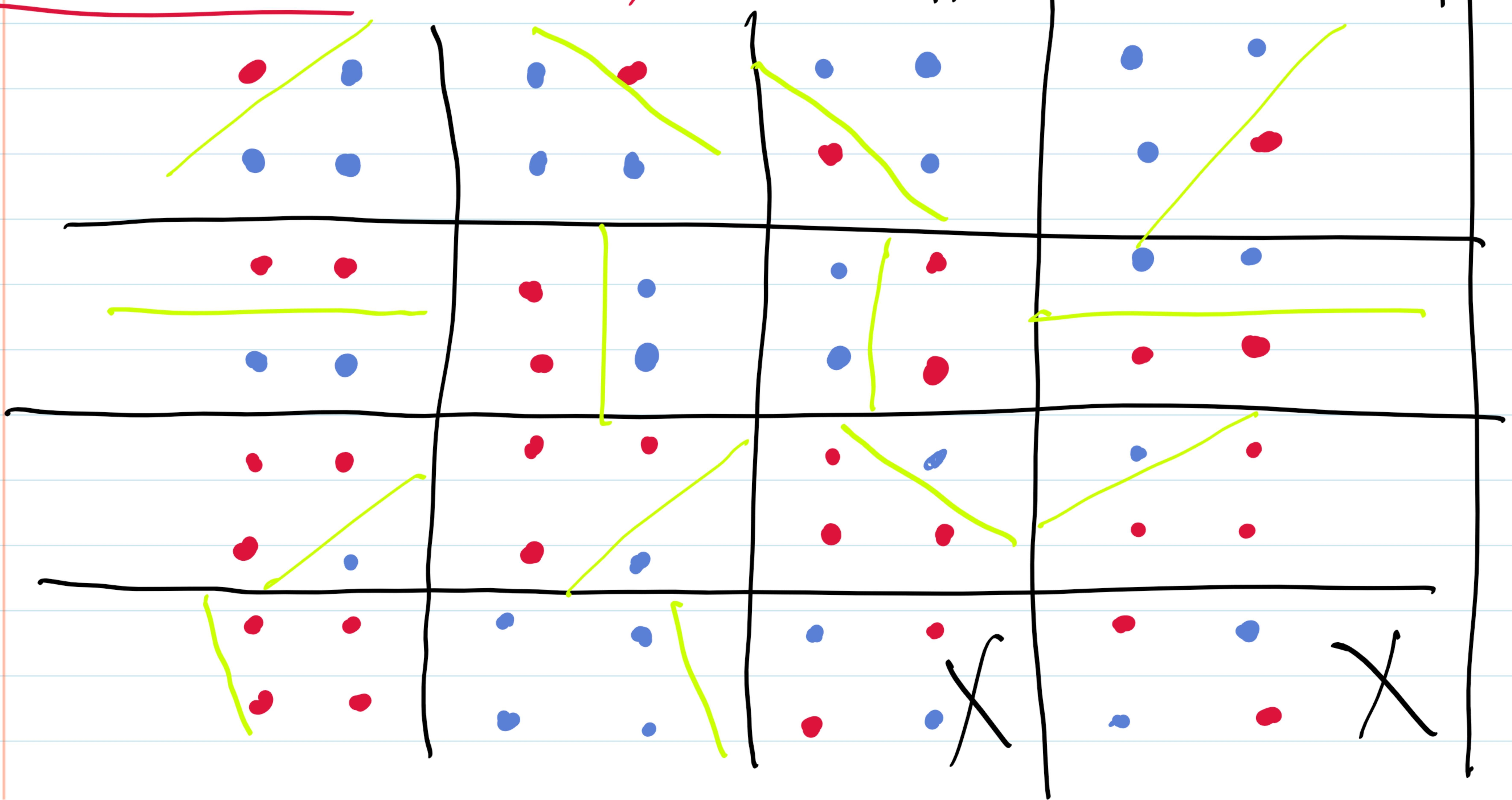
Example - 1 ($N=3$)

\mathcal{H} = linear models in 2D, $N=3$

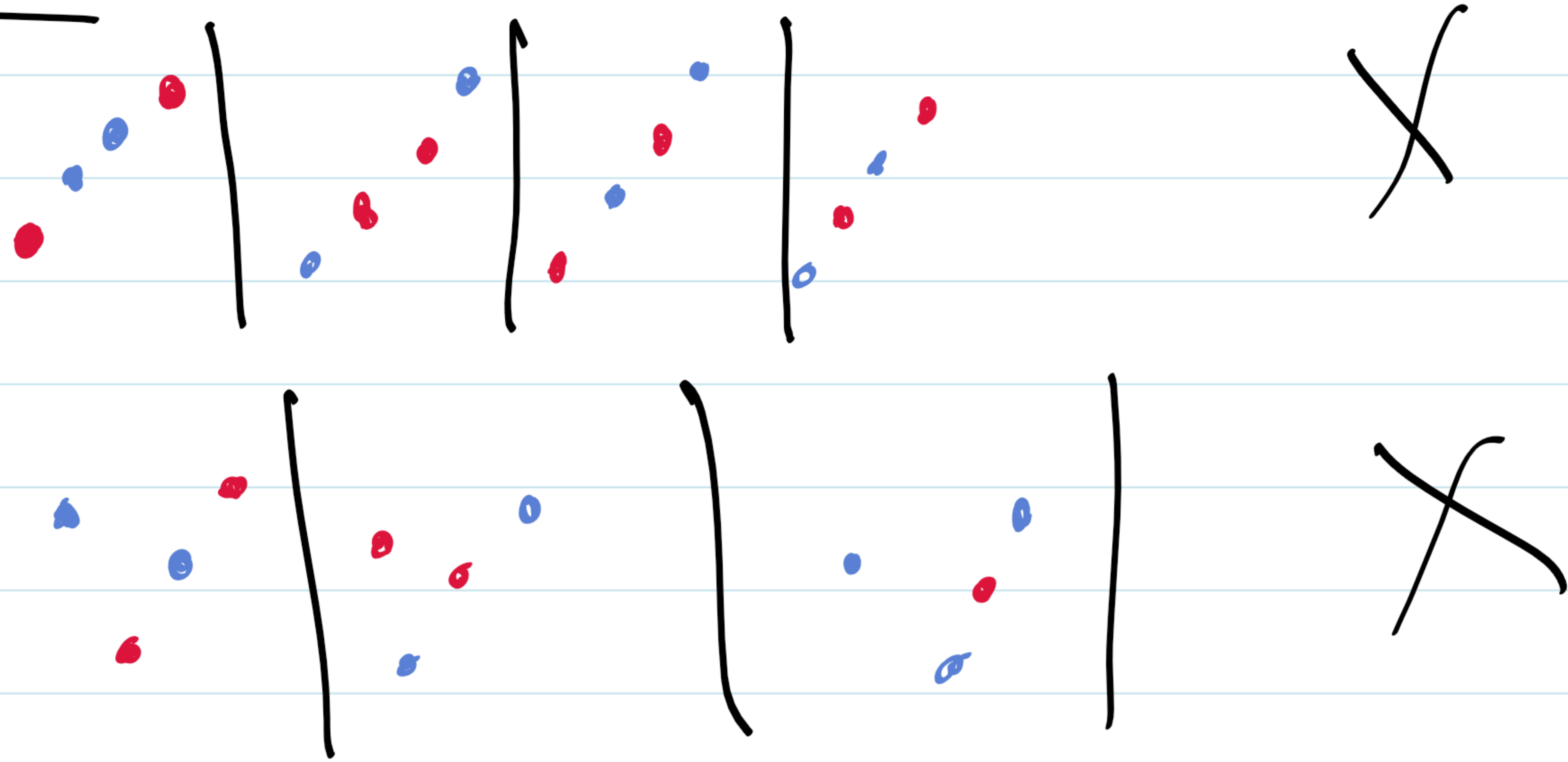
We have 8 dichotomies. (max.) Growth funⁿ ≤ 8 .

Example - 2 ($N=4$)

$$m_{\mathcal{H}}(N) = 14$$



$N=4$ Case

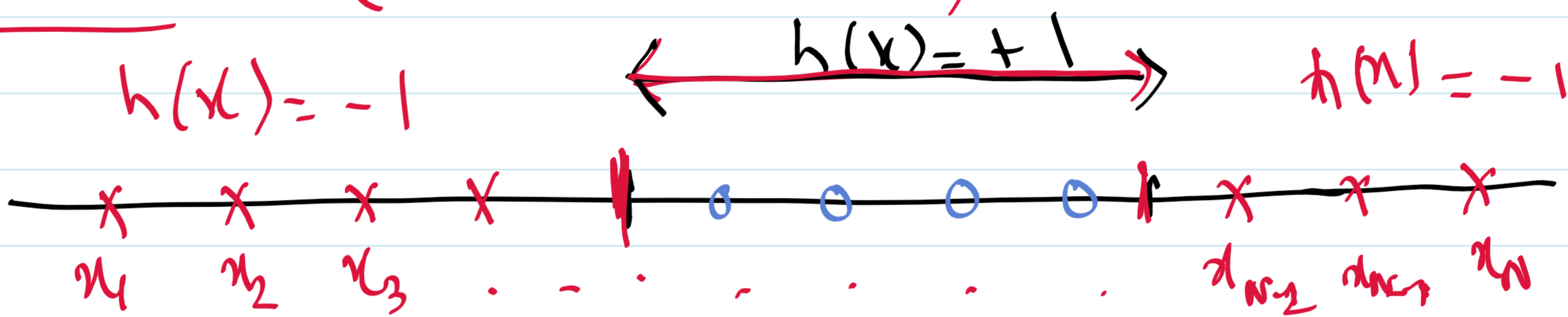


Example - 3 (Positive rays)



- \mathcal{H} consists of $h : \mathbb{R} \rightarrow \{+1, -1\}$.
- $h(x) = \text{Sign}(x - \alpha)$
- cut the line into two halves.
- $m_{\mathcal{H}}(N) = \mathbf{N+1}$

Example - 4 (Positive intervals)



$$m_{\mathcal{H}}(N) = \binom{N+1}{2} + 1 = \frac{N^2}{2} + \frac{N}{2} + 1$$