

1. (20 points) Assume that we have M slot machines in front of us. Each machine has an unknown probability of μ_m for returning one coin, and a probability of $1 - \mu_m$ for returning no coin. For each of the time step $t = 1, 2, \dots$, assume that we pull the machine $m = ((t-1) \bmod M) + 1$. After some $t > M$ time steps, we'd have pulled machine m for N_m times, and collected c_m coins from machine m . Note that $N_m \geq 1$ because $t > M$. Consider the following one-sided Hölding's inequality (which is slightly different from what we taught in class)

$$P(\mu > x + \epsilon) \leq \exp(-2\epsilon^2 N),$$

where ν, μ, ϵ, N have been defined in our class. Use the inequality above to prove that when given a fixed machine m and a fixed δ with $0 < \delta < 1$,

$$P\left(\hat{\mu}_m > \frac{c_m}{N_m} + \sqrt{\frac{\log(1 - \frac{1}{2}\log\delta)}{N_m}}\right) \leq \delta t^{-2},$$

$\nu = \frac{c_m}{N_m}$

$\mu = \frac{c_m}{N_m} + \epsilon$

$$\nu = \frac{c_m}{N_m}$$

$$\begin{aligned} \epsilon &= \sqrt{\frac{\log(1 - \frac{1}{2}\log\delta)}{N_m}}, \quad N = N_m, \quad \epsilon' = \sqrt{\frac{\log(1 - \frac{1}{2}\log\delta)}{N}} \\ \epsilon' &= \frac{\sqrt{\log(1 - \frac{1}{2}\log\delta)}}{\sqrt{N}} \\ N\epsilon'^2 &= (\log(1 - \frac{1}{2}\log\delta)) \\ -2N\epsilon'^2 &= -2\log(1 - \frac{1}{2}\log\delta) \\ e^{-2N\epsilon'^2} &= t^{-2} \cdot \delta \end{aligned}$$

$$\rightarrow P(M > \nu - \epsilon) \leq e^{-2N\epsilon'^2}$$

$$\rightarrow P(M_m > \frac{c_m}{N_m} + \sqrt{\frac{\log(1 - \frac{1}{2}\log\delta)}{N_m}}) \leq \delta t^{-2}$$

2. (20 points) Continuing from Problem 1, prove that when $M \geq 2$, for all slot machines $m = 1, 2, \dots, M$ and for all $t = M+1, M+2, \dots$, with probability at least $1 - \delta$,

$$\mu_m \leq \frac{c_m}{N_m} + \sqrt{\frac{\log t + \log M - \frac{1}{2} \log \delta}{N_m}}.$$

You can use the magical fact that

$$\sum_{t=1}^{\infty} t^{-2} = \frac{\pi^2}{6}.$$

Hint: The fact that we can upper-bound all μ_m confidently and simultaneously by $\frac{c_m}{N_m}$, plus a deviation term is the core technique for deriving the so-called upper-confidence bound algorithm for

multi-armed bandits, which is an important algorithm for the task of online and reinforcement learning. The actual algorithm differs from what we do in Problem 1 by pulling the machine with the largest upper confidence bound in each iteration, instead of periodically going through each machine. Those who are interested can certainly search for more about this.

$$\varepsilon = \sqrt{\frac{\log t + \log M - \frac{1}{2} \log \delta}{N_m}}$$

$$N_m \varepsilon^2 = \log t + \log M - \frac{1}{2} \log \delta$$

$$\exp(N_m \varepsilon^2) = M^t \cdot \delta^{\frac{1}{2}}$$

$$\exp(-2N_m \varepsilon^2) = \delta^{t^{-2} M^2}$$

$$P(M_m \leq \frac{c_m}{N_m} + \sqrt{\frac{\log t - \frac{1}{2} \log \delta}{N_m}}) = 1 - P(M_m > \frac{c_m}{N_m} + \sqrt{\frac{\log t - \frac{1}{2} \log \delta}{N_m}}) \geq 1 - \delta t^{-2}$$

$$P(M_m \leq \frac{c_m}{N_m} + \sqrt{\frac{\log t - \frac{1}{2} \log \delta + \log M}{N_m}}) \geq 1 - \delta t^{-2} M^{-2}, \quad M \geq 2$$

$$\geq 1 - \delta \cdot t^{-2} \cdot \frac{1}{4} \quad , \quad \sum t^{-2} = \frac{\pi^2}{6}$$

$$\geq 1 - \delta \cdot \frac{\pi^2}{24}$$

$$> 1 - \delta \quad \#$$

3. (20 points) Next, we illustrate what happens with multiple bins. Consider a special lottery game as follows. The game operates by having four kinds of lottery tickets placed in a big black bag, each kind with the same (super large) quantity. Exactly eight numbers 1, 2, ..., 8 are written on each ticket. The four kinds are

- A: all even numbers are colored orange, all odd numbers are colored green
- B: all even numbers are colored green, all odd numbers are colored orange
- C: all small numbers (1-4) are colored orange, all big numbers (5-8) are colored green
- D: all small numbers (1-4) are colored green, all big numbers (5-8) are colored orange

Every person is expected to draw five tickets from the bag. A small prize of 1450 is given if the five tickets contain "some number" that is purely green. What is the probability that such an event will happen?

$$\text{if } (1,2,3,4) \text{ green} \rightarrow \text{all C} \quad ; \quad C_4^4 = 1$$

$$(5,6,7,8) \text{ green} \rightarrow \text{all D} \quad ; \quad C_4^4 = 1$$

$$(1,2,5,7) \text{ green} \rightarrow \text{all A} \quad ; \quad C_4^4 = 1$$

$$(3,4,6,8) \text{ green} \rightarrow \text{all B} \quad ; \quad C_4^4 = 1$$

$$(5,7,8) \text{ green} \rightarrow \text{ACCCC, AACCC, AAACC, AARRAC} \quad C_3^3 + C_3^2 + C_3^1 + C_3^0 = 5+5+5+\frac{5 \cdot 4}{2} \times 2 = 30$$

$$(1,3) \text{ green} \rightarrow \text{AD arrange} \quad , \quad \text{the same above} = 30$$

$$(2,4) \text{ green} \rightarrow \text{BD arrange} \quad , \quad \text{the same above} = 30$$

$$(6,8) \text{ green} \rightarrow \text{BC arrange} \quad , \quad \text{the same above} = 30$$

The probability of "some number are all green" should be $\frac{12+30+4}{4^5} = \frac{84}{1024} = \frac{21}{256} \approx 8\%$

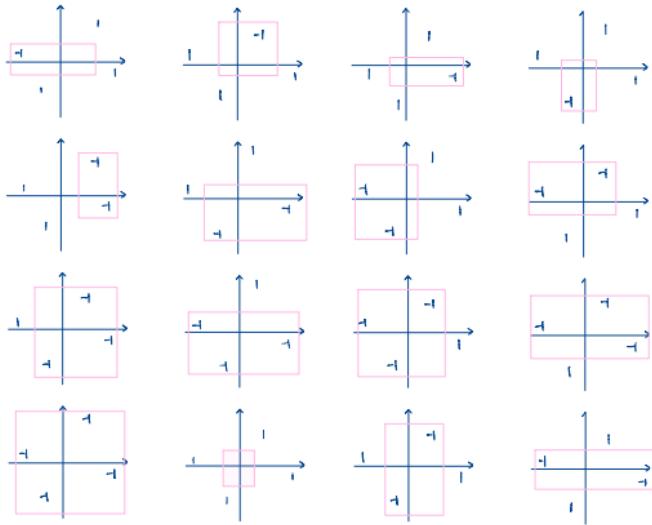
4. (20 points) Continuing from Problem 3, a bigger prize of three piggy banks will be given if the five tickets contain five green 2's. What is the probability that such an event will happen?

Hint: Each number can be viewed as a "hypothesis" and the drawn tickets can be viewed as the data. The E_{out} of each hypothesis is simply $\frac{1}{2}$ (You are welcome, :-)). Problem 4 asks you to calculate the BAD probability for hypothesis 2; Problem 3 asks you to calculate the BAD probability for all hypotheses, taking the dependence into consideration.

$$P[\text{five green 2's}] = \frac{\binom{1(2,5,7)}{1} \binom{2(4,6,8)}{1} \binom{2(4)}{1}}{4^5} = \frac{1+1+30}{4^5} = \frac{32}{4^5} = \frac{1}{32} =$$

5. (20 points) Consider the "negative rectangle" hypothesis set for $X = \mathbb{R}^2$, which includes any hypothesis that returns -1 when \mathbf{x} is within an axis-parallel rectangle and $+1$ elsewhere. Show that some set of 4 input vectors can be shattered by the hypothesis set. That is, the VC dimension of the hypothesis set is no less than 4.

$\hookrightarrow 2^4 = 16$



Result: $\text{mn}(4) = 2^4$ for 4 inputs,
that is to say : 4 inputs can be shattered,
dvc isn't less than 4.

6. (20 points) Consider a hypothesis set \mathcal{H} for $X = \mathbb{R}$ containing hypothesis with $2M+1$ ($M \geq 1$) parameters. Each hypothesis $h(x)$ in \mathcal{H} are defined by $s, a_1, b_1, a_2, b_2, \dots, a_M, b_M$ that satisfies

- $s \in \{+1, -1\}$
- $a_m < b_m$, for $1 \leq m \leq M$;
- $b_m < a_{m+1}$, for $1 \leq m \leq M-1$,

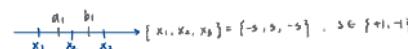
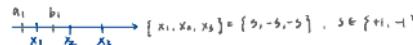
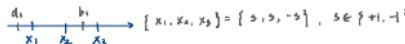
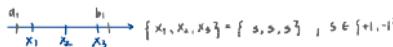
with

$$h_{s,a,b}(x) = \begin{cases} s, & \text{if } a_m \leq x \leq b_m \text{ for some } 1 \leq m \leq M \\ -s, & \text{otherwise} \end{cases}$$

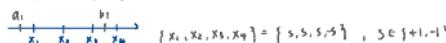
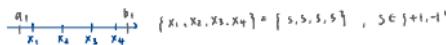
What is the VC dimension of \mathcal{H} ? Prove your answer.

Hint: The positive intervals introduced in Lecture 5 correspond to $s=+1$ with $M=1$.

For $N=3, M=1, m_H(3)=8$, shattered.



For $N=4, M=1, m_H(4)=14$, can't shattered.



If h can shattered N points $\rightarrow m_H(N)=2^N$, every x_i for $0 \leq i < N$ has two choices $\{s=+1\}$ or $\{s=-1\}$

That is, if h can shattered some N points $\rightarrow h$ has $(N-1)$ parameters to build N intervals

such that x_1, x_2, \dots, x_N won't be forced to be in the same intervals.

x_i has the freedom to choose $\{s=+1\}$ or $\{s=-1\}$

For $N=5, M=2, m_H(5)=5$, shattered



using 4 parameters, it could shattered 5 inputs.

The result can also refer to the conclusion in class:

parameters creates degree of freedom $\rightarrow \text{dvc} \approx \text{free parameters}$.

In conclusion, hypothesis with $(2M+1)$ parameters



$(2M+1)$ intervals.

\rightarrow at most can shattered $2M+1$ data

$\rightarrow \text{dvc} = 2M+1 \approx$

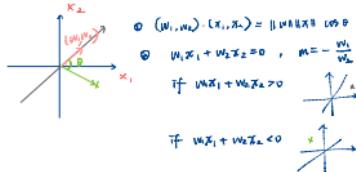
7. (20 points) What is the growth function of origin-passing perceptrons on $\mathcal{X} = \mathbb{R}^2$? Those perceptrons are

$$\mathcal{H}_0 = \{h: h(\mathbf{x}) = \text{sign}(w_1x_1 + w_2x_2) \text{ i.e. perceptrons that pass the origin}\}$$

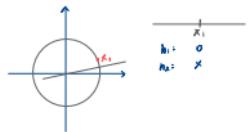
Prove your answer.

Hint: Consider putting your input vectors on the unit circle.

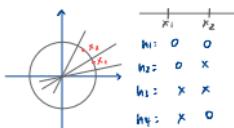
w_1, w_2 : weight
 x_1, x_2 : variables (two characteristics)



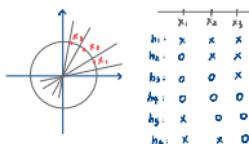
when $N=1$, $M_H(1)=2$



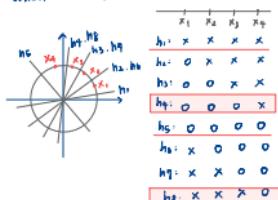
when $N=2$, $M_H(2)=4$



when $N=3$, $M_H(3)=6$



when $N=4$, $M_H(4)=8 = M_H(3)+2$



$$N+1, M_H(N)=M_H(N-1)+2$$

x_9 can be seen as the point following $x_8 \rightarrow \text{sign}(x_9) = \text{sign}(x_8)$, h_n for $n=1, 2, \dots, b$

when $h_9, h_{10}, \text{ sign}(x_9) = -\text{sign}(x_8)$,
which will create two more situation
 $\rightarrow M_H(5)=2+2=M_H(4)$

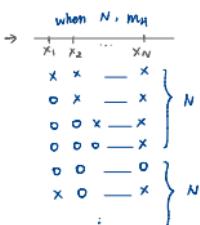
$$\rightarrow M_H(5)=2+2=M_H(4)$$

Using mathematical induction, $M_H(3)=6 \cdot M_H(4)=8$

Assume that $M_H(N)=2N$

$$M_H(N+1)=M_H(N)+2$$

$$2(N+1)=2N+2, M_H(N)=2N \text{ ok.}$$



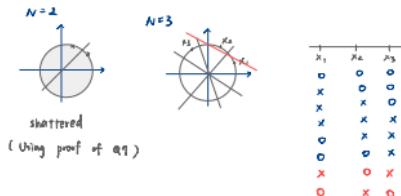
8. (20 points) For $\mathcal{X} = \mathbb{R}^2$, consider a hypothesis set $\mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1$ that is a union of two types of perceptrons:

$$\mathcal{H}_0 = \{h: h(\mathbf{x}) = \text{sign}(w_1x_1 + w_2x_2) \text{ i.e. perceptrons that pass the origin}\}$$

$$\mathcal{H}_1 = \{h: h(\mathbf{x}) = \text{sign}(w_1(x_1 - 1) + w_2(x_2 - 1)) \text{ i.e. perceptrons that pass } (1, 1)\}$$

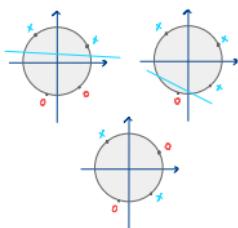
What is the VC dimension of \mathcal{H} ? Prove your answer.

$\text{dvc} \rightarrow \text{find the min number of } \mathbf{x} \text{ that could be shattered.}$



shattered by $\mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1$

Next, we assume $\mathcal{H}' = \mathcal{H}_0' \cup \mathcal{H}_1'$ with no limitation on passing $(0,0)$ and $(1,1)$
in other words, \mathcal{H}' (looser) is the upper bound of \mathcal{H} (limited)



$$\mathcal{H}_0'(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = \{0000, x000, xx00, xxx0, 0xxx, 00xx, 000x, xxxx\}$$

$$\mathcal{H}_1'(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = \{x00x, ox00, xx0x, ox0x\}$$

There aren't any \mathcal{H} such that $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\} = \{0x00, x00x\}$

when $N=4$, $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\}$ couldn't be shattered by \mathcal{H}'

\rightarrow the breakpoint k of \mathcal{H}' , $k=4$

Since \mathcal{H}' is the upper bound of $\mathcal{H} \rightarrow \mathcal{H}$ couldn't shatter $N=4$

$$\rightarrow \text{dvc} = 3 \pm$$

9. (20 points) In class, we taught about the learning model of "positive and negative rays" (which is simply one-dimensional perceptron) for one-dimensional data. The model contains hypotheses of the form:

$$h_{s,\theta}(x) = s \cdot \text{sign}(x - \theta)$$

You can take $\text{sign}(0) = -1$ for simplicity but it should not matter much for the following problems. The model is frequently named the "decision stump" model and is one of the simplest learning models. As shown in class, for one-dimensional data, the VC dimension of the decision stump model is 2.

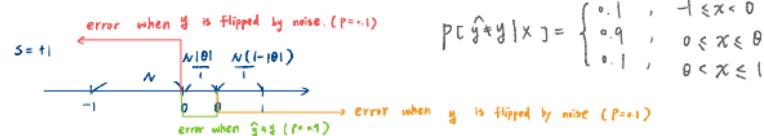
In the following problems, you are asked to play with decision stumps on an artificial data set. First, start by generating a one-dimensional data by the procedure below:

a) Generate x by a uniform distribution in $[-1, 1]$.

b) Generate y by $y = \text{sign}(x) + \text{noise}$, where the noise flips the sign with 10% probability.

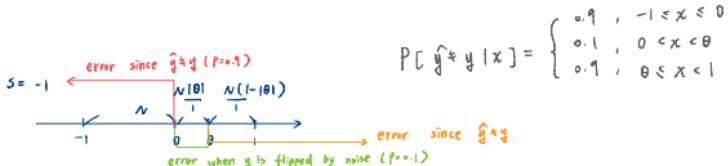
With the (x, y) generation process above, prove that for any $h_{s,\theta}$ with $s \in \{-1, +1\}$ and $\theta \in [-1, 1]$,

$$E_{\text{err}}(h_{s,\theta}) = 0.5 - 0.4s + 0.4s \cdot |\theta|$$

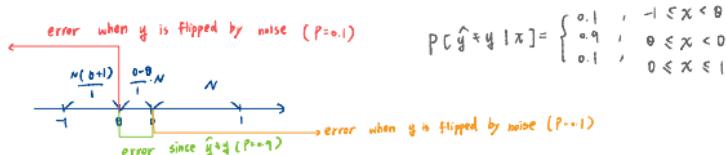


$$E_{\text{err}}(0, \theta) = \frac{s/N + 0.9 \cdot (N|\theta|) + 0.1 \cdot N(1-\theta)}{2N}$$

$$= \frac{1/N + 0.8 N \theta}{2N} = 0.1 + 0.4 \theta$$



$$E_{\text{err}}(0, \theta) = -0.9 N + 0.1 \cdot (N \cdot 0) + 0.9 N (1-\theta) = \frac{1.8 N + 0.8 N \theta}{2N} = 0.9 + 0.4 \theta$$



$$S=1 \wedge \theta > 0$$

$$E_{\text{err}}(0, \theta) = \frac{0.1 N(\theta+1) + 0.9(-\theta) N + 0.1 N}{2N} = \frac{0.2 N + 0.8(-\theta)}{2N} = 0.1 + 0.4(-\theta)$$

$$\Rightarrow E_{\text{err}}(0, \theta) = 0.1 + 0.4 \theta, \theta > 0 \rightarrow E_{\text{err}}(0, \theta) = 0.5 + 0.4 \theta + 0.4 |\theta|$$

$$E_{\text{err}}(0, -\theta) = 0.9 + 0.4 \theta, \theta > 0$$

$$E_{\text{err}}(0, +\theta) = 0.1 + 0.4(-\theta), \theta < 0$$

By Hoeffding's inequality, for any fixed $h_{s,\theta}(x)$, if N is large enough

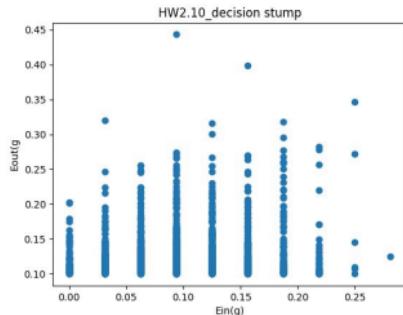
$$E_{\text{err}}(n) \approx E_{\text{err}}(N)$$

$$\text{in conclusion, } E_{\text{err}}(h, \theta) = 0.5 + 0.4s + 0.4|\theta| \neq$$

10. (20 points, *) In fact, the decision stump model is one of the few models that we could minimize E_{in} efficiently by enumerating all possible thresholds. In particular, for N examples, there are at most $2N$ dichotomies (see the slides for positive rays), and thus at most $2N$ different E_{in} values. We can then easily choose the hypothesis that leads to the lowest E_{in} by the following decision stump learning algorithm.

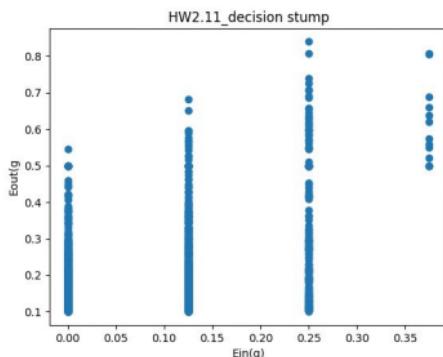
- (1) sort all N examples x_n to a sorted sequence x'_1, x'_2, \dots, x'_N such that $x'_1 \leq x'_2 \leq x'_3 \leq \dots \leq x'_N$
 - (2) for each $\theta \in \{-1\} \cup \{\frac{x'_i + x'_{i+1}}{2} : 1 \leq i \leq N-1 \text{ and } x'_i \neq x'_{i+1}\}$ and $s \in \{-1, +1\}$, calculate $E_{\text{in}}(h_{s,\theta})$
 - (3) return the $h_{s,\theta}$ with the minimum E_{in} as g ; if multiple hypotheses reach the minimum E_{in} , return the one with the smallest $s \cdot \theta$.
- (Hint: CS-majored students are encouraged to think about whether the second step can be carried out efficiently, i.e. $O(N)$, using dxxxxc pxxxxxxxxg instead of the naive implementation of $O(N^2)$.)

Generate a data set of size 32 by the procedure above and run the one-dimensional decision stump algorithm on the data set to get g . Record $E_{\text{in}}(g)$ and compute $E_{\text{out}}(g)$ with the formula in Problem 9. Repeat the experiment 2000 times. Plot a scatter plot of $(E_{\text{in}}(g), E_{\text{out}}(g))$, and calculate the median of $E_{\text{out}}(g) - E_{\text{in}}(g)$.



Median = 0.0383

11. (20 points, *) Repeat Problem 10, but generate a data set of size 8 by the procedure instead. Plot a scatter plot of $(E_{in}(g), E_{out}(g))$, and calculate the median of $E_{out}(g) - E_{in}(g)$. Compare the scatter plot and the median value with those of Problem 10. Describe your findings.



Median = 0.1240

The median of data set = 6 is smaller than the median of data set = 32.

The consequence is similar to Hoeffding's inequality, which indicates that when N become larger (N= the number of data), Ein will be closer to Eout.

Focusing on the right side of the plot, when $Ein > 0.35$, $Eout > 0.4$; to sum up, if Ein is large, $Eout$ is highly possible to be large.

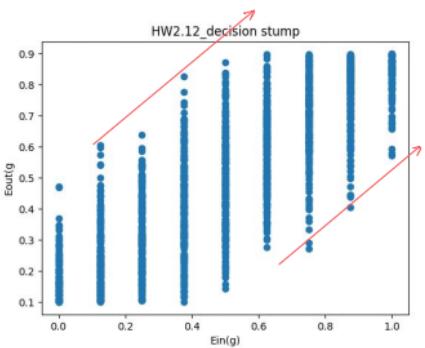
In most of cases, Ein is slightly smaller than $Eout$, which also explain the meaning if Both.

Ein = the mean of hypothesis makes error on known data;

$Eout$ = the expectation of hypothesis making error on unknown data.

It is more difficult for hypothesis to predict unknown data.

12. (20 points, *) Repeat Problem 11, generate a data set of size 8 by the procedure above. Instead of running the decision stump algorithm, return a randomly chosen $h_{s,\theta}$ as g , with s uniformly sampled from $\{-1, +1\}$ and θ uniformly sampled from $[-1, 1]$. Record $E_{in}(g)$ and compute $E_{out}(g)$ with formula in Problem 9. Repeat the experiment 2000 times. Plot a scatter plot of $(E_{in}(g), E_{out}(g))$, and calculate the median of $E_{out}(g) - E_{in}(g)$. Compare the scatter plot and the median value with those of Problem 11. Describe your findings.



$$\text{Median} = 0.0$$

Which is much lower than 2.10, 2.11

In this case, E_{out} has positive correlation to E_{in}

$$E_{in} = 0, 0 < E_{out} < 0.5$$

$$E_{in} = 1.0, 0.5 < E_{out} < 1.0$$

13. (Bonus 20 points) Consider \mathcal{H} being perceptrons in $\mathcal{X} = \mathbb{R}^d$. It is known, by the so-called Cover's Theorem, that the growth function is

$$m_{\mathcal{H}}(N) = 2 \sum_{i=0}^d \binom{N-1}{i}.$$

See, for instance,

https://www.mit.edu/course/other/12courses/www/vision_and_learning/perceptron_notes.pdf

for its proof.

Now, assume that we require the perceptrons to pass *all* k anchor points for a_1, a_2, \dots, a_k , each being in \mathbb{R}^d with $0 \leq k < d$. We shall call those perceptrons \mathcal{H}' . What is the growth function $m_{\mathcal{H}'}(N)$? Prove your answer.

Note: Problem 7 is a special case for $k = 1$ and $a_1 = 0$.

- ① The requirement to pass k anchor points imposes k constraints on \mathcal{H}' , which reduces the number of "free parameters", thereby reducing its capacity to shatter sets.
 \rightarrow Each k anchor points makes $dvc = d \rightarrow dvc = d-k$
- ② $m_{\mathcal{H}'}(N) =$ the number of ways to choose subsets of those N points by using \mathcal{H}'
 for each subset of size i ($0 \leq i \leq d-k$)
 if we want to choose subset from $N-1$ points $\rightarrow \binom{N-1}{i}$ ways.

$$\Rightarrow m_{\mathcal{H}'}(N) = 2 \sum_{i=0}^{d-k} \binom{N-1}{i}$$

