

1a.) To get an algorithm $A(S_{\text{train}})$ w/ 0% training error, we our algorithm will look through the training set to find all the labels. There are 2 possible cases: our labels contains an x with a label of 1 or our labels do not contain an x w/ label of 1. (they are all 0.) Then, if there is a label 1, we choose $h^* = h_2$ where z is equal to the x value of the point with label 1. Otherwise, we choose $h^* = h_1$, the hypothesis with all x -input labeled as 0. Either way, our training error will be 0.

$$1b) P(L_{(0, n^*)}(A(S_{\text{train}})) > \varepsilon) \leq 0$$

~~loss~~ on training data

Since S_{train} is a random variable, loss is also a random variable

$$P(L_{(0, n^*)}(A(S_{\text{train}})) > \varepsilon)$$

$L_{(0, n^*)}(A(S_{\text{train}})) > 0$ is a larger event than $L_{(0, n^*)}(A(S_{\text{train}})) > \varepsilon$ b/c $0 \geq \varepsilon$.

so we will look to bound $P(L_{(0, n^*)}(A(S_{\text{train}})) > 0)$

$$\text{want: } 0 \geq P(L_{(0, n^*)}(A(S_{\text{train}})) > 0) \geq P(L_{(0, n^*)}(A(S_{\text{train}})) > \varepsilon)$$

Each point in the S_{train} has a $1 - \frac{1}{n}$ chance of being the incorrect value ($x \neq z$). Thus for m points, the probability that our ~~h*~~ $A(S_{\text{train}})$ is NOT the real hypothesis h^* is $(1 - \frac{1}{n})^m$

$$\text{so } \sigma \geq (1 - \frac{1}{n})^m \geq P(L(A(S_{\text{train}})) > \varepsilon)$$

and we know that $P_{x \sim D}(h^*(x) \neq A(S_{\text{train}})(x)) \geq \frac{1}{N}$

because for two different hypotheses, the probability they are different is equal to the probability $x=z$ where $h(z)=1$

$$\text{so } P_{x \sim D}(h^*(x) \neq A(S_{\text{train}})(x)) = \frac{1}{N}$$

If we plug in $m = \frac{\log(1/\sigma)}{\varepsilon}$

$$\sigma \geq (1 - \frac{1}{n}) \frac{\log(1/\sigma)}{\varepsilon} \rightarrow \varepsilon \geq -\log(1 - \frac{1}{n})$$

$$\log \sigma \geq \frac{\log(1/\sigma)}{\varepsilon} \log(1 - \frac{1}{n})$$

if $\left(\frac{1}{n} > \varepsilon\right)$ then $\frac{1}{n} \geq -\log(1 - \frac{1}{n})$ and (generalization error $> \varepsilon$) $\frac{1}{n} \geq \log(1 - \frac{1}{n}) \rightarrow$

$$-\frac{1}{N} \leq \log(1 - \frac{1}{N})$$

$$\text{and } e^{-\frac{1}{N}} \geq 1 - \frac{1}{N}$$

using the given equation,
we know $1-x \leq e^{-x}$

therefore, using $x = \frac{1}{N}$,
we know $1 - \frac{1}{N} \leq e^{-\frac{1}{N}}$,

and thus we have bounded
 $(1 - \frac{1}{N})^M \leq 0$,

$$2.) H = \{ \text{sgn}(ax^2 + bx + c); a, b, c \in \mathbb{R} \}$$

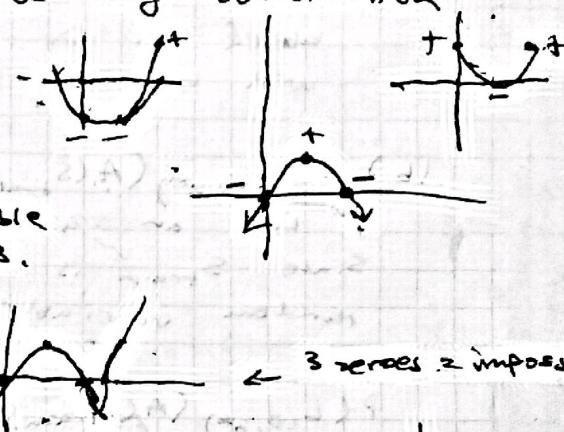
given points $(0,0), (1,1), (2,0) \quad x=0, 1, 2$

we can construct a polynomial function of degree 2 such that the sign at these 3 points can be any combination of + and -.

$$\text{Thus } VC(H) \geq 3$$

Given 4 pts, $x=0, 1, 2, 3$ a combination of $(0,0), (1,1), (2,0), (3,1)$ is impossible and requires a polynomial of degree 3.
Thus $VC(H) \leq 4$

$$\text{Therefore } VC(H) = 3$$



i	Label	Hypothesis 1			Hypothesis 2		
		D_0	$f_1(x > 2, y > 6)$	$f_2(x > 2, y > 6)$	D_1	$f_3(x > 10, y > 11)$	$f_4(x > 10, y > 11)$
1	-	0.1	-	+x	-	0.1	-
2	-	.1	-	-	-	0.1	-
3	+	.1	+	+x	+	0.1	-
4	-	.1	-	-	-	0.1	-
5	-	.1	=	+x	-	0.1	-
6	-	.1	+x	+x	+	0.1	-
7	+	.1	+x	+x	-	0.1	-
8	-	.1	-	-	-	0.1	-
9	+	.1	-x	+	-	0.1	-
10	+	.1	+x	+x	+	0.1	-

← 3 zeroes is impossible



$$3c) \alpha_0 = \frac{1}{2} \log_2 \left(\frac{1-2}{2} \right) = 1$$

$$D_0 = \frac{1}{1024} \cdot \sum_{i=1}^{1024} 2^{-1} \\ = \begin{cases} \frac{0.25}{2^4} = 0.0625 \\ \frac{0.25}{2^4} = 0.25 \end{cases}$$

$$2 \left(\frac{1}{2^4} \right) + 8 \left(\frac{0.25}{2^4} \right) = 1$$

$$y_i = h_+(x_i)$$

otherwise

$$\begin{aligned} x &\approx 10.1 \\ \epsilon &\approx 0.25 + 2 \times 0.0625 \\ &= 0.375 \end{aligned}$$

$$\begin{aligned} y &\approx 10.1 \\ \epsilon &\approx 4 \times 0.0625 = 0.25 \end{aligned}$$

$$\alpha_1 = \frac{1}{2} \log_2 \left(\frac{1-0.25}{0.25} \right) = 0.792$$

$$3d) H(x) = \text{sgn}(1 \times [x > 2] + 0.79 \times [y > 11])$$

- * 4a.i) One v All: Learn one classifier for each i , so K classifiers.
 All v All: Learn one classifier for each $\langle i, j \rangle$ pair.
 So for $i \neq j$, learn $K-1$ classifiers. Thus G learns $\frac{K(K-1)}{2}$ classifiers. $C(K, 2) = \frac{K(K-1)}{2}$ classifiers.
- 4a.ii) One v All: We use all m examples for each classifier.
 All v All: We use $\frac{2m}{k}$ examples for each classifier, $\frac{m}{k}$ from each class.
- 4a.iii) One v All: For each classifier, determine ~~if it's more likely~~ probability it is that class label v not that label.
 The one with the highest probability is the final choice.
 All vs All: For each classifier, determine the "winner".
 Take these winners as votes, and determine label with the most votes.
- 4a.iv) Using One v All: train K classifiers, each using m examples: $O(mk)$
 All v All: $\frac{K(K-1)}{2}$ classifiers. $\frac{2m}{k}$ examples = $m(k-1)$
 $= O(mk)$
- 4b.) I prefer one vs all because of its easier implementation and it is more computationally efficient since we are only training K classifiers.
- 4c.) Yes, using a kernel Perceptron does change my above analysis; since Kernel Perceptron has $O(m^2)$ time
 One v All becomes: $O(m^2k)$
 All v All becomes: $O(m^2)$
 Thus, I would use All v All since it is faster now.
- 4d.) Using the black box to train each classifier, a classifier using one vs all would have $O(dm^2)$.
 With K classifiers, this becomes $O(Kdm^2)$
 A classifier using all v all would have $O(d(\frac{2m}{k})^2)$
 With $\frac{K(K-1)}{2}$ classifiers, this becomes $O(dm^2)$
 Thus all v all becomes more efficient.
- 4e.) One vs all becomes: $O(Kd^2m)$
 All vs all becomes: $O(Kd^2m)$
 The time complexities are the same, so they both have the same efficiency.

4f.) Counting: we run the prediction for each classifier:
This takes $O(k^2)$ time

Knockout: Since there are k classes and we eliminate a class each round, the time complexity is $O(k)$.