

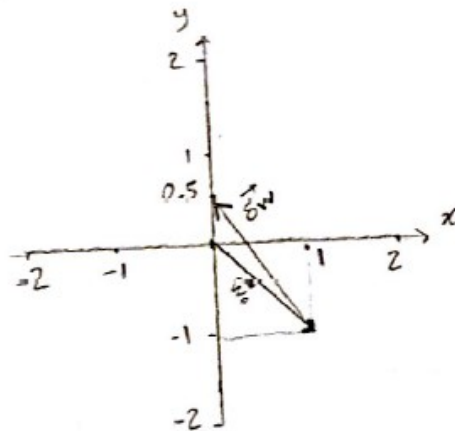
Q1: a) $\vec{w}_{t+1} = \vec{w}_t + y_n \vec{x}_n$ (if $\text{sign}(\vec{w}^T \vec{x}_n) \neq y_n$)

The influence of this update depends on ability of the separable plane to classify (\vec{x}_n) correctly before the update. If it already did a good job, the weights won't update. But if it misclassified that point, the update considers the new point and it will be better at classifying (\vec{x}_n, y_n)

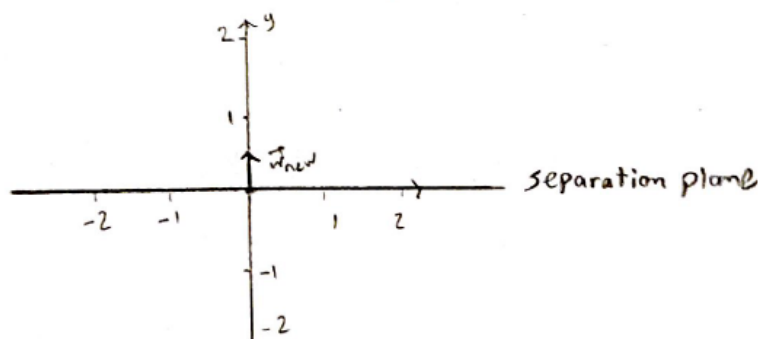
@ t: $\text{sgn}(\vec{w}^T \vec{x}_n) - y_n$
 @ t+1: $\text{sgn}(\vec{w}^T \vec{x}_n + y_n |\vec{x}_n|^2) - y_n$

$\left. \begin{array}{l} \text{if } y_n = -1, w^T x_n > 0 \Rightarrow w^T x_n + y_n |\vec{x}_n|^2 < w^T x_n \\ \text{if } y_n = +1, w^T x_n < 0 \Rightarrow w^T x_n + y_n |\vec{x}_n|^2 > w^T x_n \end{array} \right\} \begin{array}{l} \text{argument} \\ \text{gets} \\ \text{better} \\ \text{probably} \end{array}$

b) i - $w^T \begin{bmatrix} -1 \\ 1.5 \end{bmatrix} = -1 - 1.5 = -2.5 \Rightarrow \text{sgn}(-2.5) \neq y = 1 \rightarrow \delta \vec{w} = \eta \times 1 \times \begin{bmatrix} -1 \\ 1.5 \end{bmatrix} \xrightarrow{\eta=1} \begin{bmatrix} -1 \\ 1.5 \end{bmatrix}$

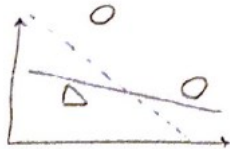


ii - $\vec{w}_{new} = \vec{w} + \delta \vec{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} -1 \\ 1.5 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}$



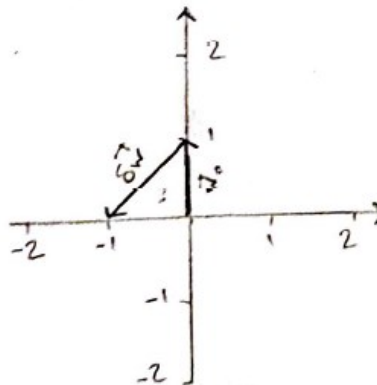
- iii - Use $[1, 0.5] \rightarrow w^T x = 0.25 \rightarrow \text{sgn}(w^T x) = 1 = y \rightarrow$ doesn't update
 Use $[-1, -0.5] \rightarrow w^T x = -0.25 \rightarrow \text{sgn}(w^T x) = -1 = y \rightarrow$ doesn't update
 Since the points are linearly separable, the updates will converge to a state that all points would be correctly classified (as in this case).

c) NO. There are an infinite number of lines that can separate these dots. Amongst them, the line that maintains a large threshold is better than others, and can generalize better.



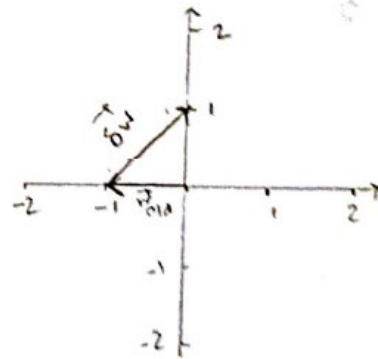
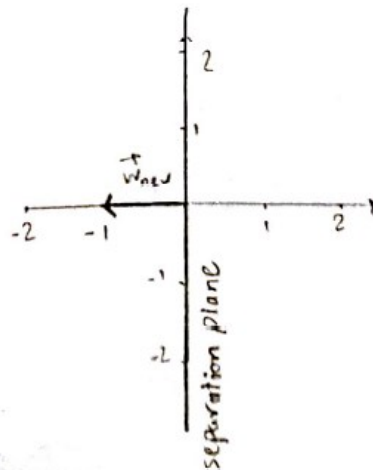
→ The dotted line is better than the solid one.

d) i- $\text{sgn}(w^T x) = \text{sgn}(-1) = -1 \neq y = 1 \Rightarrow \delta w = y_n x_n = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$



ii - $\vec{w}_{new} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$

Now use $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow \text{sgn}(w^T x) = \text{sgn}(-1) = -1 \neq y = 1 \rightarrow \delta \vec{w} = yx = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow \vec{w}_{new} = \begin{bmatrix} -1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$



iii - NO, for example $\text{sgn}(w_{new}^T \begin{bmatrix} 1 \\ -1 \end{bmatrix}) = -1 \neq y = 1$

After two updates, it is same as the initial weights (w_0).

As the datapoints are not linearly separable, the algorithm won't converge.

Q2

a) Assume that $\vec{w}^1 = 0 \Rightarrow n=0 \rightarrow \vec{w}^{n+1} \cdot \vec{w}^* = \vec{w}^1 \cdot \vec{w}^* = 0$

$$n = k \rightarrow \vec{w}^{k+1} \cdot \vec{w}^* = (\vec{w}^k + y_n \vec{x}_n^T) \cdot \vec{w}^* \\ = \vec{w}^k \cdot \vec{w}^* + y_n \vec{x}_n^T \vec{w}^*$$

From linear separability, we know that $y_n \vec{x}_n^T \vec{w}^* > \gamma$; so $\Rightarrow \vec{w}^{k+1} \cdot \vec{w}^* > \vec{w}^k \cdot \vec{w}^* + \gamma$
 $\Rightarrow \vec{w}^2 \cdot \vec{w}^* > \gamma, \vec{w}^3 \cdot \vec{w}^* > 2\gamma, \dots, \boxed{\vec{w}^{k+1} \cdot \vec{w}^* > k\gamma}$

b) $\vec{w}^{k+1} \cdot \vec{w}^* = \|\vec{w}^{k+1}\| \|\vec{w}^*\| \cos \theta \leq \|\vec{w}^{k+1}\| \|\vec{w}^*\| \Rightarrow \boxed{\|\vec{w}^{k+1}\| > k\gamma} \quad \textcircled{I}$

c) $\|\vec{w}^{k+1}\|^2 = \|\vec{w}^k + y_n \vec{x}_n\|^2 \leq \|\vec{w}^k\|^2 + \|y_n \vec{x}_n\|^2 = \|\vec{w}^k\|^2 + \|\vec{x}_n\|^2 R^2$
 $\Rightarrow \|\vec{w}^1\|^2 \leq R^2, \|\vec{w}^3\|^2 \leq R^2 + R^2 = 2R^2, \dots, \boxed{\|\vec{w}^{k+1}\|^2 \leq kR^2} \quad \textcircled{II}$

d) $\textcircled{I}, \textcircled{II} \rightarrow (k\gamma)^2 < \|\vec{w}^{k+1}\|^2 \leq kR^2 \Rightarrow kR^2 > k^2\gamma^2 \Rightarrow k < \frac{R^2}{\gamma^2}$

Q3: Flipping a coin follows binomial distribution. $\rightarrow \begin{cases} P(\text{heads}) = \theta \\ P(\text{tails}) = 1 - \theta \end{cases} \Rightarrow P(\theta) = \theta^i (1-\theta)^{10-i}$
 $i \rightarrow \# \text{ heads}$

MLE: $m = \# \text{ observations}$ $\hat{\theta} = \arg \max_{\theta} P(D) = (1-\theta)^{m-i} \theta^i$

$$\Rightarrow \log P(D) = i \log(\theta) + (m-i) \log(1-\theta)$$

$$\Rightarrow \frac{d \log P(D)}{d\theta} = \frac{i}{\theta} + \frac{-m+i}{1-\theta} = 0 \rightarrow i - i\hat{\theta} = m\hat{\theta} - i\hat{\theta} \Rightarrow \hat{\theta} = \frac{i}{m}$$

"log" is a monotonic function

$$\Rightarrow \arg \max_{\theta} P(D) = \arg \max_{\theta} \log P(D)$$

* if observation is the set of ten flips, we can use $m_{\text{new}} = 10m$

Q5: $\nabla_w l(w) = \begin{cases} (x^T w - f(x)) \vec{x} + 2\lambda \vec{w} & |f(x) - f(x)| \leq \delta \\ \delta \operatorname{sgn}(x^T w - f(x)) \vec{x} + 2\lambda \vec{w} & \text{otherwise} \end{cases}$

$$\frac{d \sum w_i^2}{dw} = 2 \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix} = 2 \vec{w}$$

$$w^{k+1} = w^k - \eta \nabla_w l(w) = w^k - \eta \times \begin{cases} (\vec{x}_n^T \vec{w}^k - f(\vec{x}_n)) \vec{x}_n + 2\lambda \vec{w}^k & |f(\vec{x}_n) - f(\vec{x}_n)| \leq \delta \\ \delta \operatorname{sgn}(\vec{x}_n^T \vec{w}^k - f(\vec{x}_n)) \vec{x}_n + 2\lambda \vec{w}^k & \text{otherwise} \end{cases}$$

Q6 :

a) Done. It is uploaded.

b) $0.8478 = 84.78\%$

c) No, it remains as same as before.

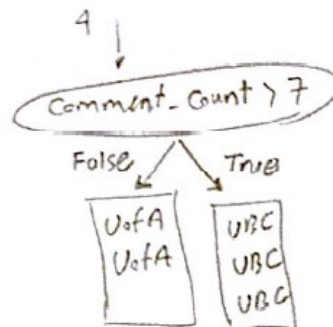
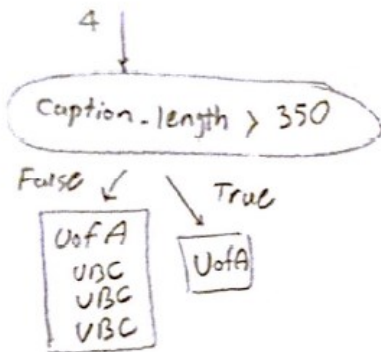
d) Although we change the fractions and possibilities, but what is important for us is the probability that is larger. As the number in the denominator is by far larger than the added values, this won't change the predicted label.

Q7: a) $-\frac{1}{3} \log_2(\frac{1}{3}) - \frac{2}{3} \log_2(\frac{2}{3}) - \frac{0}{3} \log_2(\frac{0}{3}) = 0.918$

b) comment-count is better.
Because there exist a threshold
that can fully separate universities.

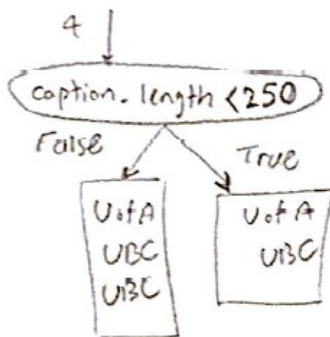
caption-length:	UofA	UBC
	163	320
	397	256
		222

comment-count:	UofA	UBC
	2	8
	0	18
		16



⇒ The nodes are pure.

c)



E(false)

$$-\frac{1}{3} \log_2(\frac{1}{3}) - \frac{2}{3} \log_2(\frac{2}{3}) = 0.918$$

E(True)

$$2 \times (-\frac{1}{2} \log_2(\frac{1}{2})) = 1$$

$$E(\text{node 4}) = -\frac{2}{5} \log_2(\frac{2}{5}) - \frac{3}{5} \log_2(\frac{3}{5}) \approx 0.971$$

$$\Rightarrow 0.971 - 0.918 = 0.053 \rightarrow \text{less than } 0.1$$

⇒ It will be pruned.