# ML homework 1: question 5

Q. What is a possible application of active learning?

The answer generated by chatGPT:

In medical imaging diagnosis, active learning is used to improve the labeling of medical images (e.g., MRI or CT scans) for diagnosing diseases like cancer or identifying abnormalities.

My opinion is:

I agree with ChatGPT's explanation because, firstly, active learning stands out from other protocols due to its unique ability to query only when the algorithm is "confused" or uncertain about a particular instance. This characteristic is especially valuable in applications like medical imaging, where identifying abnormalities (like in CT scans) is critical.

The primary advantage of this protocol is that, it significantly reduces the labeling burden. In many medical applications, labeling data requires the expertise of specialists such as doctors. Hiring professionals with such expertise to label large amounts of data can be both time-consuming and expensive, making it impractical for large datasets. By using active learning, only the most uncertain cases (such as CT scans that are difficult to distinguish from normal ones) need to be verified by the experts.

Secondly, active learning operates in a sequential manner. After the model identifies and queries the most uncertain case, an expert provides the label, and the model is immediately updated. This allows the algorithm to refine its understanding, so similar cases are less likely to be queried again. In contrast to batch learning, which updates the model with a fixed set of labeled data all at once, this sequential approach ensures that the expert is not asked to label similar cases repeatedly, therefore simplifying the process.

# ML homework 1: question 6

Q. Can learning used to predict earthquakes?

Machine learning can be used to predict earthquakes by analyzing seismic data to identify patterns that may precede earthquake events.

## My opinion is:

My answer is yes, but only partially. Recent research has shown that by analyzing the acoustic signals generated in laboratory experiments, machine learning can predict the time remaining until an "earthquake" occurs with high accuracy. Other works usually aimed to focus on predicting specific features of an earthquake, like predicting the magnitude and depth.

While we can gain some insights into earthquakes, it seemed that there is currently no method to generalize these research findings to real-world scenarios. Moreover, there is likely no model capable of accurately predicting when and where an earthquake will happen, and how strong it will be. I think that the challenges are due to reasons like:

1. There are various phenomena that may be signals that represent the occurence of an earthquake, while we can indicate that an earthquake is going to happen directly by some of the phenomena like the p-wave, some of the other phenomena may have ambiguity in whether it can infer an earthquake or not, like a research I saw stated that the sound data they used are initially considered noise.

2. Machine learning models require sufficiently large datasets, if the number of data is too small, then there might be high generalization error, therefore there is difficulty obtaining enough usable data, even though there might be some data from the government, we still have to get data from different regions with varying characteristics in order to ensure that the dataset is comprehensive.

Thus, I would not say that machine learning is entirely ineffective in predicting earthquakes, but there are still many challenges to overcome. As we can see, people are dedicated to this field, so it is not impossible, and perhaps one day, machine learning will have practical applications in earthquake prediction.

1. Let $x_n^{orig} = \begin{bmatrix} x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^d$     for $n \in \{1, ..., N\}$

    then $\overset{\rightharpoonup}{x}_n = \begin{bmatrix} 1 & x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^{d+1}$     for $n \in \{1, ..., N\}$

    and $\overset{\rightharpoonup}{x}_n' = \begin{bmatrix} 2 & x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^{d+1}$     for $n \in \{1, ..., N\}$

Suppose for the the $\overset{\rightharpoonup}{x}_n$'s that are formed conventionally,

    the first $p-1$ inputs are predicted correctly, then:

$$ \text{sign}(\overset{\rightharpoonup}{w_0}{}^T \overset{\rightharpoonup}{x}_i) = y_i \quad \text{for } i = 0 \sim p-1 $$

If we assume $\text{sign}(0) = -1$ for simplicity, and $\overset{\rightharpoonup}{w_0} = \overset{\rightharpoonup}{0}$, then this means that

$$ \text{sign}(\overset{\rightharpoonup}{0}{}^T \overset{\rightharpoonup}{x}_i) = \text{sign}(0) = -1 = y_i \quad \text{for } i = 1 \sim p-1 $$

If we further assume $\overset{\rightharpoonup}{x}_p$ receives incorrect prediction, then

$$ \text{sign}(\overset{\rightharpoonup}{0}{}^T \overset{\rightharpoonup}{x}_p) \neq y_p $$

$$ \Rightarrow \quad \text{sign}(0) = -1 \neq y_p \quad \therefore y_p = 1 $$

Thus we update $\overset{\rightharpoonup}{w}_1$ as :

$$ \overset{\rightharpoonup}{w}_1 = \overset{\rightharpoonup}{w_0} + y_p \overset{\rightharpoonup}{x}_p = \overset{\rightharpoonup}{0} + 1 \cdot \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix} = \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix} $$

Since $w_0^1 = \vec{0}$, and $y_i$ remain unchanged no matter we use $\vec{x_i}$ or $\vec{x_i}^1$, so:

$$\text{sign}(\vec{w_0^1}^T \vec{x_i}^1) = \text{sign}(\vec{0}^T \vec{x_i}^1) = 1 = y_i \quad \text{for } i = 1, ..., p-1$$

We next want to know if $\vec{x_p}^1$ will get correct prediction or not:

$$\text{sign}(\vec{w_0^1}^T \vec{x_p}^1) = \text{sign}(\vec{0}^T \vec{x_p}^1) = 1 \neq 1 = y_p$$

Thus $\vec{x_p}^1$ is also mispredicted, and we update the weight vector as:

$$\vec{w_1^1} = \vec{w_0^1} + y_p \vec{x_p}^1 = \vec{0} + 1 \cdot \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix} = \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix}$$

Consider the next input $\vec{x_{p+1}}$ and $\vec{x_{p+1}}^1$:

$$\boxed{\text{sign}(\vec{w_1}^T \vec{x_{p+1}})} = \text{sign}\left( \begin{bmatrix} 1 & x_{p1} & x_{p2} & \cdots & x_{pd} \end{bmatrix} \begin{bmatrix} x_{p+1, 1} \\ x_{p+1, 2} \\ \vdots \\ x_{p+1, d} \end{bmatrix} \right)$$

$$= \text{sign}\left( 1 + \sum_{i=1}^{d} x_{pi} \, x_{p+1, i} \right)$$

$$\boxed{\text{sign}(\vec{w_1^1}^T \vec{x_{p+1}}^1)} = \text{sign}\left( \begin{bmatrix} y & x_{p1} & x_{p2} & \cdots & x_{pd} \end{bmatrix} \begin{bmatrix} x_{p+1, 1} \\ x_{p+1, 2} \\ \vdots \\ x_{p+1, d} \end{bmatrix} \right)$$

$$= \text{sign}\left( y + \sum_{i=1}^{d} x_{pi} \, x_{p+1, i} \right)$$

Therefore, consider the situation that if $\sum_{i=1}^{d} x_{pi} \, x_{p+1,i} = -3$

then :

$$\text{sign}(\vec{w}_1^{\,T} \vec{x}_{p+1}) = \text{sign}(1 + (-3)) = \text{sign}(-2) = 1$$

$$\text{sign}(\vec{w}_1^{\,'T} \vec{x}_{p+1}) = \text{sign}(4 + (-3)) = \text{sign}(1) = 1$$

$\Rightarrow$ This will cause $\vec{w}_1$ to update if $y_{p+1} = 1$, $\vec{w}_1'$ to update if $y_{p+1} = 1$.

If the case is $y_{p+1} = 1$, and there are only $p+1$ points in the dataset,

$$\left( \text{i.e.} \quad \{\vec{x}_n, y_n\}_{n=1}^{p+1} \,/\, \{\vec{x}_n, y_n\}_{n=1}^{p+1} \right)$$

then $\vec{w}_{p+1} = \vec{w}_1$

and we continue to update $\vec{w}_1'$ and check for $\vec{x}_1'$:

$$\vec{w}_2' = \vec{w}_1' + y_{p+1} \cdot \vec{x}_{p+1}' = \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix} + (-1) \cdot \begin{bmatrix} x_{p+1,1} \\ x_{p+1,2} \\ \vdots \\ x_{p+1,d} \end{bmatrix} = \begin{bmatrix} x_{p1} - x_{p+1,1} \\ x_{p2} - x_{p+1,2} \\ \vdots \\ x_{pd} - x_{p+1,d} \end{bmatrix}$$

$$\text{sign}(\vec{w}_2'^{\,T} \vec{x}_1') = \text{sign}\left( \begin{bmatrix} 0 & x_{p1} - x_{p+1,1} & \cdots & x_{pd} - x_{p+1,d} \end{bmatrix} \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1d} \end{bmatrix} \right)$$

$$= \text{sign}\left( \sum_{i=1}^{d} (x_{pi} - x_{p+1,i}) \, x_{1i} \right)$$

If $\sum_{i=1}^{d} (x_{pi} - x_{pn,i}) x_{1i} \leq 0$, and since $y_1 = -1$ as we assumed earlier,

we have correct prediction using $\vec{w_2}^{1}$.

We can imagine that we construct the inputs to make :

$$\text{sign} (\vec{w_2}^{1T} \vec{x_i}) = y_i \quad \text{for } i = 1, \cdots, p+1$$

Then $\vec{w}_{PLA}^{1} = \vec{w_2}^{1} = \begin{bmatrix} x_{p1} - x_{pn,1} \\ \vdots \\ x_{pn} - x_{pn,d} \end{bmatrix}^{0}$

or we can see it is trivial that each time we update $\vec{w_i}^{1}$,

we add / substract 2 to $w_{i0}$, thus

$$\vec{w}_{PLA}^{1} = \begin{bmatrix} 2k \\ \vdots \end{bmatrix} \quad \text{for } k \in \mathbb{Z}$$

and $\vec{w}_{PLA} = \begin{bmatrix} 1 \\ \vdots \end{bmatrix}$, so $\vec{w}_{PLA}^{1} \neq \vec{w}_{PLA}$. $\square$

$f.$

Let $\overset{\text{orig}}{x_n} = [x_{n_1} \ x_{n_2} \ \cdots \ x_{nd}] \in \mathbb{R}^d$

$\Rightarrow \begin{cases} \vec{x_n} = [1 \ x_{n_1} \ x_{n_2} \ \cdots \ x_{na}] \in \mathbb{R}^{d+1} \\ \\ \vec{x_n}' = [3 \ 3x_{n_1} \ 3x_{n_2} \cdots 3x_{nd}] \in \mathbb{R}^{d+1} \end{cases}$

For arbitrary $(\vec{x_n}, y_n) \in \{(\vec{x_n}, y_n)\}_{n=1}^{N}$, and assume $\text{sign}(0) = 1$

if $\quad \text{sign}(\vec{w_0}^T \vec{x_p}) \neq y_p$, which means $\vec{x_p}$ is the first mispredicted input

then $y_p = 1$

and

$$\vec{w_1} = \vec{w_0} + y_n \vec{x_p} = \vec{x_p}$$

If we further assume the second misprediction occurs at $(\vec{x_q}, y_8)$, then

$\quad \text{sign}(\vec{w_1}^T \vec{x_8}) \neq y_8$

$\Rightarrow \text{sign}(\vec{x_p}^T \vec{x_8}) \neq y_8$

$\Rightarrow \text{sign}\left([1 \ x_{p_1} \ x_{p_2} \cdots \ x_{pd}] \begin{bmatrix} 1 \\ x_{8_1} \\ x_{8_2} \\ \vdots \\ x_{8d} \end{bmatrix}\right) \neq y_8$

$\Rightarrow \text{sign}\left(1 + \sum_{i=1}^{d} x_{pi} x_{8i}\right) \neq y_8$

∵ $\vec{x}_0 \sim \vec{x}_{p-1}$ is predicted correctly by $\vec{w}_0$

∴ $\text{sign}(\vec{w}_0^T \vec{x}_i) = y_i$ for $i = 1 \sim p-1$

⇒ $\text{sign}(\vec{w}_0^T \vec{x}_i) = y_i$

⇒ $\text{sign}\left( \begin{bmatrix} 0 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix} \right) = y_i$

⇒ $\text{sign}(0) = 1 = y_i$ for $i = 1 \sim p-1$

∵ $\text{sign}(\vec{w}_0^T \vec{x}_p) \neq y_p$ as we assumed ∴ $y_p = 1$

Thus, $\vec{w}_1 = \vec{w}_0 + y_p \vec{x}_p = \vec{0} + 1 \cdot \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix} = \begin{bmatrix} x_{p1} \\ x_{p2} \\ \vdots \\ x_{pd} \end{bmatrix}$

Consider the same situation but replace all $\vec{x}_n$ to $\vec{x}_n'$ :

for $\vec{x}_i'$ $i = 1 \sim p-1$, since $\text{sign}(\vec{w}_0^T \vec{x}_i') = 1$ and $y_i$ remain unchanged

∴ $\vec{x}_1' \sim \vec{x}_{p-1}'$ is predicted correctly

For $\vec{x}_p'$, we have also have $\text{sign}(\vec{w}_0^T \vec{x}_p') = 1 \neq y_p = 1$

Therefore,

$$\vec{w}_1' = \vec{w}_0' + y_p \vec{x}_p' = \vec{0} + 1 \cdot \begin{bmatrix} 3x_{p1} \\ 3x_{p2} \\ \vdots \\ 3x_{pd} \end{bmatrix} = \begin{bmatrix} 3x_{p1} \\ 3x_{p2} \\ \vdots \\ 3x_{pd} \end{bmatrix}$$

∵ $\vec{x}_{p+1} \sim \vec{x}_{q-1}$ is predicted correctly

∴ $\text{sign}(\vec{x}_p^T \vec{x}_i) = y_i$ for $i = p+1 \sim q-1$

We want to know if $\tilde{x}_{p+1}' \sim \tilde{x}_{q}'$ will be predicted correctly or not

as $\quad \tilde{w}_1' = [3 \;\; 3\tilde{x}_{p1} \cdots 3\tilde{x}_{pd}]^T$

$\Rightarrow \; \text{sign} (\tilde{w}_1'^T \tilde{x}_{p+1}') = \text{sign} ( [3 \;\; 3\tilde{x}_{p1} \cdots 3\tilde{x}_{pd}] \begin{bmatrix} 3 \\ 3\tilde{x}_{p+1,1} \\ 3\tilde{x}_{p+1,2} \\ \vdots \\ 3\tilde{x}_{p+1,d} \end{bmatrix} )$

$\qquad = \text{sign} ( 9 (1 + \tilde{x}_{p1}\tilde{x}_{p+1,1} + \cdots + \tilde{x}_{pd}\tilde{x}_{p+1,d} ) )$

$\qquad = \text{sign} ( 1 + \tilde{x}_{p1}\tilde{x}_{p+1,1} + \cdots + \tilde{x}_{pd}\tilde{x}_{p+1,d} )$

$\qquad = \text{sign} ( \tilde{x}_p^T \tilde{x}_{p+1} )$

$\qquad = \text{sign} ( \tilde{w}_1^T \tilde{x}_{p+1} )$

Similarly, we know that:

$\qquad \text{sign} ( \tilde{w}_1'^T \tilde{x}_i' ) = \text{sign} (\tilde{w}_1^T \tilde{x}_i) \quad \text{for} \;\; i = p+1 \sim q-1$

Thus, $\tilde{x}_i'$ is predicted correctly by $\tilde{w}_1'$ for $i = p+1 \sim q-1$

Next we want to check if $\text{sign} (\tilde{w}_1'^T \tilde{x}_q') = y_q$ or not:

$\text{sign} (\tilde{w}_1'^T \tilde{x}_q') = \text{sign} ( [3 \;\; 3\tilde{x}_{p1} \cdots 3\tilde{x}_{pd}] \begin{bmatrix} 3 \\ 3\tilde{x}_{q1} \\ 3\tilde{x}_{q2} \\ \vdots \\ 3\tilde{x}_{qd} \end{bmatrix} )$

$\qquad = \text{sign} ( 9 (1 + \tilde{x}_{p1}\tilde{x}_{q1} + \cdots + \tilde{x}_{pd}\tilde{x}_{qd} ) )$

$\qquad = \text{sign} ( 9 (1 + \sum_{i=1}^{d} \tilde{x}_{pi}\tilde{x}_{qi} ) )$

$\qquad = \text{sign} ( 1 + \sum_{i=1}^{d} \tilde{x}_{pi}\tilde{x}_{qi} )$

$$= \text{sign} \left( \vec{x_p}^T \vec{x_g} \right)$$

$$= \text{sign} \left( \vec{w_i}^T \vec{x_g} \right)$$

Thus, if $\text{sign} \left( \vec{w_i}^T \vec{x_g} \right) \neq y_g$, then $\text{sign} \left( \vec{w_i}'^T \vec{x_g}' \right) \neq y_g$

Therefore, we can know that $\vec{w}_{PLA} = \vec{w}_{PLA}'$.

9. Let article: $\vec{x_t} = \begin{bmatrix} 1 & x_{t1} & x_{t2} & \cdots & x_{td} \end{bmatrix}^T$ $t \in \{1, ..., \textcircled{N}\}$

dictionary: $D$, where $|D| = d \geq m$

$$D = \begin{bmatrix} w_1 & \cdots & w_{d_+} & w_{d_++1} & \cdots & w_{d_++d_-} \end{bmatrix}^T \quad {}^{w_d}$$

more hatred like words          less hatred like words

$\forall x_{ti}$ $i \in \{1,...,d\}$  $x_{ti} = 1$ if $w_i \in$ article $t$

$$\begin{cases} z_+(\vec{x_t}) = \sum_{i=1}^{d_+} x_{ti} \\ \\ z_-(\vec{x_t}) = \sum_{i=d_++1}^{d} x_{ti} \end{cases}$$

The mistake bound of PLA $= \dfrac{R^2}{\rho^2}$, where 

$\begin{cases} R^2 = \max\limits_{n} \| \vec{x_n} \|^2 \\ \\ \rho = \min\limits_{n} y_n \dfrac{\vec{w_f}^T}{\|\vec{w_f}\|} \vec{x_n} \quad \because \|\vec{w_f}\| = \sqrt{\vec{w_f}^T \vec{w_f}} \\ \qquad\qquad\qquad\qquad \therefore \|\vec{w_f}\|^2 = \vec{w_f}^T \vec{w_f} \\ \\ \Rightarrow \rho^2 = \min\limits_{n} \dfrac{(\vec{w_f}^T \vec{x_n})^2}{\vec{w_f}^T \vec{w_f}} \end{cases}$

Since we count number of hatred words to determine if an article is hatred or not, and the bias is -3.5, $\vec{wf}$ is:

$$\vec{wf} = \begin{bmatrix} -3.5 & +1 & +1 & \cdots & +1 & \overbrace{-1 & -1 & \cdots & -1}^{d} \end{bmatrix}^T \in \mathbb{R}^{d+1}$$

$\underbrace{\quad}_{d+}$  $\underbrace{\quad}_{d-}$

$$\left( \vec{x} = \begin{bmatrix} x_0 & x_1 & \cdots & x_d \end{bmatrix}^T = \begin{bmatrix} 1 & x_1 & \cdots & \boxed{x_d} \end{bmatrix}^T \right)$$

↳ 如果 word i 在 article $\vec{x}$ 中 ⇒ $x_i = 1$
     for $i = 1, \cdots, d$

$\therefore \|\vec{wf}\| = \sqrt{\vec{wf}^T \vec{wf}} = \sqrt{12.25 + 1^2 \times d_+ + (-1)^2 \times d_-} = \sqrt{12.25 + d}$

$\Rightarrow \ell^2 = \min_n \dfrac{(\vec{wf}^T \vec{x}_n)^2}{\vec{wf}^T \vec{wf}} = \min_n \dfrac{(\vec{wf}^T \vec{x}_n)^2}{12.25 + d} = \dfrac{0.25}{12.25 + d}$

suppose $d_+ > d_-$
则 $\min$ 发生在 $\vec{x} =$

$$\begin{bmatrix} 1 & 1 & \underbrace{\cdots 1}_{} & 1 & 1 & \underbrace{1 \cdots 0}_{} & 1 & \cdots & 1 \end{bmatrix}^T$$

$\underbrace{\quad}_{\min(d_+, d_-)}$  $\underbrace{\quad}_{\min(d_+, d_-)}$
$\quad \| \quad$  $\quad \| \quad$
$\quad d_- \quad$  $\quad d_- \quad$

$R^2 = \max_n \|\vec{x}_n\|^2$

$= \max_n \vec{x}_n^T \vec{x}_n$

取 $\vec{x}_n = \begin{bmatrix} 1 & \underbrace{1 \ 1 \ \cdots \ 1}_{m \text{ 位}} & 0 & \cdots & 0 \end{bmatrix}^T$
     $\ \| \ $
     $x_0$

$\therefore \min_n (\vec{wf}^T \vec{x}_n)^2 = \min_n (-3.5 + d_- \cdot (+1) + 0 \cdot d_-)^2$
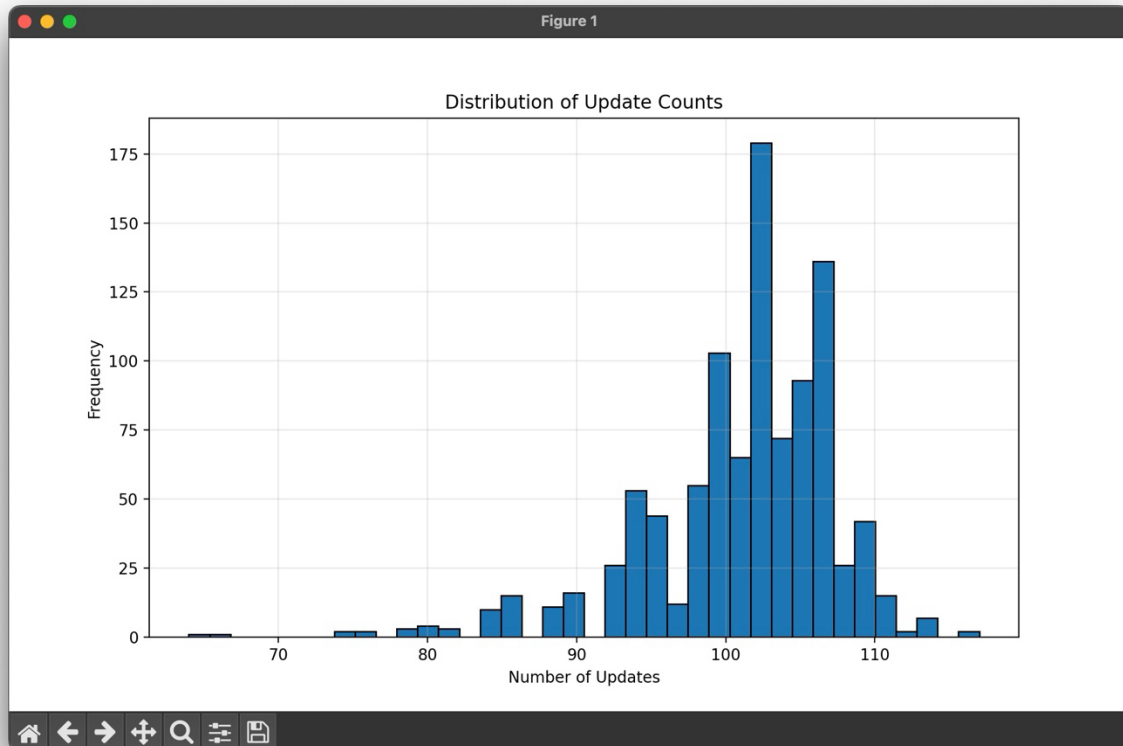$= \min_n (-0.5)^2 = 0.25$

$\therefore R^2 = m+1$

Therefore, the upper bound of mistake number $= \dfrac{m+1}{\dfrac{0.25}{12.25 + d}}$

□

# ML homework 1: Question 10

## Part 1: Interpretation

The resulting histogram generated by my code is as below:



From the histogram, we can see that the peak is about 104, and most of the number of updates are around 95 to 110. As what we learned in class, the upper bound T is proportional to $R^2$, if we consider the high dimensionality of the inputs, we might expect T to be a considerably large value, which does not match the result of the distribution. I think the reason would be the sparseness of the inputs, which makes $R^2$ much smaller than a 47206 dimension input should probably have.

Another interpretation is that, since the input vector with the 1 added as the zero-th term is of dimension 47206, we knew that the weight vector is of same dimension, which means we have

high flexibility to modify our perceptron. In such a high dimensional space, finding a hyperplane to separate the data points seemed to be feasible, and this reflects in our histogram, where most number of updates are relatively low, indicating that we can achieve 99% of accuracy in around one hundred updates.

However, there seemed to be some extreme cases, like the experiments that have update number higher than 110. This is probably due to the random choosing process that we use allows replacement, also, the dataset is small (200) compared to the dimension (47206), therefore, if in some situation we chose the points near the decision boundary many times, it may be difficult for the PLA to get 1000 consecutive correct predictions.

For the cases that the number of updates are smaller than 90, it might because the similar reason as the extreme large cases we stated above. Since the dataset is small, it is more possible that we're "fortunate" enough to choose the points with large distance repeatedly, so that PLA would converge quickly.

## Part 2: Code snippet

The following pictures are some of my code snippets. Since I wrote lots of comments, providing only the first page would show quite little code that provides real usage:

```python
def dot_product_sign(sparse_vector, w):
    total = 0
    for item in sparse_vector:
        index, value = item.split(':')
        index = int(index)
        total += float(value) * w[index]
    if total > 0:
        return 1
    else:
        return -1
```

```
52    for index, line in enumerate(lines):
53        if index >= N:
54            break
55
56        parts = line.split()
57        label = int(parts[0])
58        y_arr.append(1 if label == 1 else -1)
59        sparse_vector = ["0:1"] + parts[1:]
60        sparse_vector_arr.append(sparse_vector)
```
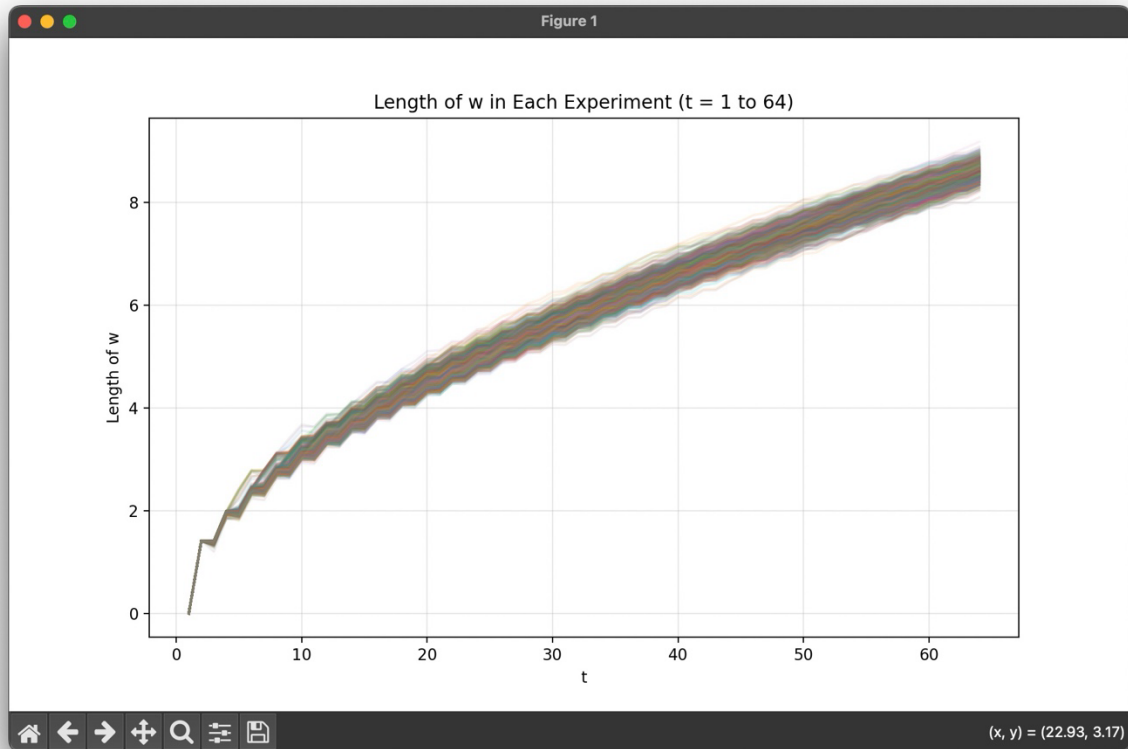
The following part is contained in two loops, where the outer one does the experiment 1000 times, and the inner one loops while the correct amount is smaller than 1000:

```
91          for random_no in random_check_indices:
92
93              sparse_vector = sparse_vector_arr[random_no]
94              label = y_arr[random_no]
95              sign = dot_product_sign(sparse_vector, w)
96              if sign != label:
97                  for item in sparse_vector:
98                      index, value = item.split(':')
99                      index = int(index)
100                     w[index] += label * float(value)
101                 correct_amount = 0
102                 current_update_times += 1
103             else:
104                 correct_amount += 1
105
106
107             if correct_amount == 1000:
108                 break
109
110             # if the current update times is a multiple of 1000, this means that the current "random_check_indices" array has used out,
111             # so we update the array with 1000 new random indices
112             if current_update_times % 1000 == 0:
113                 random_check_indices = [next(random_sequence) for _ in range(5 * N)]
114
115      update_times_arr.append(current_update_times)
```

# ML homework 1: question 11



## Findings:

The minimum number of updates $T_{min}$ (or min_T in my code) I got is 64.

From this graph, we can see that the growth of the length of the weight vector is slow, and the length until the 64-th update did not exceed 10. This might reflect what we learned in class: the square of the growing length in each update is bounded by $R^2$.

However, we can still see that compare to the later updates, the length grows faster in the first 10 updates. This may due to the fact that the initial value of the weight vector is a zero vector, so that it is far from the ultimate $w_{PLA}$, hence making bigger changes in the first few updates. The growth became more stable after some updates, it is probably because after the first few

updates, we have constructed a hyperplane that can classify most of the examples correctly, so the amplitude of the adjustments became small.

Moreover, we can see that the variation of the 1000 experiments is larger in the first few updates, but as the number of update (t) increases, the variations became small, indicating that we get similar weight lengths after a certain number of updates. I think that the reason is because we got 99% accuracy after 1000 consecutive correct predictions, and the dataset we use is the same, so that if the data points are linear separable, then there exists a hyperplane that can correctly classify all the points, even if there is some noise, we still can get a hyperplane that can separate most of the points. Therefore, as the high accuracy we got in each experiment, the weight vector we got will represent similar hyperplanes, hence resulting a similar value of the weight vector length.

## Part of code:

Most of the code in this question is similar to question 10, I just add some modifications to record the history length of $\vec{w}$, take the minimum value from the resulting value of Ts, and change the figure plotting part.

Due to the same problem as question 10, there's too many comments in my code, so I would provide the snapshots of the parts that include the modifications as below:

```python
w_lengths = []                    # use list w_lengths to store the lengths for w in each experiment

# aim: repeat the experiment 1000 times
for experiment_no in tqdm(range(1000)):
    w = np.zeros(47206)
    correct_amount = 0            # the amount of correct predictions, this value is reseted to 0 if there's a wrong prediction
    current_update_times = 0      # the total amount of updates until 1000 consecutive correct predictions
    seed = experiment_no

    w_length_history = [np.linalg.norm(w)]  # initialize w_length_history with initial w length
```

```python
        random_check_indices = [next(random_sequence) for _ in range(5 * N)]
        # subaim: for each index in the random check indices
        for random_no in random_check_indices:

            sparse_vector = sparse_vector_arr[random_no]
            label = y_arr[random_no]
            sign = dot_product_sign(sparse_vector, w)
            if sign != label:
                for item in sparse_vector:
                    index, value = item.split(':')
                    index = int(index)
                    w[index] += label * float(value)
                correct_amount = 0
                current_update_times += 1
                w_length_history.append(np.linalg.norm(w))                    # append the length of w after each update
            else:
                correct_amount += 1


            if correct_amount == 1000:
                break

            # if the current update times is a multiple of 1000, this means that the current "random_check_indices" array has used out,
            # so we update the array with 1000 new random indices
            if current_update_times % 1000 == 0:
                random_check_indices = [next(random_sequence) for _ in range(5 * N)]

    update_times_arr.append(current_update_times)
    w_lengths.append(w_length_history)
```

```python
min_T = min(update_times_arr)

# aim: plot the length of w for each experiment
plt.figure(figsize=(10, 6))
for length in w_lengths:
    plot_length = min(min_T, len(length))
    plt.plot(range(1, plot_length + 1), length[:plot_length], alpha=0.1)
plt.xlabel('t')
plt.ylabel('Length of w')
plt.title(f'Length of w in Each Experiment (t = 1 to {min_T})')
plt.grid(True, alpha=0.3)
plt.show()

print(f"Minimum number of updates: {min_T}")
```