

Homework # 1

姓名	學號
操之晴	R13922A04

5. The answer ChatGPT gave me was “medical image analysis”, such as for cancer detection. It’s convincing because “active learning” is especially useful when labeled data is scarce or expensive to obtain. Imagine a situation where we need a few expert radiologists to label every medical image, it can be both time-consuming and costly.

In this case, active learning can help by allowing the model to identify the most uncertain or difficult cases for labeling, reducing the number of images that need to be labeled while still maintaining high performance. This makes the process more efficient and improves accuracy while minimizing the effort required from experts.

For example, the model can select images that are most likely to be misclassified or those that provide the most value in distinguishing between cancerous and healthy tissues.

That’s why I agree with this answer. Medical image analysis is a possible application of active learning.

6. The answer ChatGPT gave me was “Yes”, and it listed several ways machine learning can be applied to earthquake prediction: Seismic Data Analysis, Early Warning Systems, Identifying Risk Zones and Aftershock Prediction.

I partially agree with this answer. Since the problem was about “earthquakes prediction”, the focus should be on predicting earthquakes that haven’t happened yet. In that sense, only the first point “Seismic Data Analysis” is directly relevant.

Machine learning models can analyze large amounts of historical seismic data to detect patterns that might signal an upcoming earthquake. This could help find precursor signals, such as specific types of seismic waves, that suggest a higher risk of earthquake. However, predicting the exact time and location of an earthquake is still a huge challenge. That’s why more attention is given to the other tasks listed above.

7. (convention 0: $\text{sign}(0) = 0$)

Suppose followings are two data points and its label:

data points x	$x_1 = (-1, 2)$	$x_2 = (1, -1)$
Label y	+1	-1
$x_0 = 1$	$(1, -1, 2)$	$(1, 1, -1)$
	$w_1 = w_0 + y_1 x_1 = (1, -1, 2)$	$w_2 = w_1 = (1, -1, 2)$
$x'_0 = 2$	$(2, -1, 2)$	$(2, 1, -1)$
	$w_1 = w_0 + y_1 x_1 = (2, -1, 2)$	$w_2 = w_1 + y_2 x_2 = (0, -2, 3)$

From above, $w_{PLA} = (1, -1, -1)$ and $w'_{PLA} = (0, -2, -2)$. When we take $x_3 = (2, 1)$, in $x_0 = 1$, $\text{sign}((1, -1, 2) \cdot (2, 1)) = 1$. In $x'_0 = 2$, $\text{sign}((0, -2, 3) \cdot (2, 1)) = -1$. They return in different binary classification output. Therefore, they are not equivalent.

8. When we take $x'_n = 3x_n$ $w_0 = 0$ convention 1: $\text{sign}(0) = 1$

$$w_1 = w_0 + y_1 x_1 \quad w'_1 = w_0 + y_1 x'_1 = y_1 3x_1 = 3w_1$$

For next x_2 , we have two possible cases:

Case 1: $\text{sign}(w_1 \cdot x_2) \neq y_2$

$$w_2 = w_1 + y_2 x_2$$

Since $\text{sign}(w'_1 \cdot x_2) = \text{sign}(3w_1 \cdot x_2) = \text{sign}(w_1 \cdot x_2) \neq y_2$

$$w'_2 = w'_1 + y_2 x'_2 = 3w_1 + y_2 3x_2 = 3w_2$$

Case 2: $\text{sign}(w_1 \cdot x_2) = y_2$

$$w_2 = w_1$$

Since $\text{sign}(w'_1 \cdot x_2) = \text{sign}(3w_1 \cdot x_2) = \text{sign}(w_1 \cdot x_2) = y_2$

$$w'_2 = w'_1 = 3w_1 = 3w_2$$

We can see that no matter the next sample is misclassified or not, $w' = 3w$. Which leads to $w'_{PLA} = 3w_{PLA}$. Then for every possible sample x_k in \mathbb{R}^d ,

$$\text{sign}(w'_{PLA} \cdot x'_k) = \text{sign}(3w_{PLA} \cdot 3x_k) = \text{sign}(9w_{PLA} \cdot x_k) = \text{sign}(w_{PLA} \cdot x_k)$$

Therefore, w_{PLA} and w'_{PLA} are equivalent.

$$9. \quad T \leq \frac{R^2}{\rho^2} \quad z_+ : \text{hatred} - \text{like} \quad z_- : \text{less hatred} - \text{like}$$

$$f(x) = \text{sign}(z_+(x) - z_-(x) - 3.5)$$

$$R^2 = \max \|x_n\|^2 = \sqrt{(m+1)}^2 = m+1$$

$$\rho = \min y_n \frac{w_f^T}{\|w_f\|} x_n$$

$$\rho^2 = \frac{1}{\sqrt{d+3.5^2}^2} \min (y_n w_f^T x_n)^2 = \frac{1}{d+12.25} \min (z_+(x) - z_-(x) - 3.5)^2$$

$$\text{When } z_+ = z_-, \text{ we get minimum } \rho^2 = \frac{1}{d+12.25} 3.5^2 = \frac{12.25}{d+12.25} = \frac{49}{4d+49}$$

$$\text{Thus, the upper bound is } T \leq \frac{R^2}{\rho^2} = \frac{1}{49} (4d+49)(m+1)$$

10.

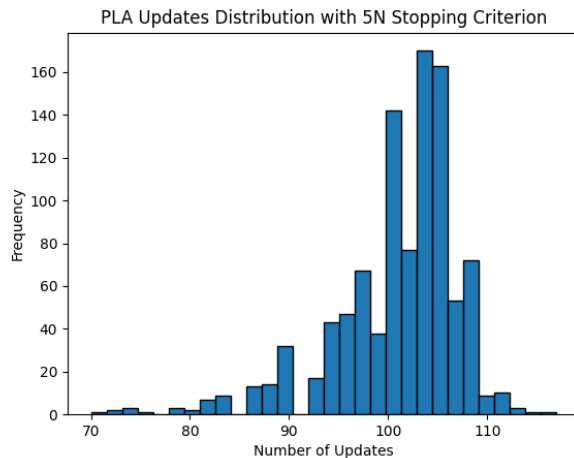


Figure 10a. updates histogram

```

q10.py > ...
1 import numpy as np
2 import random
3 from sklearn.datasets import load_svmlight_file
4 import matplotlib.pyplot as plt
5
6 # parameters
7 experiment = 1000
8 sample = 200
9 #max_update = 5000
10 update_counts = [] # save every experiment update counts
11
12 # loading datas
13 def load_data(path, sample):
14     x, y = load_svmlight_file(path) # loading LIBSVM form
15     x = x[sample:].toarray() # to numpy array
16     y = y[sample:]
17     return x, y
18
19 x, y = load_data('rcv1_train.binary', sample)
20
21 # adding bias x0 = 1 to every xn
22 x = np.hstack((np.ones((x.shape[0], 1)), x))
23
24 # PLA
25 def pla(x, y):
26     w = np.zeros(x[0].shape) # initialize w = 0
27     updates = 0
28     sample = x.shape[0]
29     max_SN = 5 * sample
30     no_error = 0
31
32     while True:
33         i = np.random.choice(sample) # with replacement
34         if np.sign(np.dot(w, x[i])) != y[i]: # sign different
35             w = w + y[i] * x[i]
36             updates += 1
37             no_error = 0
38         else:
39             no_error += 1
40         if no_error >= max_SN:

```

Figure 10b. snapshot of my first page of code

a. Distribution Range

From figure 10a, we can see that the histogram has a positive skew, with most experiments needing around 100 to 110 updates. This suggests that it typically takes about that many updates for the PLA to converge and correctly classify all the samples. With this tight range, I found PLA is consistent in its performance on this dataset, usually requiring a similar number of updates.

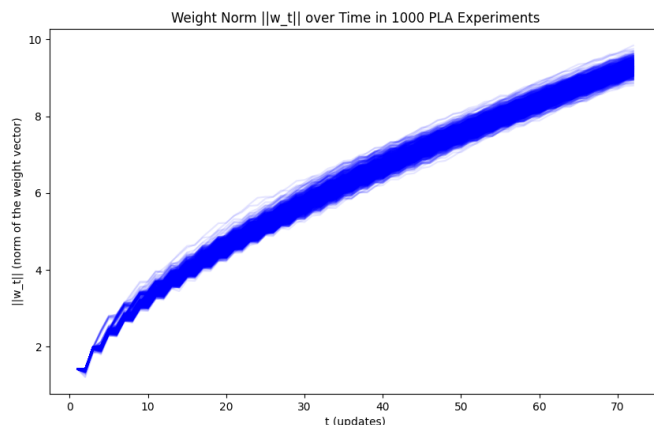
b. Some Outliers

A few experiments needed more than 110 updates, likely because the algorithm ran into more difficult-to-classify samples, which required extra updates to resolve.

In summary, most experiments converge in around 100 updates, showing that this is the typical behavior of the PLA on this dataset. Although a few cases needed significantly more updates due to the randomness in sample selection, those were rare.

*For problems 10 to 12, the full code has been made available on [GitHub](#).
Feel free to check it out if required or interested!*

11.

Figure 11a. $\|w_t\|$ vs. t

```

q11.py > ...
1  import numpy as np
2  import random
3  from sklearn.datasets import load_svmlight_file
4  import matplotlib.pyplot as plt
5
6  # parameters
7  experiment = 1000
8  sample = 200
9  #max_update = 5000
10 update_counts = [] # save every experiment update counts
11 all_w_norm = [] # save every experiment w norms
12 min_update = np.inf # save minimum updates
13
14
15 # loading datas
16 def load_data(path, sample):
17     x, y = load_svmlight_file(path) # loading LIBSVM form
18     x = x[sample:].toarray() # to numpy array
19     y = y[sample:]
20     return x, y
21
22 x, y = load_data('rcv1_train.binary', sample)
23
24 # adding bias x0 = 1 to every xn
25 x = np.hstack([np.ones((x.shape[0], 1)), x])
26
27 # PLA
28 def pla(x, y):
29     w = np.zeros(x[0].shape) # initialize w = 0
30     updates = 0
31     sample = x.shape[0]
32     max_5N = 5 * sample
33     no_error = 0
34     w_norm = []
35
36     while True:
37         i = np.random.choice(sample) # with replacement
38         if np.sign(np.dot(w, x[i])) != y[i]: # sign different
39             w = w + y[i] * x[i]
40             updates += 1

```

Figure 11b. snapshot of my first page of code

a. Rapid Growth in Early Stages

At the beginning, $\|w_t\|$ grows quickly. I assume this is because in the early stages of PLA, each update brings large adjustments since the algorithm is still far from finding the correct classification boundary.

b. Slowing Growth as Convergence Nears

As the number of updates increases, the growth of $\|w_t\|$ slows down. This shows PLA is getting closer to the correct classification boundary, where smaller adjustments are needed.

Once the algorithm correctly classifies all samples, the weights stop changing, so $\|w_t\|$ stops growing.

c. Similarity in All Experiments

All curves show a similar pattern of above: rapid growth followed by slower growth and eventually stabilization.

In summary, $\|w_t\|$ grows quickly at the start of the learning process and then stabilizes as the PLA converges. Also $\|w_t\|$ shows a similar pattern across different experiments.

12.

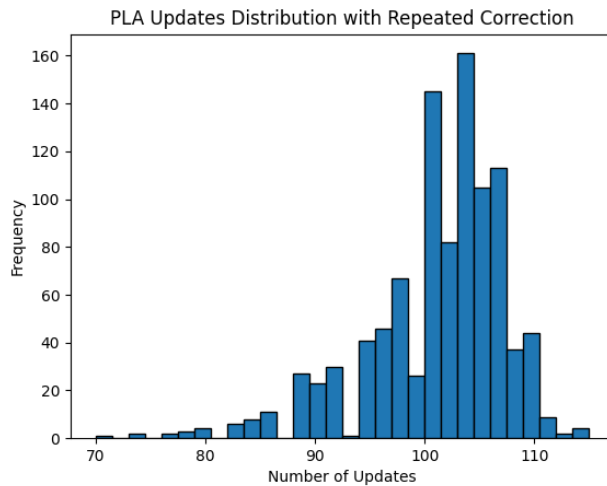


Figure 12a. updates histogram

```

q12.py > ...
1 import numpy as np
2 import random
3 from sklearn.datasets import load_svmlight_file
4 import matplotlib.pyplot as plt
5
6 # parameters
7 experiment = 1000
8 sample = 200
9 #max_update = 5000
10 update_counts = [] # save every experiment update counts
11
12 # loading datas
13 def load_data(path, sample):
14     x, y = load_svmlight_file(path) # loading LIBSVM form
15     x = x[sample].toarray() # to numpy array
16     y = y[sample]
17     return x, y
18
19 x, y = load_data('rcv1_train.binary', sample)
20
21 # adding bias x0 = 1 to every xn
22 x = np.hstack((np.ones((x.shape[0], 1)), x))
23
24 # PLA
25 def pla(x, y):
26     w = np.zeros(x[0].shape) # initialize w = 0
27     updates = 0
28     sample = x.shape[0]
29     max_SN = 5 * sample
30     no_error = 0
31
32     while no_error < max_SN:
33         i = np.random.choice(sample) # with replacement
34         while np.sign(np.dot(w, x[i])) != y[i]: # sign different
35             w = w + y[i] * x[i]
36             updates += 1
37         no_error += 1
38     return w, updates
39
40 for i in range(experiment):

```

Figure 12b. snapshot of my first page of code

Median number of updates: 102.0

Figure 12c. median number of updates

a. Results Comparison between 10 and 12

In problem 10, we get our median number of updates 102, which is identical to problem 12. Initially, I hypothesized that problem 10 might result in fewer updates overall because it doesn't stick with a wrong sample to fully correct it before moving on. This could lead to a more compact distribution. On the other hand, prob.12 might have a more spread-out update distribution. However, based on the median and the figures, problem 12 only updates once for every sample, which end up with similar update counts and distribution patterns as problem 10.

b. Advantages

By fully correcting each misclassified sample before moving on, the algorithm reduces the chances of revisiting the same problem repeatedly in later iterations. This increase classification accuracy per sample.

It also allows the algorithm to focus on a specific issue and resolve it entirely, which might result in fewer overall updates needed to achieve convergence on difficult samples.

c. Drawbacks

While the algorithm might fully correct each sample, this approach could slow down the overall learning process, as it may spend more time on difficult samples. The trade-off here is between focusing on immediate corrections versus spreading updates more evenly across different samples.

13. Proof 1 :

For a misclassified sample x_n , we have $y_{n(t)} (w_t^T x_{n(t)}) \leq 0$

From the new update rule, we get

$$w_{t+1} = w_t + \frac{1}{10} y_{n(t)} x_{n(t)} \cdot \left[\frac{-10 y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 \right] \quad (1)$$

We need to check whether $y_{n(t)} (w_{t+1}^T x_{n(t)}) > 0$ or not. So, we replace it with (1).

$$\begin{aligned} y_{n(t)} (w_{t+1}^T x_{n(t)}) &= y_{n(t)} \left(\left(w_t + \frac{1}{10} y_{n(t)} x_{n(t)} \cdot \left[\frac{-10 y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 \right] \right)^T x_{n(t)} \right) \\ &= y_{n(t)} w_t^T x_{n(t)} + \left(\frac{1}{10} y_{n(t)}^2 x_{n(t)} \cdot \left[\frac{-10 y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 \right] \right)^T x_{n(t)} \end{aligned}$$

Since $y_{n(t)}^2 = 1$, we get

$$\begin{aligned} &y_{n(t)} w_t^T x_{n(t)} + \left[x_{n(t)} \cdot \left[\frac{-y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + \frac{1}{10} \right] \right]^T x_{n(t)} \\ &= y_{n(t)} w_t^T x_{n(t)} + \|x_{n(t)}\|^2 \cdot \left[\frac{-y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + \frac{1}{10} \right] \\ &= y_{n(t)} w_t^T x_{n(t)} - y_{n(t)} w_t^T x_{n(t)} + \frac{\|x_{n(t)}\|^2}{10} = \frac{\|x_{n(t)}\|^2}{10} \end{aligned}$$

Since $\|x_{n(t)}\|^2$ is always positive, $y_{n(t)} (w_{t+1}^T x_{n(t)})$ ends in positive. The updated weight vector will always correctly classify after the update.

Proof 2 :