

## ML Assignment #1

409410120 資工三 劉哲嘉

### Execution description:

檔案說明： $f(x) = y = 3x + 7$  隨機座標範圍： $-100 < x, y < 100$

- (1) points30.py 創造 15 筆正例(1) 15 筆反例(-1) 以  $f(x)$  線性分割
- (2) points2000.py 創造 1000 筆正例(1) 1000 筆反例(-1) 以  $f(x)$  線性分割
- (3) miss.py 創造 1000 筆正例(1) 1000 筆反例(-1) 兩者都有各 50 筆標示錯誤共 100 筆錯誤
- (4) pla.py PLA 演算法，當資料線性可分割會自動停止，不可分割最多迭代 10 萬次
- (5) pocket.py Pocket 演算法，當資料線性可分割會自動停止，不可分割最多迭代 10 萬次

使用 makefile：(in terminal)至檔案目錄下

執行：make data 獲得資料集 points30.txt points2000.txt miss.txt

執行：make pla30 PLA 演算法 資料集：points30.txt

執行：make pla2000 PLA 演算法 資料集：points2000.txt

執行：make plaMiss PLA 演算法 資料集：miss.txt

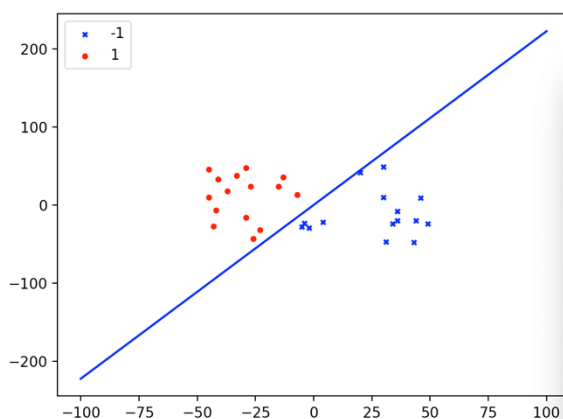
執行：make pok30 Pocket 演算法 資料集：points30.txt

執行：make pok2000 Pocket 演算法 資料集：points2000.txt

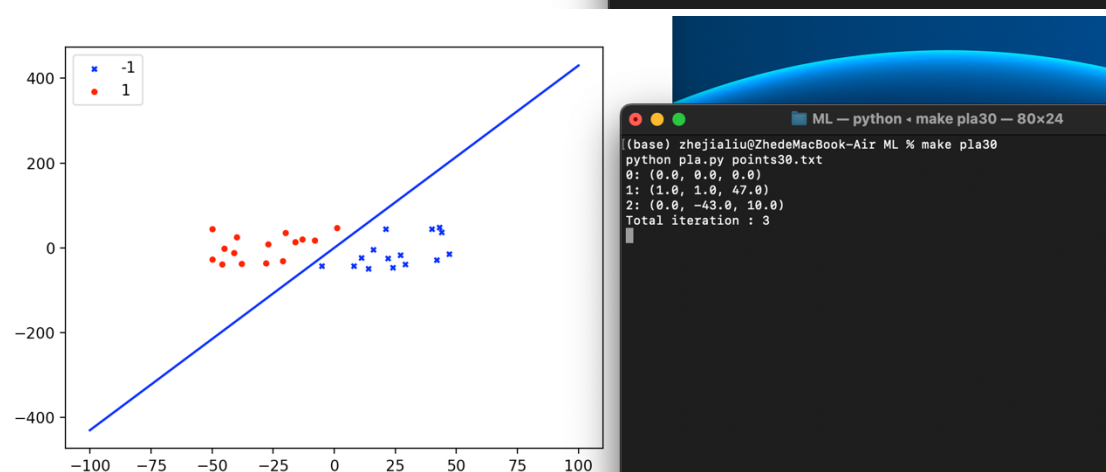
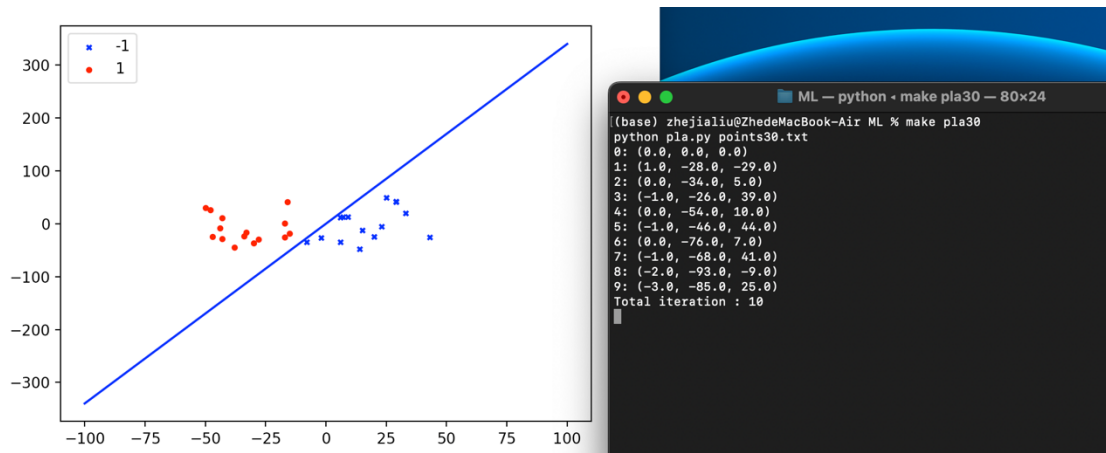
執行：make pokMiss Pocket 演算法 資料集：miss.txt

### Experimental results:

Problem 2 : Generate the data samples three times and calculate the average number of iterations when PLA halts.



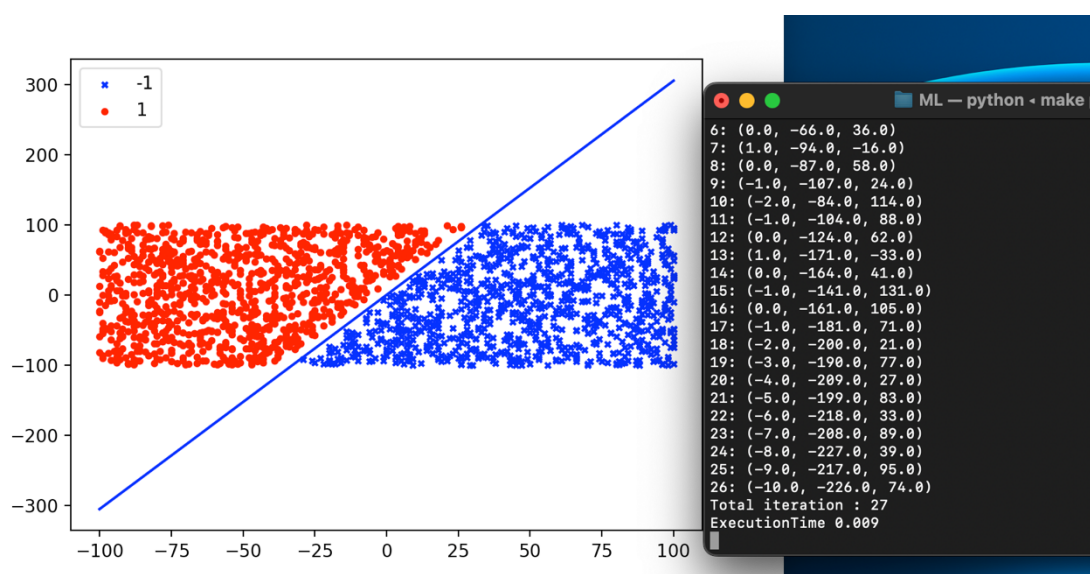
```
ML - python - make pla30 - 80x24
(base) zhejialiu@ZhedsMacBook-Air ML % make pla30
python pla.py points30.txt
0: (0.0, 0.0, 0.0)
1: (1.0, -37.0, 18.0)
2: (0.0, -57.0, -24.0)
3: (-1.0, -53.0, -1.0)
4: (-2.0, -49.0, 22.0)
Total iteration : 5
```



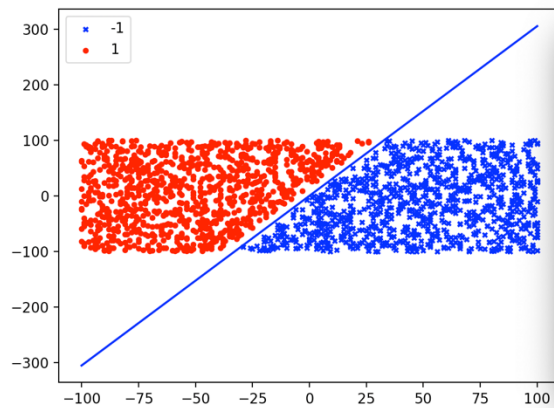
PLA Average iterations :  $(5+10+3)/3 = 6$

Problem 3 : Implement Pocket Algorithm and compare the execution time to PLA on the same dataset ( 2000 data )

PLA : 0.009 s



Pocket : 0.093 s



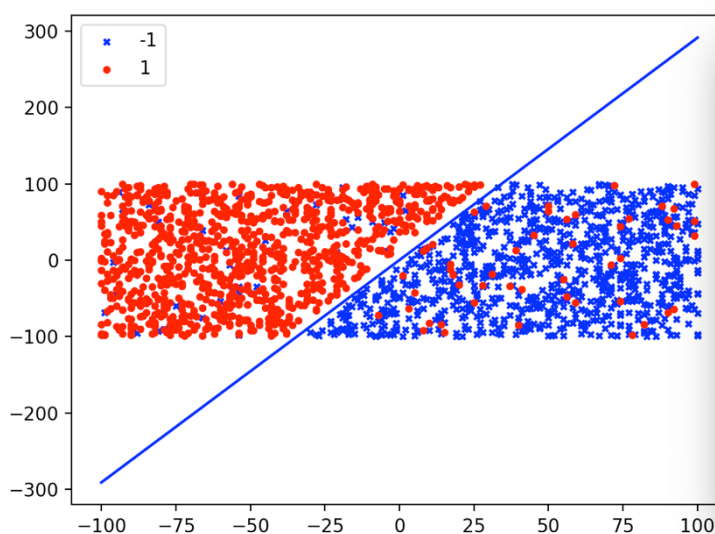
```
ML - python - make pok2000 - 80x24
3: (4.0, -161.0, -11.0) minW (3.0, -113.0, 78.0) errors: 173
Update Min!!!
4: (3.0, -165.0, 77.0) minW (3.0, -165.0, 77.0) errors: 30
5: (2.0, -201.0, -10.0) minW (3.0, -165.0, 77.0) errors: 30
Update Min!!!
6: (3.0, -198.0, 60.0) minW (3.0, -198.0, 60.0) errors: 1
7: (2.0, -169.0, 152.0) minW (3.0, -198.0, 60.0) errors: 1
8: (1.0, -178.0, 131.0) minW (3.0, -198.0, 60.0) errors: 1
9: (2.0, -225.0, 32.0) minW (3.0, -198.0, 60.0) errors: 1
10: (1.0, -213.0, 86.0) minW (3.0, -198.0, 60.0) errors: 1
11: (0.0, -242.0, 8.0) minW (3.0, -198.0, 60.0) errors: 1
12: (-1.0, -230.0, 95.0) minW (3.0, -198.0, 60.0) errors: 1
13: (-2.0, -244.0, -3.0) minW (3.0, -198.0, 60.0) errors: 1
14: (-3.0, -245.0, 95.0) minW (3.0, -198.0, 60.0) errors: 1
15: (-4.0, -279.0, 5.0) minW (3.0, -198.0, 60.0) errors: 1
16: (-3.0, -274.0, 65.0) minW (3.0, -198.0, 60.0) errors: 1
17: (-2.0, -248.0, 161.0) minW (3.0, -198.0, 60.0) errors: 1
18: (-1.0, -263.0, 137.0) minW (3.0, -198.0, 60.0) errors: 1
19: (0.0, -296.0, 70.0) minW (3.0, -198.0, 60.0) errors: 1
Update Min!!!
20: (-1.0, -290.0, 95.0) minW (-1.0, -290.0, 95.0) errors: 0
Accuracy: 100.0%
ExecutionTime 0.093
```

PLA : 0.009 s < Pocket : 0.093 s

在可線性分割的資料集下 PLA 的速度會比 Pocket 快

Problem 4 : Mislabel 50 positive and 50 negative samples by incorrect label.

Report the accuracy of Pocket Algorithm by this setting and the setting in Problem 3.



```
ML - p
99979: (-78.0, -24.0, 46.0)
99980: (-77.0, -97.0, -51.0)
99981: (-76.0, -133.0, 16.0)
99982: (-75.0, -118.0, 99.0)
99983: (-74.0, -26.0, 35.0)
99984: (-73.0, -116.0, -50.0)
99985: (-72.0, -117.0, -1.0)
99986: (-73.0, -109.0, 49.0)
99987: (-72.0, -147.0, -47.0)
99988: (-71.0, -122.0, 52.0)
99989: (-70.0, -34.0, 123.0)
99990: (-69.0, -94.0, 78.0)
99991: (-68.0, -123.0, 26.0)
99992: (-67.0, -100.0, 124.0)
99993: (-66.0, -22.0, 26.0)
99994: (-65.0, -87.0, -64.0)
99995: (-64.0, -110.0, -2.0)
99996: (-65.0, -87.0, 78.0)
99997: (-66.0, -124.0, 22.0)
99998: (-65.0, -104.0, 101.0)
99999: (-66.0, -86.0, 47.0)
Accuracy: 95.0%
ExecutionTime 406.696
```

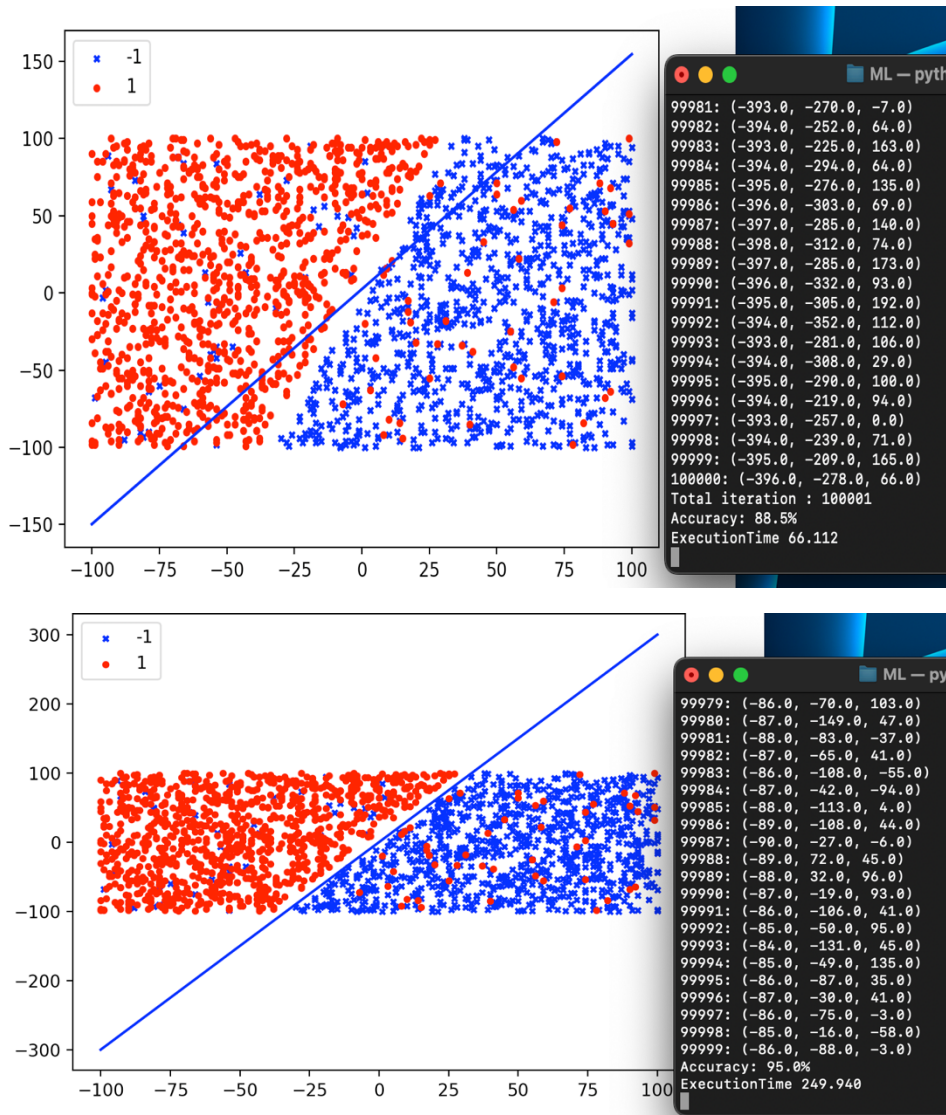
Accuracy = ( ( samples - errors ) / samples ) \* 100%

在 problem3 中資料集是線性可分割的，因此 Accuracy 可以達到 100%，而在 mislabel 的資料集中有 100 筆的錯誤標記，Accuracy 僅有 95%

$$(2000 - 100) / 2000 = 0.95$$

## Conclusion:

PLA 在線性可分割的資料集可以找到  $f(x)$  將資料二分，Pocket 也同樣可以找到  $f(x)$  但因為 Pocket 需要將一最佳權重記錄下來，每一次都需要計算錯誤資料量相比較，所以在線性可分割的資料集中速度會比 PLA 慢，而當資料集不可分割的狀態，經由實驗觀察在同樣迭代 10 萬次的情況，Pocket 的結果會比 PLA 的結果更好



在 mislabel 的資料集

PLA 的 Accuracy = 88.5%

Pocket 的 Accuracy = 95%

若已知資料集為線性可分割的情況 PLA 速度快且可以二分，但若不可分割的情況雖然 PLA 的速度較快，但準確率卻低於 Pocket，因此需要針對資料集的狀況選擇不同的演算法。

## Discussion:

Pocket 演算法會將最佳的權重記錄下來，而當找到一個較佳權重時一開始的做法是一直針對這項權重 update，會發現錯誤數量一直降不下來（在線性可分割的狀況），所以改變作法，變成只先記錄下來，讓權重繼續以隨機的方式更新，最後再回傳最佳的權重值，避免單一權重持續使用的狀況，使得線性可分割的資料可以找到  $f(x)$ 。